



CUSPARSE LIBRARY

DU-06709-001_v5.5 | May 2013



TABLE OF CONTENTS

Chapter 1. Introduction.....	1
1.1. Naming Conventions.....	1
1.2. Asynchronous Execution.....	2
Chapter 2. Using the CUSPARSE API.....	3
2.1. Thread Safety.....	3
2.2. Scalar Parameters.....	3
2.3. Parallelism with Streams.....	3
Chapter 3. CUSPARSE Indexing and Data Formats.....	5
3.1. Index Base Format.....	5
3.2. Vector Formats.....	5
3.2.1. Dense Format.....	5
3.2.2. Sparse Format.....	5
3.3. Matrix Formats.....	6
3.3.1. Dense Format.....	6
3.3.2. Coordinate Format (COO).....	7
3.3.3. Compressed Sparse Row Format (CSR).....	7
3.3.4. Compressed Sparse Column Format (CSC).....	8
3.3.5. Ellpack-Itpack Format (ELL).....	9
3.3.6. Hybrid Format (HYB).....	10
3.3.7. Block Compressed Sparse Row Format (BSR).....	10
3.3.8. Extended BSR Format (BSRX).....	12
Chapter 4. CUSPARSE Types Reference.....	14
4.1. Data types.....	14
4.2. cusparseAction_t.....	14
4.3. cusparseDirection_t.....	14
4.4. cusparseHandle_t.....	15
4.5. cusparseHybMat_t.....	15
4.5.1. cusparseHybPartition_t.....	15
4.6. cusparseMatDescr_t.....	15
4.6.1. cusparseDiagType_t.....	16
4.6.2. cusparseFillMode_t.....	16
4.6.3. cusparseIndexBase_t.....	16
4.6.4. cusparseMatrixType_t.....	16
4.7. cusparseOperation_t.....	17
4.8. cusparsePointerMode_t.....	17
4.9. cusparseSolveAnalysisInfo_t.....	17
4.10. cusparseStatus_t.....	17
Chapter 5. CUSPARSE Helper Function Reference.....	19
5.1. cusparseCreate().....	19
5.2. cusparseCreateHybMat().....	19

5.3. <code>cusparseCreateMatDescr()</code>	20
5.4. <code>cusparseCreateSolveAnalysisInfo()</code>	20
5.5. <code>cusparseDestroy()</code>	20
5.6. <code>cusparseDestroyHybMat()</code>	21
5.7. <code>cusparseDestroyMatDescr()</code>	21
5.8. <code>cusparseDestroySolveAnalysisInfo()</code>	21
5.9. <code>cusparseGetLevelInfo()</code>	21
5.10. <code>cusparseGetMatDiagType()</code>	22
5.11. <code>cusparseGetMatFillMode()</code>	22
5.12. <code>cusparseGetMatIndexBase()</code>	23
5.13. <code>cusparseGetMatType()</code>	23
5.14. <code>cusparseGetPointerMode()</code>	23
5.15. <code>cusparseGetVersion()</code>	24
5.16. <code>cusparseSetMatDiagType()</code>	24
5.17. <code>cusparseSetMatFillMode()</code>	24
5.18. <code>cusparseSetMatIndexBase()</code>	25
5.19. <code>cusparseSetMatType()</code>	25
5.20. <code>cusparseSetPointerMode()</code>	26
5.21. <code>cusparseSetStream()</code>	26
Chapter 6. CUSPARSE Level 1 Function Reference	27
6.1. <code>cusparse<t>axpyi</code>	27
6.2. <code>cusparse<t>doti</code>	28
6.3. <code>cusparse<t>dotci</code>	29
6.4. <code>cusparse<t>gthr</code>	31
6.5. <code>cusparse<t>gthrz</code>	32
6.6. <code>cusparse<t>roti</code>	33
6.7. <code>cusparse<t>sctr</code>	34
Chapter 7. CUSPARSE Level 2 Function Reference	36
7.1. <code>cusparse<t>bsrmv</code>	36
7.2. <code>cusparse<t>bsrxmv</code>	39
7.3. <code>cusparse<t>csrmmv</code>	42
7.4. <code>cusparse<t>csrsv_analysis</code>	44
7.5. <code>cusparse<t>csrsv_solve</code>	46
7.6. <code>cusparse<t>hybmvm</code>	48
7.7. <code>cusparse<t>hybsv_analysis</code>	49
7.8. <code>cusparse<t>hybsv_solve</code>	51
Chapter 8. CUSPARSE Level 3 Function Reference	53
8.1. <code>cusparse<t>csrmm</code>	53
8.2. <code>cusparse<t>csrmm2</code>	56
8.3. <code>cusparse<t>csrsm_analysis</code>	60
8.4. <code>cusparse<t>csrsm_solve</code>	61
Chapter 9. CUSPARSE Extra Function Reference	64
9.1. <code>cusparse<t>csrgeam</code>	64

9.2. <code>cusparse<t>csrsgemm</code>	67
Chapter 10. CUSPARSE Preconditioners Reference	72
10.1. <code>cusparse<t>csric0</code>	72
10.2. <code>cusparse<t>csrilu0</code>	74
10.3. <code>cusparse<t>gtsv</code>	76
10.4. <code>cusparse<t>gtsv_nopivot</code>	77
10.5. <code>cusparse<t>gtsvStridedBatch</code>	79
Chapter 11. CUSPARSE Format Conversion Reference	81
11.1. <code>cusparse<t>bsr2csr</code>	81
11.2. <code>cusparse<t>coo2csr</code>	83
11.3. <code>cusparse<t>csc2dense</code>	84
11.4. <code>cusparse<t>csr2bsr</code>	85
11.5. <code>cusparse<t>csr2coo</code>	88
11.6. <code>cusparse<t>csr2csc</code>	89
11.7. <code>cusparse<t>csr2dense</code>	91
11.8. <code>cusparse<t>csr2hyb</code>	92
11.9. <code>cusparse<t>dense2csc</code>	94
11.10. <code>cusparse<t>dense2csr</code>	95
11.11. <code>cusparse<t>dense2hyb</code>	97
11.12. <code>cusparse<t>hyb2csr</code>	98
11.13. <code>cusparse<t>hyb2dense</code>	99
11.14. <code>cusparse<t>nnz</code>	100
Chapter 12. Appendix A: CUSPARSE Library C++ Example	103
Chapter 13. Appendix B: CUSPARSE Fortran Bindings	109
13.1. Example B, Fortran Application.....	110
Chapter 14. Appendix C: Acknowledgements	118
Chapter 15. Bibliography	119

Chapter 1.

INTRODUCTION

The CUSPARSE library contains a set of basic linear algebra subroutines used for handling sparse matrices. It is implemented on top of the NVIDIA® CUDA™ runtime (which is part of the CUDA Toolkit) and is designed to be called from C and C++. The library routines can be classified into four categories:

- ▶ Level 1: operations between a vector in sparse format and a vector in dense format
- ▶ Level 2: operations between a matrix in sparse format and a vector in dense format
- ▶ Level 3: operations between a matrix in sparse format and a set of vectors in dense format (which can also usually be viewed as a dense tall matrix)
- ▶ Conversion: operations that allow conversion between different matrix formats

The CUSPARSE library allows developers to access the computational resources of the NVIDIA graphics processing unit (GPU), although it does not auto-parallelize across multiple GPUs. The CUSPARSE API assumes that input and output data reside in GPU (device) memory, unless it is explicitly indicated otherwise by the string **DevHostPtr** in a function parameter's name (for example, the parameter ***resultDevHostPtr** in the function **cusparses<t>doti()**).

It is the responsibility of the developer to allocate memory and to copy data between GPU memory and CPU memory using standard CUDA runtime API routines, such as **cudaMalloc()**, **cudaFree()**, **cudaMemcpy()**, and **cudaMemcpyAsync()**.



The CUSPARSE library requires hardware with compute capability (CC) of at least 1.1 or higher. Please see the *NVIDIA CUDA C Programming Guide*, *Appendix A* for a list of the compute capabilities corresponding to all NVIDIA GPUs.

1.1. Naming Conventions

The CUSPARSE library functions are available for data types **float**, **double**, **cuComplex**, and **cuDoubleComplex**. The sparse Level 1, Level 2, and Level 3 functions follow this naming convention:

cusparses<t>[<matrix data format>]<operation>[<output matrix data format>]

where `<t>` can be **S**, **D**, **C**, **Z**, or **X**, corresponding to the data types **float**, **double**, **cuComplex**, **cuDoubleComplex**, and the generic type, respectively.

The `<matrix data format>` can be **dense**, **coo**, **csr**, **csc**, or **hyb**, corresponding to the dense, coordinate, compressed sparse row, compressed sparse column, and hybrid storage formats, respectively.

Finally, the `<operation>` can be **axpyi**, **doti**, **dotci**, **gthr**, **gthrz**, **roti**, or **sctr**, corresponding to the Level 1 functions; it also can be **mv** or **sv**, corresponding to the Level 2 functions, as well as **mm** or **sm**, corresponding to the Level 3 functions.

All of the functions have the return type **cusparseStatus_t** and are explained in more detail in the chapters that follow.

1.2. Asynchronous Execution

The CUSPARSE library functions are executed asynchronously with respect to the host and may return control to the application on the host before the result is ready. Developers can use the **cudaDeviceSynchronize()** function to ensure that the execution of a particular CUSPARSE library routine has completed.

A developer can also use the **cudaMemcpy()** routine to copy data from the device to the host and vice versa, using the **cudaMemcpyDeviceToHost** and **cudaMemcpyHostToDevice** parameters, respectively. In this case there is no need to add a call to **cudaDeviceSynchronize()** because the call to **cudaMemcpy()** with the above parameters is blocking and completes only when the results are ready on the host.

Chapter 2.

USING THE CUSPARSE API

This chapter describes how to use the CUSPARSE library API. It is not a reference for the CUSPARSE API data types and functions; that is provided in subsequent chapters.

2.1. Thread Safety

The library is thread safe and its functions can be called from multiple host threads.

2.2. Scalar Parameters

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

The few functions that return a scalar result, such as `doti()` and `nnz()`, return the resulting value by reference on the host or the device. Even though these functions return immediately, similarly to those that return matrix and vector results, the scalar result is not ready until execution of the routine on the GPU completes. This requires proper synchronization be used when reading the result from the host.

This feature allows the CUSPARSE library functions to execute completely asynchronously using streams, even when α and β are generated by a previous kernel. This situation arises, for example, when the library is used to implement iterative methods for the solution of linear systems and eigenvalue problems [3].

2.3. Parallelism with Streams

If the application performs several small independent computations, or if it makes data transfers in parallel with the computation, CUDA streams can be used to overlap these tasks.

The application can conceptually associate a stream with each task. To achieve the overlap of computation between the tasks, the developer should create CUDA streams using the function `cudaStreamCreate()` and set the stream to be used by each

individual CUSPARSE library routine by calling `cusparseSetStream()` just before calling the actual CUSPARSE routine. Then, computations performed in separate streams would be overlapped automatically on the GPU, when possible. This approach is especially useful when the computation performed by a single task is relatively small and is not enough to fill the GPU with work, or when there is a data transfer that can be performed in parallel with the computation.

When streams are used, we recommend using the new CUSPARSE API with scalar parameters and results passed by reference in the device memory to achieve maximum computational overlap.

Although a developer can create many streams, in practice it is not possible to have more than 16 concurrent kernels executing at the same time.

Chapter 3.

CUSPARSE INDEXING AND DATA FORMATS

The CUSPARSE library supports dense and sparse vector, and dense and sparse matrix formats.

3.1. Index Base Format

The library supports zero- and one-based indexing. The index base is selected through the `cusparseIndexBase_t` type, which is passed as a standalone parameter or as a field in the matrix descriptor `cusparseMatDescr_t` type.

3.2. Vector Formats

This section describes dense and sparse vector formats.

3.2.1. Dense Format

Dense vectors are represented with a single data array that is stored linearly in memory, such as the following 7×1 dense vector.

[1.0 0.0 0.0 2.0 3.0 0.0 4.0]

(This vector is referenced again in the next section.)

3.2.2. Sparse Format

Sparse vectors are represented with two arrays.

- ▶ The *data array* has the nonzero values from the equivalent array in dense format.
- ▶ The *integer index array* has the positions of the corresponding nonzero values in the equivalent array in dense format.

For example, the dense vector in section 3.2.1 can be stored as a sparse vector with one-based indexing.

$$\begin{bmatrix} 1.0 & 2.0 & 3.0 & 4.0 \\ 1 & 4 & 5 & 7 \end{bmatrix}$$

It can also be stored as a sparse vector with zero-based indexing.

$$\begin{bmatrix} 1.0 & 2.0 & 3.0 & 4.0 \\ 0 & 3 & 4 & 6 \end{bmatrix}$$

In each example, the top row is the data array and the bottom row is the index array, and it is assumed that the indices are provided in increasing order and that each index appears only once.

3.3. Matrix Formats

Dense and several sparse formats for matrices are discussed in this section.

3.3.1. Dense Format

The dense matrix \mathbf{x} is assumed to be stored in column-major format in memory and is represented by the following parameters.

m	(integer)	The number of rows in the matrix.
n	(integer)	The number of columns in the matrix.
ldx	(integer)	The leading dimension of \mathbf{x} , which must be greater than or equal to m . If ldx is greater than m , then \mathbf{x} represents a sub-matrix of a larger matrix stored in memory
x	(pointer)	Points to the data array containing the matrix elements. It is assumed that enough storage is allocated for \mathbf{x} to hold all of the matrix elements and that CUSPARSE library functions may access values outside of the sub-matrix, but will never overwrite them.

For example, $\mathbf{m} \times \mathbf{n}$ dense matrix \mathbf{x} with leading dimension **ldx** can be stored with one-based indexing as shown.

$$\begin{bmatrix} X_{1,1} & X_{1,2} & \cdots & X_{1,n} \\ X_{2,1} & X_{2,2} & \cdots & X_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ X_{m,1} & X_{m,2} & \cdots & X_{m,n} \\ \vdots & \vdots & \ddots & \vdots \\ X_{ldx,1} & X_{ldx,2} & \cdots & X_{ldx,n} \end{bmatrix}$$

Its elements are arranged linearly in memory in the order below.

$$[X_{1,1} \ X_{2,1} \ \cdots \ X_{m,1} \ \cdots \ X_{ldx,1} \ \cdots \ X_{1,n} \ X_{2,n} \ \cdots \ X_{m,n} \ \cdots \ X_{ldx,n}]$$



This format and notation are similar to those used in the NVIDIA CUDA CUBLAS library.

3.3.2. Coordinate Format (COO)

The $m \times n$ sparse matrix **A** is represented in COO format by the following parameters.

nnz	(integer)	The number of nonzero elements in the matrix.
cooValA	(pointer)	Points to the data array of length nnz that holds all nonzero values of A in row-major format.
cooRowIndA	(pointer)	Points to the integer array of length nnz that contains the row indices of the corresponding elements in array cooValA .
cooColIndA	(pointer)	Points to the integer array of length nnz that contains the column indices of the corresponding elements in array cooValA .

A sparse matrix in COO format is assumed to be stored in row-major format: the index arrays are first sorted by row indices and then within the same row by compressed column indices. It is assumed that each pair of row and column indices appears only once.

For example, consider the following 4×5 matrix **A**.

$$\begin{bmatrix} 1.0 & 4.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 2.0 & 3.0 & 0.0 & 0.0 \\ 5.0 & 0.0 & 0.0 & 7.0 & 8.0 \\ 0.0 & 0.0 & 9.0 & 0.0 & 6.0 \end{bmatrix}$$

It is stored in COO format with zero-based indexing this way.

```
cooValA = [1.0  4.0  2.0  3.0  5.0  7.0  8.0  9.0  6.0]
cooRowIndA = [0   0   1   1   2   2   2   3   3 ]
cooColIndA = [0   1   1   2   0   3   4   2   4 ]
```

In the COO format with one-based indexing, it is stored as shown.

```
cooValA = [1.0  4.0  2.0  3.0  5.0  7.0  8.0  9.0  6.0]
cooRowIndA = [1   1   2   2   3   3   3   4   4 ]
cooColIndA = [1   2   2   3   1   4   5   3   5 ]
```

3.3.3. Compressed Sparse Row Format (CSR)

The only way the CSR differs from the COO format is that the array containing the row indices is compressed in CSR format. The $m \times n$ sparse matrix **A** is represented in CSR format by the following parameters.

nnz	(integer)	The number of nonzero elements in the matrix.
csrValA	(pointer)	Points to the data array of length nnz that holds all nonzero values of A in row-major format.
csrRowPtrA	(pointer)	Points to the integer array of length m+1 that holds indices into the arrays csrColIndA and csrValA . The first m entries of this array contain the indices of the first nonzero element in the <i>i</i> th row for <i>i</i> =1, ..., <i>m</i> , while the last entry contains nnz+csrRowPtrA(0) . In general, csrRowPtrA(0) is 0 or 1 for zero- and one-based indexing, respectively.

csrColIndA	(pointer)	Points to the integer array of length nnz that contains the column indices of the corresponding elements in array csrValA .
-------------------	-----------	---

Sparse matrices in CSR format are assumed to be stored in row-major CSR format, in other words, the index arrays are first sorted by row indices and then within the same row by column indices. It is assumed that each pair of row and column indices appears only once.

Consider again the 4×5 matrix **A**.

$$\begin{bmatrix} 1.0 & 4.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 2.0 & 3.0 & 0.0 & 0.0 \\ 5.0 & 0.0 & 0.0 & 7.0 & 8.0 \\ 0.0 & 0.0 & 9.0 & 0.0 & 6.0 \end{bmatrix}$$

It is stored in CSR format with zero-based indexing as shown.

```
csrValA = [1.0  4.0  2.0  3.0  5.0  7.0  8.0  9.0  6.0]
csrRowPtrA = [0   2   4   7   9 ]
csrColIndA = [0   1   1   2   0   3   4   2   4 ]
```

This is how it is stored in CSR format with one-based indexing.

```
csrValA = [1.0  4.0  2.0  3.0  5.0  7.0  8.0  9.0  6.0]
csrRowPtrA = [1   3   5   8  10 ]
csrColIndA = [1   2   2   3   1   4   5   3   5 ]
```

3.3.4. Compressed Sparse Column Format (CSC)

The CSC format is different from the COO format in two ways: the matrix is stored in column-major format, and the array containing the column indices is compressed in CSC format. The $m \times n$ matrix **A** is represented in CSC format by the following parameters.

nnz	(integer)	The number of nonzero elements in the matrix.
cscValA	(pointer)	Points to the data array of length nnz that holds all nonzero values of A in column-major format.
cscRowIndA	(pointer)	Points to the integer array of length nnz that contains the row indices of the corresponding elements in array cscValA .
cscColPtrA	(pointer)	Points to the integer array of length n+1 that holds indices into the arrays cscRowIndA and cscValA . The first n entries of this array contain the indices of the first nonzero element in the i th row for $i=1, \dots, n$, while the last entry contains nnz+cscColPtrA(0) . In general, cscColPtrA(0) is 0 or 1 for zero- and one-based indexing, respectively.



The matrix **A** in CSR format has exactly the same memory layout as its transpose in CSC format (and vice versa).

For example, consider once again the 4×5 matrix **A**.

$$\begin{bmatrix} 1.0 & 4.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 2.0 & 3.0 & 0.0 & 0.0 \\ 5.0 & 0.0 & 0.0 & 7.0 & 8.0 \\ 0.0 & 0.0 & 9.0 & 0.0 & 6.0 \end{bmatrix}$$

It is stored in CSC format with zero-based indexing this way.

```
cscValA = [1.0  5.0  4.0  2.0  3.0  9.0  7.0  8.0  6.0]
cscRowIndA = [0   2   0   1   1   3   2   2   3 ]
cscColPtrA = [0   2   4   6   7   9 ]
```

In CSC format with one-based indexing, this is how it is stored.

```
cscValA = [1.0  5.0  4.0  2.0  3.0  9.0  7.0  8.0  6.0]
cscRowIndA = [1   3   1   2   2   4   3   3   4 ]
cscColPtrA = [1   3   5   7   8  10 ]
```

Each pair of row and column indices appears only once.

3.3.5. Ellpack-Itpack Format (ELL)

An $\mathbf{m} \times \mathbf{n}$ sparse matrix \mathbf{A} with at most \mathbf{k} nonzero elements per row is stored in the Ellpack-Itpack (ELL) format [2] using two dense arrays of dimension $\mathbf{m} \times \mathbf{k}$. The first data array contains the values of the nonzero elements in the matrix, while the second integer array contains the corresponding column indices.

For example, consider the 4×5 matrix \mathbf{A} .

$$\begin{bmatrix} 1.0 & 4.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 2.0 & 3.0 & 0.0 & 0.0 \\ 5.0 & 0.0 & 0.0 & 7.0 & 8.0 \\ 0.0 & 0.0 & 9.0 & 0.0 & 6.0 \end{bmatrix}$$

This is how it is stored in ELL format with zero-based indexing.

$$\begin{aligned} \text{data} &= \begin{bmatrix} 1.0 & 4.0 & 0.0 \\ 2.0 & 3.0 & 0.0 \\ 5.0 & 7.0 & 8.0 \\ 9.0 & 6.0 & 0.0 \end{bmatrix} \\ \text{indices} &= \begin{bmatrix} 0 & 1 & -1 \\ 1 & 2 & -1 \\ 0 & 3 & 4 \\ 2 & 4 & -1 \end{bmatrix} \end{aligned}$$

It is stored this way in ELL format with one-based indexing.

$$\begin{aligned} \text{data} &= \begin{bmatrix} 1.0 & 4.0 & 0.0 \\ 2.0 & 3.0 & 0.0 \\ 5.0 & 7.0 & 8.0 \\ 9.0 & 6.0 & 0.0 \end{bmatrix} \\ \text{indices} &= \begin{bmatrix} 1 & 2 & -1 \\ 2 & 3 & -1 \\ 1 & 4 & 5 \\ 3 & 5 & -1 \end{bmatrix} \end{aligned}$$

Sparse matrices in ELL format are assumed to be stored in column-major format in memory. Also, rows with less than **k** nonzero elements are padded in the **data** and **indices** arrays with zero and -1 , respectively.

The ELL format is not supported directly, but it is used to store the regular part of the matrix in the HYB format that is described in the next section.

3.3.6. Hybrid Format (HYB)

The HYB sparse storage format is composed of a regular part, usually stored in ELL format, and an irregular part, usually stored in COO format [1]. The ELL and COO parts are always stored using zero-based indexing. HYB is implemented as an opaque data format that requires the use of a conversion operation to store a matrix in it. The conversion operation partitions the general matrix into the regular and irregular parts automatically or according to developer-specified criteria.

For more information, please refer to the description of `cusparseHybPartition_t` type, as well as the description of the conversion routines `dense2hyb` and `csr2hyb`.

3.3.7. Block Compressed Sparse Row Format (BSR)

The only difference between the CSR and BSR formats is the format of the storage element. The former stores primitive data types (**single**, **double**, **cuComplex**, and **cuDoubleComplex**) whereas the latter stores a two-dimensional square block of primitive data types. The dimension of the square block is *blockDim*. The $m \times n$ sparse matrix **A** is equivalent to a block sparse matrix A_b with $mb = \frac{m + blockDim - 1}{blockDim}$ block rows and $nb = \frac{n + blockDim - 1}{blockDim}$ block columns. If *m* or *n* is not multiple of *blockDim*, then zeros are filled into A_b .

A is represented in BSR format by the following parameters.

blockDim	(integer)	Block dimension of matrix A .
mb	(integer)	The number of block rows of A .
nb	(integer)	The number of block columns of A .
nnzb	(integer)	The number of nonzero blocks in the matrix.
bsrValA	(pointer)	Points to the data array of length $nnzb * blockDim^2$ that holds all elements of nonzero blocks of A . The block elements are stored in either column-major order or row-major order.
bsrRowPtrA	(pointer)	Points to the integer array of length $mb+1$ that holds indices into the arrays bsrColIndA and bsrValA . The first mb entries of this array contain the indices of the first nonzero block in the <i>i</i> th block row for $i=1, \dots, mb$, while the last entry contains $nnzb + bsrRowPtrA(0)$. In general, $bsrRowPtrA(0)$ is 0 or 1 for zero- and one-based indexing, respectively.
bsrColIndA	(pointer)	Points to the integer array of length $nnzb$ that contains the column indices of the corresponding blocks in array bsrValA .

As with CSR format, (row, column) indices of BSR are stored in row-major order. The index arrays are first sorted by row indices and then within the same row by column indices.

For example, consider again the 4×5 matrix \mathbf{A} .

$$\begin{bmatrix} 1.0 & 4.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 2.0 & 3.0 & 0.0 & 0.0 \\ 5.0 & 0.0 & 0.0 & 7.0 & 8.0 \\ 0.0 & 0.0 & 9.0 & 0.0 & 6.0 \end{bmatrix}$$

If *blockDim* is equal to 2, then *mb* is 2, *nb* is 3, and matrix \mathbf{A} is split into 2×3 block matrix A_b . The dimension of A_b is 4×6 , slightly bigger than matrix A , so zeros are filled in the last column of A_b . The element-wise view of A_b is this.

$$\begin{bmatrix} 1.0 & 4.0 & 0.0 & 0.0 & 0.0 & 0.0 \\ 0.0 & 2.0 & 3.0 & 0.0 & 0.0 & 0.0 \\ 5.0 & 0.0 & 0.0 & 7.0 & 8.0 & 0.0 \\ 0.0 & 0.0 & 9.0 & 0.0 & 6.0 & 0.0 \end{bmatrix}$$

Based on zero-based indexing, the block-wise view of A_b can be represented as follows.

$$A_b = \begin{bmatrix} A_{00} & A_{01} & A_{02} \\ A_{10} & A_{11} & A_{12} \end{bmatrix}$$

The basic element of BSR is a nonzero A_{ij} block, one that contains at least one nonzero element of \mathbf{A} . Five of six blocks are nonzero in A_b .

$$A_{00} = \begin{bmatrix} 1 & 4 \\ 0 & 2 \end{bmatrix}, A_{01} = \begin{bmatrix} 0 & 0 \\ 3 & 0 \end{bmatrix}, A_{10} = \begin{bmatrix} 5 & 0 \\ 0 & 0 \end{bmatrix}, A_{11} = \begin{bmatrix} 0 & 7 \\ 9 & 0 \end{bmatrix}, A_{12} = \begin{bmatrix} 8 & 0 \\ 6 & 0 \end{bmatrix}$$

BSR format only stores the information of nonzero blocks, including block indices (*i, j*) and values A_{ij} . Also row indices are compressed in CSR format.

$$\begin{aligned} \text{bsrValA} &= [A_{00} \ A_{01} \ A_{10} \ A_{11} \ A_{12}] \\ \text{bsrRowPtrA} &= [0 \ 2 \ 5] \\ \text{bsrColIndA} &= [0 \ 1 \ 0 \ 1 \ 2] \end{aligned}$$

There are two ways to arrange the data element of block A_{ij} : row-major order and column-major order. Under column-major order, the physical storage of **bsrValA** is this.

$$\text{bsrValA} = [1 \ 0 \ 4 \ 2 \ | \ 0 \ 3 \ 0 \ 0 \ | \ 5 \ 0 \ 0 \ 0 \ | \ 0 \ 9 \ 7 \ 0 \ | \ 8 \ 6 \ 0 \ 0 \]$$

Under row-major order, the physical storage of **bsrValA** is this.

$$\text{bsrValA} = [1 \ 4 \ 0 \ 2 \ | \ 0 \ 0 \ 3 \ 0 \ | \ 5 \ 0 \ 0 \ 0 \ | \ 0 \ 7 \ 9 \ 0 \ | \ 8 \ 0 \ 6 \ 0 \]$$

Similarly, in BSR format with one-based indexing and column-major order, \mathbf{A} can be represented by the following.

$$A_b = \begin{bmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \end{bmatrix}$$

$$\text{bsrValA} = [1 \ 0 \ 4 \ 2 \ | \ 0 \ 3 \ 0 \ 0 \ | \ 5 \ 0 \ 0 \ 0 \ | \ 0 \ 9 \ 7 \ 0 \ | \ 8 \ 6 \ 0 \ 0 \]$$

```
bsrRowPtrA = [ 1   3   6]
bsrColIndA = [ 1   2   1   2   3]
```



The storage format of blocks in BSR format can be column-major or row-major, independently of the base index. However, if the developer has BSR format from the Math Kernel Library (MKL) and wants to directly copy it to BSR in CUSPARSE, then `cusparseDirection_t` is `CUSPARSE_DIRECTION_COLUMN` if the base index is one; otherwise, `cusparseDirection_t` is `CUSPARSE_DIRECTION_ROW`.

3.3.8. Extended BSR Format (BSRX)

BSRX is the same as the BSR format, but the array `bsrRowPtrA` is separated into two parts. The first nonzero block of each row is still specified by the array `bsrRowPtrA`, which is the same as in BSR, but the position next to the last nonzero block of each row is specified by the array `bsrEndPtrA`. Briefly, BSRX format is simply like a 4-vector variant of BSR format.

Matrix **A** is represented in BSRX format by the following parameters.

blockDim	(integer)	Block dimension of matrix A .
mb	(integer)	The number of block rows of A .
nb	(integer)	The number of block columns of A .
nnzb	(integer)	The size of <code>bsrColIndA</code> and <code>bsrValA</code> ; nnzb is greater than or equal to the number of nonzero blocks in the matrix A .
bsrValA	(pointer)	Points to the data array of length $nnzb * blockDim^2$ that holds all the elements of the nonzero blocks of A . The block elements are stored in either column-major order or row-major order.
bsrRowPtrA	(pointer)	Points to the integer array of length mb that holds indices into the arrays <code>bsrColIndA</code> and <code>bsrValA</code> ; <code>bsrRowPtr(i)</code> is the position of the first nonzero block of the <i>i</i> th block row in <code>bsrColIndA</code> and <code>bsrValA</code> .
bsrEndPtrA	(pointer)	Points to the integer array of length mb that holds indices into the arrays <code>bsrColIndA</code> and <code>bsrValA</code> ; <code>bsrRowPtr(i)</code> is the position next to the last nonzero block of the <i>i</i> th block row in <code>bsrColIndA</code> and <code>bsrValA</code> .
bsrColIndA	(pointer)	Points to the integer array of length nnzb that contains the column indices of the corresponding blocks in array <code>bsrValA</code> .

A simple conversion between BSR and BSRX can be done as follows. Suppose the developer has a 2×3 block sparse matrix A_b represented as shown.

$$A_b = \begin{bmatrix} A_{00} & A_{01} & A_{02} \\ A_{10} & A_{11} & A_{12} \end{bmatrix}$$

Assume it has this BSR format.

```
bsrValA of BSR = [A00 A01 A10 A11 A12]
bsrRowPtrA of BSR = [0   2   5]
bsrColIndA of BSR = [0   1   0   1   2]
```


The **bsrRowPtrA** of the BSRX format is simply the first two elements of the **bsrRowPtrA** BSR format. The **bsrEndPtrA** of BSRX format is the last two elements of the **bsrRowPtrA** of BSR format.

$$\begin{aligned}\text{bsrRowPtrA of BSRX} &= [0 \quad 2] \\ \text{bsrEndPtrA of BSRX} &= [2 \quad 5]\end{aligned}$$

The power of the BSRX format is that the developer can specify a submatrix in the original BSR format by modifying **bsrRowPtrA** and **bsrEndPtrA** while keeping **bsrColIndA** and **bsrValA** unchanged.

For example, to create another block matrix $\tilde{A} = \begin{bmatrix} O & O & O \\ O & A_{11} & O \end{bmatrix}$ that is slightly different from A , the developer can keep **bsrColIndA** and **bsrValA**, but reconstruct \tilde{A} by properly setting of **bsrRowPtrA** and **bsrEndPtrA**. The following 4-vector characterizes \tilde{A} .

$$\begin{aligned}\text{bsrValA of } \tilde{A} &= [A_{00} \ A_{01} \ A_{10} \ A_{11} \ A_{12}] \\ \text{bsrColIndA of } \tilde{A} &= [0 \quad 1 \quad 0 \quad 1 \quad 2] \\ \text{bsrRowPtrA of } \tilde{A} &= [0 \quad 3] \\ \text{bsrEndPtrA of } \tilde{A} &= [0 \quad 4]\end{aligned}$$

Chapter 4.

CUSPARSE TYPES REFERENCE

4.1. Data types

The `float`, `double`, `cuComplex`, and `cuDoubleComplex` data types are supported. The first two are standard C data types, while the last two are exported from `cuComplex.h`.

4.2. `cusparseAction_t`

This type indicates whether the operation is performed only on indices or on data and indices.

Value	Meaning
<code>CUSPARSE_ACTION_SYMBOLIC</code>	the operation is performed only on indices.
<code>CUSPARSE_ACTION_NUMERIC</code>	the operation is performed on data and indices.

4.3. `cusparseDirection_t`

This type indicates whether the elements of a dense matrix should be parsed by rows or by columns (assuming column-major storage in memory of the dense matrix) in function `cusparse[S|D|C|Z]nnz`. Besides storage format of blocks in BSR format is also controlled by this type.

Value	Meaning
<code>CUSPARSE_DIRECTION_ROW</code>	the matrix should be parsed by rows.
<code>CUSPARSE_DIRECTION_COLUMN</code>	the matrix should be parsed by columns.

4.4. `cusparseHandle_t`

This is a pointer type to an opaque CUSPARSE context, which the user must initialize by calling prior to calling `cusparseCreate()` any other library function. The handle created and returned by `cusparseCreate()` must be passed to every CUSPARSE function.

4.5. `cusparseHybMat_t`

This is a pointer type to an opaque structure holding the matrix in HYB format, which is created by `cusparseCreateHybMat` and destroyed by `cusparseDestroyHybMat`.

4.5.1. `cusparseHybPartition_t`

This type indicates how to perform the partitioning of the matrix into regular (ELL) and irregular (COO) parts of the HYB format.

The partitioning is performed during the conversion of the matrix from a dense or sparse format into the HYB format and is governed by the following rules. When `CUSPARSE_HYB_PARTITION_AUTO` is selected, the CUSPARSE library automatically decides how much data to put into the regular and irregular parts of the HYB format. When `CUSPARSE_HYB_PARTITION_USER` is selected, the width of the regular part of the HYB format should be specified by the caller. When `CUSPARSE_HYB_PARTITION_MAX` is selected, the width of the regular part of the HYB format equals to the maximum number of non-zero elements per row, in other words, the entire matrix is stored in the regular part of the HYB format.

The *default* is to let the library automatically decide how to split the data.

Value	Meaning
<code>CUSPARSE_HYB_PARTITION_AUTO</code>	the automatic partitioning is selected (<i>default</i>).
<code>CUSPARSE_HYB_PARTITION_USER</code>	the user specified treshold is used.
<code>CUSPARSE_HYB_PARTITION_MAX</code>	the data is stored in ELL format.

4.6. `cusparseMatDescr_t`

This structure is used to describe the shape and properties of a matrix.

```
typedef struct {
    cusparseMatrixType_t MatrixType;
    cusparseFillMode_t FillMode;
    cusparseDiagType_t DiagType;
    cusparseIndexBase_t IndexBase;
} cusparseMatDescr_t;
```

4.6.1. `cusparseDiagType_t`

This type indicates if the matrix diagonal entries are unity. The diagonal elements are always assumed to be present, but if `CUSPARSE_DIAG_TYPE_UNIT` is passed to an API routine, then the routine will assume that all diagonal entries are unity and will not read or modify those entries. Note that in this case the routine assumes the diagonal entries are equal to one, regardless of what those entries are actually set to in memory.

Value	Meaning
<code>CUSPARSE_DIAG_TYPE_NON_UNIT</code>	the matrix diagonal has non-unit elements.
<code>CUSPARSE_DIAG_TYPE_UNIT</code>	the matrix diagonal has unit elements.

4.6.2. `cusparseFillMode_t`

This type indicates if the lower or upper part of a matrix is stored in sparse storage.

Value	Meaning
<code>CUSPARSE_FILL_MODE_LOWER</code>	the lower triangular part is stored.
<code>CUSPARSE_FILL_MODE_UPPER</code>	the upper triangular part is stored.

4.6.3. `cusparseIndexBase_t`

This type indicates if the base of the matrix indices is zero or one.

Value	Meaning
<code>CUSPARSE_INDEX_BASE_ZERO</code>	the base index is zero.
<code>CUSPARSE_INDEX_BASE_ONE</code>	the base index is one.

4.6.4. `cusparseMatrixType_t`

This type indicates the type of matrix stored in sparse storage. Notice that for symmetric, Hermitian and triangular matrices only their lower or upper part is assumed to be stored.

Value	Meaning
<code>CUSPARSE_MATRIX_TYPE_GENERAL</code>	the matrix is general.
<code>CUSPARSE_MATRIX_TYPE_SYMMETRIC</code>	the matrix is symmetric.
<code>CUSPARSE_MATRIX_TYPE_HERMITIAN</code>	the matrix is Hermitian.
<code>CUSPARSE_MATRIX_TYPE_TRIANGULAR</code>	the matrix is triangular.

4.7. `cusparseOperation_t`

This type indicates which operations need to be performed with the sparse matrix.

Value	Meaning
<code>CUSPARSE_OPERATION_NON_TRANSPOSE</code>	the non-transpose operation is selected.
<code>CUSPARSE_OPERATION_TRANSPOSE</code>	the transpose operation is selected.
<code>CUSPARSE_OPERATION_CONJUGATE_TRANSPOSE</code>	the conjugate transpose operation is selected.

4.8. `cusparsePointerMode_t`

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

Value	Meaning
<code>CUSPARSE_POINTER_MODE_HOST</code>	the scalars are passed by reference on the host.
<code>CUSPARSE_POINTER_MODE_DEVICE</code>	the scalars are passed by reference on the device.

4.9. `cusparseSolveAnalysisInfo_t`

This is a pointer type to an opaque structure holding the information collected in the analysis phase of the solution of the sparse triangular linear system. It is expected to be passed unchanged to the solution phase of the sparse triangular linear system.

4.10. `cusparseStatus_t`

This is a status type returned by the library functions and it can have the following values.

<code>CUSPARSE_STATUS_SUCCESS</code>	The operation completed successfully.
<code>CUSPARSE_STATUS_NOT_INITIALIZED</code>	<p>The CUSPARSE library was not initialized. This is usually caused by the lack of a prior call, an error in the CUDA Runtime API called by the CUSPARSE routine, or an error in the hardware setup.</p> <p>To correct: call <code>cusparseCreate()</code> prior to the function call; and check that the hardware, an appropriate version of the driver, and the CUSPARSE library are correctly installed.</p>

CUSPARSE_STATUS_ALLOC_FAILED	<p>Resource allocation failed inside the CUSPARSE library. This is usually caused by a <code>cudaMalloc()</code> failure.</p> <p>To correct: prior to the function call, deallocate previously allocated memory as much as possible.</p>
CUSPARSE_STATUS_INVALID_VALUE	<p>An unsupported value or parameter was passed to the function (a negative vector size, for example).</p> <p>To correct: ensure that all the parameters being passed have valid values.</p>
CUSPARSE_STATUS_ARCH_MISMATCH	<p>The function requires a feature absent from the device architecture; usually caused by the lack of support for atomic operations or double precision.</p> <p>To correct: compile and run the application on a device with appropriate compute capability, which is 1.1 for 32-bit atomic operations and 1.3 for double precision.</p>
CUSPARSE_STATUS_MAPPING_ERROR	<p>An access to GPU memory space failed, which is usually caused by a failure to bind a texture.</p> <p>To correct: prior to the function call, unbind any previously bound textures.</p>
CUSPARSE_STATUS_EXECUTION_FAILED	<p>The GPU program failed to execute. This is often caused by a launch failure of the kernel on the GPU, which can be caused by multiple reasons.</p> <p>To correct: check that the hardware, an appropriate version of the driver, and the CUSPARSE library are correctly installed.</p>
CUSPARSE_STATUS_INTERNAL_ERROR	<p>An internal CUSPARSE operation failed. This error is usually caused by a <code>cudaMemcpyAsync()</code> failure.</p> <p>To correct: check that the hardware, an appropriate version of the driver, and the CUSPARSE library are correctly installed. Also, check that the memory passed as a parameter to the routine is not being deallocated prior to the routine's completion.</p>
CUSPARSE_STATUS_MATRIX_TYPE_NOT_SUPPORTED	<p>The matrix type is not supported by this function. This is usually caused by passing an invalid matrix descriptor to the function.</p> <p>To correct: check that the fields in <code>cusparseMatDescr_t descrA</code> were set correctly.</p>

Chapter 5.

CUSPARSE HELPER FUNCTION REFERENCE

The CUSPARSE helper functions are described in this section.

5.1. cusparseCreate()

```
cusparseStatus_t  
cusparseCreate(cusparseHandle_t *handle)
```

This function initializes the CUSPARSE library and creates a handle on the CUSPARSE context. It must be called before any other CUSPARSE API function is invoked. It allocates hardware resources necessary for accessing the GPU.

Output

handle	the pointer to the handle to the CUSPARSE context.
---------------	--

Status Returned

CUSPARSE_STATUS_SUCCESS	the initialization succeeded.
CUSPARSE_STATUS_NOT_INITIALIZED	the CUDA Runtime initialization failed.
CUSPARSE_STATUS_ALLOC_FAILED	the resources could not be allocated.
CUSPARSE_STATUS_ARCH_MISMATCH	the device compute capability (CC) is less than 1.1. The CC of at least 1.1 is required.

5.2. cusparseCreateHybMat()

```
cusparseStatus_t  
cusparseCreateHybMat(cusparseHybMat_t *hybA)
```

This function creates and initializes the **hybA** opaque data structure.

Input

hybA	the pointer to the hybrid format storage structure.
-------------	---

Status Returned

<code>CUSPARSE_STATUS_SUCCESS</code>	the structure was initialized successfully.
<code>CUSPARSE_STATUS_ALLOC_FAILED</code>	the resources could not be allocated.

5.3. `cusparseCreateMatDescr()`

```
cusparseStatus_t
cusparseCreateMatDescr(cusparseMatDescr_t *descrA)
```

This function initializes the matrix descriptor. It sets the fields **MatrixType** and **IndexBase** to the *default* values `CUSPARSE_MATRIX_TYPE_GENERAL` and `CUSPARSE_INDEX_BASE_ZERO`, respectively, while leaving other fields uninitialized.

Input

<code>descrA</code>	the pointer to the matrix descriptor.
---------------------	---------------------------------------

Status Returned

<code>CUSPARSE_STATUS_SUCCESS</code>	the descriptor was initialized successfully.
<code>CUSPARSE_STATUS_ALLOC_FAILED</code>	the resources could not be allocated.

5.4. `cusparseCreateSolveAnalysisInfo()`

```
cusparseStatus_t
cusparseCreateSolveAnalysisInfo(cusparseSolveAnalysisInfo_t *info)
```

This function creates and initializes the solve and analysis structure to *default* values.

Input

<code>info</code>	the pointer to the solve and analysis structure.
-------------------	--

Status Returned

<code>CUSPARSE_STATUS_SUCCESS</code>	the structure was initialized successfully.
<code>CUSPARSE_STATUS_ALLOC_FAILED</code>	the resources could not be allocated.

5.5. `cusparseDestroy()`

```
cusparseStatus_t
cusparseDestroy(cusparseHandle_t handle)
```

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

Input

<code>handle</code>	the handle to the CUSPARSE context.
---------------------	-------------------------------------

Status Returned

<code>CUSPARSE_STATUS_SUCCESS</code>	the shutdown succeeded.
--------------------------------------	-------------------------

<code>CUSPARSE_STATUS_NOT_INITIALIZED</code>	the library was not initialized.
--	----------------------------------

5.6. `cusparseDestroyHybMat()`

```
cusparseStatus_t
cusparseDestroyHybMat(cusparseHybMat_t hybA)
```

This function destroys and releases any memory required by the **hybA** structure.

Input

<code>hybA</code>	the hybrid format storage structure.
-------------------	--------------------------------------

Status Returned

<code>CUSPARSE_STATUS_SUCCESS</code>	the resources were released successfully.
--------------------------------------	---

5.7. `cusparseDestroyMatDescr()`

```
cusparseStatus_t
cusparseDestroyMatDescr(cusparseMatDescr_t descrA)
```

This function releases the memory allocated for the matrix descriptor.

Input

<code>descrA</code>	the matrix descriptor.
---------------------	------------------------

Status Returned

<code>CUSPARSE_STATUS_SUCCESS</code>	the resources were released successfully.
--------------------------------------	---

5.8. `cusparseDestroySolveAnalysisInfo()`

```
cusparseStatus_t
cusparseDestroySolveAnalysisInfo(cusparseSolveAnalysisInfo_t info)
```

This function destroys and releases any memory required by the structure.

Input

<code>info</code>	the solve and analysis structure.
-------------------	-----------------------------------

Status Returned

<code>CUSPARSE_STATUS_SUCCESS</code>	the resources were released successfully.
--------------------------------------	---

5.9. `cusparseGetLevelInfo()`

```
cusparseStatus_t
cusparseGetLevelInfo(cusparseHandle_t handle,
                    cusparseSolveAnalysisInfo_t info,
```

```
int *nlevels,
int **levelPtr,
int **levelInd)
```

This function returns the number of levels and the assignment of rows into the levels computed by either the `csrsv_analysis`, `csrsm_analysis` or `hybsv_analysis` routines.

Input

handle	handle to the CUSPARSE library context.
info	the pointer to the solve and analysis structure.

Output

nlevels	number of levels.
levelPtr	integer array of nlevels+1 elements that contains the start of every level and the end of the last level plus one.
levelInd	integer array of m (number of rows in the matrix) elements that contains the row indices belonging to every level.

Status Returned

CUSPARSE_STATUS_SUCCESS	the structure was initialized successfully.
CUSPARSE_STATUS_NOT_INITIALIZED	the library or the solve analysis structure was not initialized.

5.10. `cusparseGetMatDiagType()`

```
cusparseDiagType_t
cusparseGetMatDiagType(const cusparseMatDescr_t descrA)
```

This function returns the **DiagType** field of the matrix descriptor **descrA**.

Input

descrA	the matrix descriptor.
---------------	------------------------

Returned

	One of the enumerated <code>diagType</code> types.
--	--

5.11. `cusparseGetMatFillMode()`

```
cusparseFillMode_t
cusparseGetMatFillMode(const cusparseMatDescr_t descrA)
```

This function returns the **FillMode** field of the matrix descriptor **descrA**.

Input

descrA	the matrix descriptor.
---------------	------------------------

Returned

	One of the enumerated fillMode types.
--	---------------------------------------

5.12. `cusparseGetMatIndexBase()`

```
cusparseIndexBase_t
cusparseGetMatIndexBase(const cusparseMatDescr_t descrA)
```

This function returns the **IndexBase** field of the matrix descriptor **descrA**.

Input

descrA	the matrix descriptor.
---------------	------------------------

Returned

	One of the enumerated indexBase types.
--	--

5.13. `cusparseGetMatType()`

```
cusparseMatrixType_t
cusparseGetMatType(const cusparseMatDescr_t descrA)
```

This function returns the **MatrixType** field of the matrix descriptor **descrA**.

Input

descrA	the matrix descriptor.
---------------	------------------------

Returned

	One of the enumerated matrix types.
--	-------------------------------------

5.14. `cusparseGetPointerMode()`

```
cusparseStatus_t
cusparseGetPointerMode(cusparseHandle_t handle,
                       cusparsePointerMode_t *mode)
```

This function obtains the pointer mode used by the CUSPARSE library. Please see the section on the **cusparsePointerMode_t** type for more details.

Input

handle	the handle to the CUSPARSE context.
---------------	-------------------------------------

Output

mode	One of the enumerated pointer mode types.
-------------	---

Status Returned

CUSPARSE_STATUS_SUCCESS	the pointer mode was returned successfully.
--------------------------------	---

<code>CUSPARSE_STATUS_NOT_INITIALIZED</code>	the library was not initialized.
--	----------------------------------

5.15. `cusparseGetVersion()`

```
cusparseStatus_t
cusparseGetVersion(cusparseHandle_t handle, int *version)
```

This function returns the version number of the CUSPARSE library.

Input

handle	the handle to the CUSPARSE context.
---------------	-------------------------------------

Output

version	the version number of the library.
----------------	------------------------------------

Status Returned

<code>CUSPARSE_STATUS_SUCCESS</code>	the version was returned successfully.
<code>CUSPARSE_STATUS_NOT_INITIALIZED</code>	the library was not initialized.

5.16. `cusparseSetMatDiagType()`

```
cusparseStatus_t
cusparseSetMatDiagType(cusparseMatDescr_t descrA,
                       cusparseDiagType_t diagType)
```

This function sets the **DiagType** field of the matrix descriptor **descrA**.

Input

diagType	One of the enumerated diagType types.
-----------------	--

Output

descrA	the matrix descriptor.
---------------	------------------------

Status Returned

<code>CUSPARSE_STATUS_SUCCESS</code>	the field DiagType was set successfully.
<code>CUSPARSE_STATUS_INVALID_VALUE</code>	An invalid diagType parameter was passed.

5.17. `cusparseSetMatFillMode()`

```
cusparseStatus_t
cusparseSetMatFillMode(cusparseMatDescr_t descrA,
                       cusparseFillMode_t fillMode)
```

This function sets the **FillMode** field of the matrix descriptor **descrA**.

Input

fillMode	One of the enumerated fillMode types.
-----------------	---------------------------------------

Output

descrA	the matrix descriptor.
---------------	------------------------

Status Returned

CUSPARSE_STATUS_SUCCESS	the FillMode field was set successfully.
CUSPARSE_STATUS_INVALID_VALUE	An invalid fillMode parameter was passed.

5.18. `cusparseSetMatIndexBase()`

```
cusparseStatus_t
cusparseSetMatIndexBase(cusparseMatDescr_t descrA,
                        cusparseIndexBase_t base)
```

This function sets the **IndexBase** field of the matrix descriptor **descrA**.

Input

base	One of the enumerated indexBase types.
-------------	--

Output

descrA	the matrix descriptor.
---------------	------------------------

Status Returned

CUSPARSE_STATUS_SUCCESS	the IndexBase field was set successfully.
CUSPARSE_STATUS_INVALID_VALUE	An invalid base parameter was passed.

5.19. `cusparseSetMatType()`

```
cusparseStatus_t
cusparseSetMatType(cusparseMatDescr_t descrA, cusparseMatrixType_t type)
```

This function sets the **MatrixType** field of the matrix descriptor **descrA**.

Input

type	One of the enumerated matrix types.
-------------	-------------------------------------

Output

descrA	the matrix descriptor.
---------------	------------------------

Status Returned

CUSPARSE_STATUS_SUCCESS	the MatrixType field was set successfully.
CUSPARSE_STATUS_INVALID_VALUE	An invalid type parameter was passed.

5.20. `cusparseSetPointerMode()`

```
cusparseStatus_t
cusparseSetPointerMode(cusparseHandle_t handle,
                      cusparsePointerMode_t mode)
```

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

Input

handle	the handle to the CUSPARSE context.
mode	One of the enumerated pointer mode types.

Status Returned

CUSPARSE_STATUS_SUCCESS	the pointer mode was set successfully.
CUSPARSE_STATUS_INVALID_VALUE	the library was not initialized.

5.21. `cusparseSetStream()`

```
cusparseStatus_t
cusparseSetStream(cusparseHandle_t handle, cudaStream_t streamId)
```

This function sets the stream to be used by the CUSPARSE library to execute its routines.

Input

handle	the handle to the CUSPARSE context.
streamId	the stream to be used by the library.

Status Returned

CUSPARSE_STATUS_SUCCESS	the stream was set successfully.
CUSPARSE_STATUS_INVALID_VALUE	the library was not initialized.

Chapter 6.

CUSPARSE LEVEL 1 FUNCTION REFERENCE

This chapter describes sparse linear algebra functions that perform operations between dense and sparse vectors.

6.1. `cusparse<t>axpyi`

```
cusparseStatus_t
cusparseSaxpyi(cusparseHandle_t handle, int nnz,
               const float      *alpha,
               const float      *xVal, const int *xInd,
               float            *y, cusparseIndexBase_t idxBase)
cusparseStatus_t
cusparseDaxpyi(cusparseHandle_t handle, int nnz,
               const double     *alpha,
               const double     *xVal, const int *xInd,
               double           *y, cusparseIndexBase_t idxBase)
cusparseStatus_t
cusparseCaxpyi(cusparseHandle_t handle, int nnz,
               const cuComplex  *alpha,
               const cuComplex  *xVal, const int *xInd,
               cuComplex        *y, cusparseIndexBase_t idxBase)
cusparseStatus_t
cusparseZaxpyi(cusparseHandle_t handle, int nnz,
               const cuDoubleComplex *alpha,
               const cuDoubleComplex *xVal, const int *xInd,
               cuDoubleComplex *y, cusparseIndexBase_t idxBase)
```

This function multiplies the vector **x** in sparse format by the constant α and adds the result to the vector **y** in dense format. This operation can be written as

$$y = y + \alpha * x$$

in other words,

```
for i=0 to nnz-1
    y[xInd[i]-idxBase] = y[xInd[i]-idxBase] + alpha*xVal[i]
```

This function requires no extra storage. It is executed asynchronously with respect to the host and it may return control to the application on the host before the result is ready.

Input

handle	handle to the CUSPARSE library context.
nnz	number of elements in vector x .
alpha	<type> scalar used for multiplication.
xVal	<type> vector with nnz non-zero values of vector x .
xInd	integer vector with nnz indices of the non-zero values of vector x .
y	<type> vector in dense format.
idxBase	CUSPARSE_INDEX_BASE_ZERO or CUSPARSE_INDEX_BASE_ONE

Output

y	<type> updated vector in dense format (that is unchanged if nnz == 0).
----------	---

Status Returned

CUSPARSE_STATUS_SUCCESS	the operation completed successfully.
CUSPARSE_STATUS_NOT_INITIALIZED	the library was not initialized.
CUSPARSE_STATUS_INVALID_VALUE	the idxBase is neither CUSPARSE_INDEX_BASE_ZERO nor CUSPARSE_INDEX_BASE_ONE.
CUSPARSE_STATUS_ARCH_MISMATCH	the device does not support double precision.
CUSPARSE_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU.

6.2. `cusparse<t>doti`

```

cusparseStatus_t
cusparseSdoti(cusparseHandle_t handle, int nnz,
              const float          *xVal,
              const int *xInd, const float          *y,
              float          *resultDevHostPtr,
              cusparseIndexBase_t idxBase)
cusparseStatus_t
cusparseDdoti(cusparseHandle_t handle, int nnz,
              const double         *xVal,
              const int *xInd, const double         *y,
              double         *resultDevHostPtr,
              cusparseIndexBase_t idxBase)
cusparseStatus_t
cusparseCdoti(cusparseHandle_t handle, int nnz,
              const cuComplex      *xVal,
              const int *xInd, const cuComplex      *y,
              cuComplex      *resultDevHostPtr,
              cusparseIndexBase_t idxBase)
cusparseStatus_t
cusparseZdoti(cusparseHandle_t handle, int nnz, const
              cuDoubleComplex *xVal,
              const int *xInd, const cuDoubleComplex *y,
              cuDoubleComplex *resultDevHostPtr,

```



```
cusparseIndexBase_t idxBase)
```

This function returns the dot product of a vector **x** in sparse format and vector **y** in dense format. This operation can be written as

$$result = y^T x$$

in other words,

```
for i=0 to nnz-1
    resultDevHostPtr += xVal[i]*y[xInd[i-idxBase]]
```

This function requires some temporary extra storage that is allocated internally. It is executed asynchronously with respect to the host and it may return control to the application on the host before the result is ready.

Input

handle	handle to the CUSPARSE library context.
nnz	number of elements in vector x .
xVal	<type> vector with nnz non-zero values of vector x .
xInd	integer vector with nnz indices of the non-zero values of vector x .
y	<type> vector in dense format.
resultDevHostPtr	pointer to the location of the result in the device or host memory.
idxBase	CUSPARSE_INDEX_BASE_ZERO or CUSPARSE_INDEX_BASE_ONE

Output

resultDevHostPtr	scalar result in the device or host memory (that is zero if nnz == 0).
-------------------------	---

Status Returned

CUSPARSE_STATUS_SUCCESS	the operation completed successfully.
CUSPARSE_STATUS_NOT_INITIALIZED	the library was not initialized.
CUSPARSE_STATUS_ALLOC_FAILED	the reduction buffer could not be allocated.
CUSPARSE_STATUS_INVALID_VALUE	the idxBase is neither CUSPARSE_INDEX_BASE_ZERO nor CUSPARSE_INDEX_BASE_ONE.
CUSPARSE_STATUS_ARCH_MISMATCH	the device does not support double precision.
CUSPARSE_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU.
CUSPARSE_STATUS_INTERNAL_ERROR	an internal operation failed.

6.3. cusparse<t>dotci

```
cusparseStatus_t
```

```

cusparseCdotci(cusparseHandle_t handle, int nnz,
               const cuComplex      *xVal,
               const int *xInd, const cuComplex      *y,
               cuComplex      *resultDevHostPtr, cusparseIndexBase_t
               idxBase)
cusparseStatus_t
cusparseZdotci(cusparseHandle_t handle, int nnz,
               const cuDoubleComplex *xVal,
               const int *xInd, const cuDoubleComplex *y,
               cuDoubleComplex *resultDevHostPtr, cusparseIndexBase_t
               idxBase)

```

This function returns the dot product of a complex conjugate of vector **x** in sparse format and vector **y** in dense format. This operation can be written as

$$result = y^H x$$

in other words,

```

for i=0 to nnz-1
    resultDevHostPtr += xVal[i]*y[xInd[i-idxBase]]

```

This function requires some temporary extra storage that is allocated internally. It is executed asynchronously with respect to the host and it may return control to the application on the host before the result is ready.

Input

handle	handle to the CUSPARSE library context.
nnz	number of elements in vector x .
xVal	<type> vector with nnz non-zero values of vector x .
xInd	integer vector with nnz indices of the non-zero values of vector x .
y	<type> vector in dense format.
resultDevHostPtr	pointer to the location of the result in the device or host memory.
idxBase	CUSPARSE_INDEX_BASE_ZERO or CUSPARSE_INDEX_BASE_ONE

Output

resultDevHostPtr	scalar result in the device or host memory (that is zero if nnz == 0).
-------------------------	---

Status Returned

CUSPARSE_STATUS_SUCCESS	the operation completed successfully.
CUSPARSE_STATUS_NOT_INITIALIZED	the library was not initialized.
CUSPARSE_STATUS_ALLOC_FAILED	the reduction buffer could not be allocated.
CUSPARSE_STATUS_INVALID_VALUE	the idxBase is neither CUSPARSE_INDEX_BASE_ZERO nor CUSPARSE_INDEX_BASE_ONE.

<code>CUSPARSE_STATUS_ARCH_MISMATCH</code>	the device does not support double precision.
<code>CUSPARSE_STATUS_EXECUTION_FAILED</code>	the function failed to launch on the GPU.
<code>CUSPARSE_STATUS_INTERNAL_ERROR</code>	an internal operation failed.

6.4. `cusparse<T>gthr`

```

cusparseStatus_t
cusparseSgthr(cusparseHandle_t handle, int nnz,
              const float          *y,
              float                *xVal, const int *xInd,
              cusparseIndexBase_t idxBase)
cusparseStatus_t
cusparseDgthr(cusparseHandle_t handle, int nnz,
              const double         *y,
              double               *xVal, const int *xInd,
              cusparseIndexBase_t idxBase)
cusparseStatus_t
cusparseCgthr(cusparseHandle_t handle, int nnz,
              const cuComplex      *y,
              cuComplex            *xVal, const int *xInd,
              cusparseIndexBase_t idxBase)
cusparseStatus_t
cusparseZgthr(cusparseHandle_t handle, int nnz,
              const cuDoubleComplex *y,
              cuDoubleComplex *xVal, const int *xInd,
              cusparseIndexBase_t idxBase)

```

This function gathers the elements of the vector **y** listed in the index array **xInd** into the data array **xVal**.

This function requires no extra storage. It is executed asynchronously with respect to the host and it may return control to the application on the host before the result is ready.

Input

handle	handle to the CUSPARSE library context.
nnz	number of elements in vector x .
y	<type> vector in dense format (of size $\geq \max(\mathbf{xInd}) - \mathbf{idxBase} + 1$).
xInd	integer vector with nnz indices of the non-zero values of vector x .
idxBase	<code>CUSPARSE_INDEX_BASE_ZERO</code> or <code>CUSPARSE_INDEX_BASE_ONE</code>

Output

xVal	<type> vector with nnz non-zero values that were gathered from vector y (that is unchanged if nnz == 0).
-------------	---

Status Returned

<code>CUSPARSE_STATUS_SUCCESS</code>	the operation completed successfully.
--------------------------------------	---------------------------------------

<code>CUSPARSE_STATUS_NOT_INITIALIZED</code>	the library was not initialized.
<code>CUSPARSE_STATUS_INVALID_VALUE</code>	the <code>idxBase</code> is neither <code>CUSPARSE_INDEX_BASE_ZERO</code> nor <code>CUSPARSE_INDEX_BASE_ONE</code> .
<code>CUSPARSE_STATUS_ARCH_MISMATCH</code>	the device does not support double precision.
<code>CUSPARSE_STATUS_EXECUTION_FAILED</code>	the function failed to launch on the GPU.

6.5. `cusparse<t>gthrz`

```

cusparseStatus_t
cusparseSgthrz(cusparseHandle_t handle, int nnz, float          *y,
               float          *xVal, const int *xInd,
               cusparseIndexBase_t idxBase)
cusparseStatus_t
cusparseDgthrz(cusparseHandle_t handle, int nnz, double        *y,
               double        *xVal, const int *xInd,
               cusparseIndexBase_t idxBase)
cusparseStatus_t
cusparseCgthrz(cusparseHandle_t handle, int nnz, cuComplex      *y,
               cuComplex      *xVal, const int *xInd,
               cusparseIndexBase_t idxBase)
cusparseStatus_t
cusparseZgthrz(cusparseHandle_t handle, int nnz, cuDoubleComplex *y,
               cuDoubleComplex *xVal, const int *xInd,
               cusparseIndexBase_t idxBase)

```

This function gathers the elements of the vector **y** listed in the index array **xInd** into the data array **xVal**. Also, it zeroes out the gathered elements in the vector **y**.

This function requires no extra storage. It is executed asynchronously with respect to the host and it may return control to the application on the host before the result is ready.

Input

handle	handle to the CUSPARSE library context.
nnz	number of elements in vector x .
y	<type> vector in dense format (of <code>size ≥ max(xInd) - idxBase + 1</code>).
xInd	integer vector with nnz indices of the non-zero values of vector x .
idxBase	<code>CUSPARSE_INDEX_BASE_ZERO</code> or <code>CUSPARSE_INDEX_BASE_ONE</code>

Output

xVal	<type> vector with nnz non-zero values that were gathered from vector y (that is unchanged if <code>nnz == 0</code>).
y	<type> vector in dense format with elements indexed by xInd set to zero (it is unchanged if <code>nnz == 0</code>).

Status Returned

<code>CUSPARSE_STATUS_SUCCESS</code>	the operation completed successfully.
<code>CUSPARSE_STATUS_NOT_INITIALIZED</code>	the library was not initialized.
<code>CUSPARSE_STATUS_INVALID_VALUE</code>	the <code>idxBase</code> is neither <code>CUSPARSE_INDEX_BASE_ZERO</code> nor <code>CUSPARSE_INDEX_BASE_ONE</code> .
<code>CUSPARSE_STATUS_ARCH_MISMATCH</code>	the device does not support double precision.
<code>CUSPARSE_STATUS_EXECUTION_FAILED</code>	the function failed to launch on the GPU.

6.6. `cusparse<t>roti`

```

cusparseStatus_t
cusparseSroti(cusparseHandle_t handle, int nnz, float *xVal,
              const int *xInd,
              float *y, const float *c, const float *s,
              cusparseIndexBase_t idxBase)
cusparseStatus_t
cusparseDroti(cusparseHandle_t handle, int nnz, double *xVal,
              const int *xInd,
              double *y, const double *c, const double *s,
              cusparseIndexBase_t idxBase)

```

This function applies Givens rotation matrix

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

to sparse **x** and dense **y** vectors. In other words,

```

for i=0 to nnz-1
    y[xInd[i]-idxBase] = c * y[xInd[i]-idxBase] - s*xVal[i]
    x[i]                = c * xVal[i] + s * y[xInd[i]-idxBase]

```

Input

handle	handle to the CUSPARSE library context.
nnz	number of elements in vector x .
xVal	<type> vector with nnz non-zero values of vector x .
xInd	integer vector with nnz indices of the non-zero values of vector x .
y	<type> vector in dense format.
c	cosine element of the rotation matrix.
s	sine element of the rotation matrix.
idxBase	<code>CUSPARSE_INDEX_BASE_ZERO</code> or <code>CUSPARSE_INDEX_BASE_ONE</code>

Output

xVal	<type> updated vector in sparse format (that is unchanged if nnz == 0).
y	<type> updated vector in dense format (that is unchanged if nnz == 0).

Status Returned

CUSPARSE_STATUS_SUCCESS	the operation completed successfully.
CUSPARSE_STATUS_NOT_INITIALIZED	the library was not initialized.
CUSPARSE_STATUS_INVALID_VALUE	the idxBase is neither CUSPARSE_INDEX_BASE_ZERO nor CUSPARSE_INDEX_BASE_ONE .
CUSPARSE_STATUS_ARCH_MISMATCH	the device does not support double precision.
CUSPARSE_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU.

6.7. `cusparse<t>sctr`

```

cusparseStatus_t
cusparseSsctr(cusparseHandle_t handle, int nnz,
              const float *xVal,
              const int *xInd, float *y,
              cusparseIndexBase_t idxBase)
cusparseStatus_t
cusparseDsctr(cusparseHandle_t handle, int nnz,
              const double *xVal,
              const int *xInd, double *y,
              cusparseIndexBase_t idxBase)
cusparseStatus_t
cusparseCsctr(cusparseHandle_t handle, int nnz,
              const cuComplex *xVal,
              const int *xInd, cuComplex *y,
              cusparseIndexBase_t idxBase)
cusparseStatus_t
cusparseZsctr(cusparseHandle_t handle, int nnz,
              const cuDoubleComplex *xVal,
              const int *xInd, cuDoubleComplex *y,
              cusparseIndexBase_t idxBase)

```

This function scatters the elements of the vector **x** in sparse format into the vector **y** in dense format. It modifies only the elements of **y** whose indices are listed in the array **xInd**.

This function requires no extra storage. It is executed asynchronously with respect to the host and it may return control to the application on the host before the result is ready.

Input

handle	handle to the CUSPARSE library context.
nnz	number of elements in vector x .
xVal	<type> vector with nnz non-zero values of vector x .

xInd	integer vector with nnz indices of the non-zero values of vector x .
y	<type> dense vector (of size $\geq \max(\mathbf{xInd}) - \mathbf{idxBase} + 1$).
idxBase	CUSPARSE_INDEX_BASE_ZERO or CUSPARSE_INDEX_BASE_ONE

Output

y	<type> vector with nnz non-zero values that were scattered from vector x (that is unchanged if nnz == 0).
----------	--

Status Returned

CUSPARSE_STATUS_SUCCESS	the operation completed successfully.
CUSPARSE_STATUS_NOT_INITIALIZED	the library was not initialized.
CUSPARSE_STATUS_INVALID_VALUE	the idxBase is neither CUSPARSE_INDEX_BASE_ZERO nor CUSPARSE_INDEX_BASE_ONE..
CUSPARSE_STATUS_ARCH_MISMATCH	the device does not support double precision.
CUSPARSE_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU.

Chapter 7.

CUSPARSE LEVEL 2 FUNCTION REFERENCE

This chapter describes the sparse linear algebra functions that perform operations between sparse matrices and dense vectors.

In particular, the solution of sparse triangular linear systems is implemented in two phases. First, during the analysis phase, the sparse triangular matrix is analyzed to determine the dependencies between its elements by calling the appropriate **csrsv_analysis()** function. The analysis is specific to the sparsity pattern of the given matrix and to the selected **cusparseOperation_t** type. The information from the analysis phase is stored in the parameter of type **cusparseSolveAnalysisInfo_t** that has been initialized previously with a call to **cusparseCreateSolveAnalysisInfo()**.

Second, during the solve phase, the given sparse triangular linear system is solved using the information stored in the **cusparseSolveAnalysisInfo_t** parameter by calling the appropriate **csrsv_solve()** function. The solve phase may be performed multiple times with different right-hand-sides, while the analysis phase needs to be performed only once. This is especially useful when a sparse triangular linear system must be solved for a set of different right-hand-sides one at a time, while its coefficient matrix remains the same.

Finally, once all the solves have completed, the opaque data structure pointed to by the **cusparseSolveAnalysisInfo_t** parameter can be released by calling **cusparseDestroySolveAnalysisInfo()**. For more information please refer to [3].

7.1. cusparse<t>bsrmv

```
cusparseStatus_t
cusparseSbsrmv(cusparseHandle_t handle, cusparseDirection_t dir,
               cusparseOperation_t trans, int mb, int nb, int nnzb,
               const float *alpha, const cusparseMatDescr_t descr,
               const float *bsrVal, const int *bsrRowPtr, const int *bsrColInd,
               int blockDim, const float *x,
               const float *beta, float *y)
cusparseStatus_t
cusparseDbsrmv(cusparseHandle_t handle, cusparseDirection_t dir,
               cusparseOperation_t trans, int mb, int nb, int nnzb,
               const double *alpha, const cusparseMatDescr_t descr,
               const double *bsrVal, const int *bsrRowPtr, const int *bsrColInd,
```



```

    int blockDim, const double *x,
    const double *beta, double *y)
cusparseStatus_t
cusparseCbsrmv(cusparseHandle_t handle, cusparseDirection_t dir,
cusparseOperation_t trans, int mb, int nb, int nnzb,
    const cuComplex *alpha, const cusparseMatDescr_t descr,
    const cuComplex *bsrVal, const int *bsrRowPtr, const int *bsrColInd,
    int blockDim, const cuComplex *x,
    const cuComplex *beta, cuComplex *y)
cusparseStatus_t
cusparseZbsrmv(cusparseHandle_t handle, cusparseDirection_t dir,
cusparseOperation_t trans, int mb, int nb, int nnzb,
    const cuDoubleComplex *alpha, const cusparseMatDescr_t descr,
    const cuDoubleComplex *bsrVal, const int *bsrRowPtr, const int
    *bsrColInd,
    int blockDim, const cuDoubleComplex *x,
    const cuDoubleComplex *beta, cuDoubleComplex *y)

```

This function performs the matrix-vector operation

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

where A is $(mb * \text{blockDim}) \times (nb * \text{blockDim})$ sparse matrix (that is defined in BSR storage format by the three arrays **bsrVal**, **bsrRowPtr**, and **bsrColInd**), x and y are vectors, α and β are scalars, and

$$\text{op}(A) = \begin{cases} A & \text{if trans == CUSPARSE_OPERATION_NON_TRANSPOSE} \\ A^T & \text{if trans == CUSPARSE_OPERATION_TRANSPOSE} \\ A^H & \text{if trans == CUSPARSE_OPERATION_CONJUGATE_TRANSPOSE} \end{cases}$$

Several comments on bsrnv:

1. Only **CUSPARSE_OPERATION_NON_TRANSPOSE** is supported, i.e.

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

2. Only **CUSPARSE_MATRIX_TYPE_GENERAL** is supported.

3. The size of vector x should be $(nb * \text{blockDim})$ at least and the size of vector y should be $(mb * \text{blockDim})$ at least. Otherwise the kernel may return **CUSPARSE_STATUS_EXECUTION_FAILED** because of out-of-array-bound.

Example: suppose the user has a CSR format and wants to try bsrnv, the following code demonstrates csr2csc and csmv on single precision.

```

// Suppose that A is m x n sparse matrix represented by CSR format,
// hx is a host vector of size n, and hy is also a host vector of size m.
// m and n are not multiple of blockDim.
// step 1: transform CSR to BSR with column-major order
int base, nnz;
cusparseDirection_t dirA = CUSPARSE_DIRECTION_COLUMN;
int mb = (m + blockDim-1)/blockDim;
int nb = (n + blockDim-1)/blockDim;
cudaMalloc((void**) &bsrRowPtrC, sizeof(int) * (mb+1));
cusparseXcsr2bsrNnz(handle, dirA, m, n,
    descrA, csrRowPtrA, csrColIndA, blockDim,
    descrC, bsrRowPtrC);
cudaMemcpy(&nnzb, bsrRowPtrC+mb, sizeof(int), cudaMemcpyDeviceToHost);
cudaMemcpy(&base, bsrRowPtrC, sizeof(int), cudaMemcpyDeviceToHost);
nnzb -= base;
cudaMalloc((void**) &bsrColIndC, sizeof(int) * nnzb);
cudaMalloc((void**) &bsrValC, sizeof(float) * (blockDim*blockDim) * nnzb);

```

```

cusparsScsr2bsr(handle, dirA, m, n,
    descrA, csrValA, csrRowPtrA, csrColIndA, blockDim,
    descrC, bsrValC, bsrRowPtrC, bsrColIndC);
// step 2: allocate vector x and vector y large enough for bsrmv
cudaMalloc((void**)&x, sizeof(float)*(nb*blockDim));
cudaMalloc((void**)&y, sizeof(float)*(mb*blockDim));
cudaMemcpy(x, hx, sizeof(float)*n, cudaMemcpyHostToDevice);
cudaMemcpy(y, hy, sizeof(float)*m, cudaMemcpyHostToDevice);
// step 3: perform bsrmv
cusparsSbsrmv(handle, dirA, transA, mb, nb, alpha, descrC, bsrValC, bsrRowPtrC,
    bsrColIndC, blockDim, x, beta, y);

```

Input

handle	handle to the CUSPARSE library context.
dir	storage format of blocks, either <code>CUSPARSE_DIRECTION_ROW</code> or <code>CUSPARSE_DIRECTION_COLUMN</code> .
trans	the operation $op(A)$. Only <code>CUSPARSE_OPERATION_NON_TRANSPOSE</code> is supported.
mb	number of block rows of matrix A .
nb	number of block columns of matrix A .
nnzb	number of non-zero blocks of matrix A .
alpha	<type> scalar used for multiplication.
descr	the descriptor of matrix A . The supported matrix type is <code>CUSPARSE_MATRIX_TYPE_GENERAL</code> . Also, the supported index bases are <code>CUSPARSE_INDEX_BASE_ZERO</code> and <code>CUSPARSE_INDEX_BASE_ONE</code> .
bsrVal	<type> array of <code>nnz (= csrRowPtrA(mb) - csrRowPtrA(0))</code> non-zero blocks of matrix A .
bsrRowPtr	integer array of <code>mb + 1</code> elements that contains the start of every block row and the end of the last block row plus one.
bsrColInd	integer array of <code>nnz (= csrRowPtrA(mb) - csrRowPtrA(0))</code> column indices of the non-zero blocks of matrix A .
blockDim	block dimension of sparse matrix A , larger than zero.
x	<type> vector of <code>nb * blockDim</code> elements.
beta	<type> scalar used for multiplication. If <code>beta</code> is zero, <code>y</code> does not have to be a valid input.
y	<type> vector of <code>mb * blockDim</code> elements.

Output

y	<type> updated vector.
----------	------------------------

Status Returned

<code>CUSPARSE_STATUS_SUCCESS</code>	the operation completed successfully.
<code>CUSPARSE_STATUS_NOT_INITIALIZED</code>	the library was not initialized.
<code>CUSPARSE_STATUS_INVALID_VALUE</code>	invalid parameters were passed (<code>m,n,nnz<0</code> , <code>trans != CUSPARSE_OPERATION_NON_TRANSPOSE</code> , <code>blockDim<1</code> , <code>dir</code> is not row-major or column-major, or <code>IndexBase</code> of <code>descr</code> is not base-0 or base-1).
<code>CUSPARSE_STATUS_ARCH_MISMATCH</code>	the device does not support double precision.
<code>CUSPARSE_STATUS_EXECUTION_FAILED</code>	the function failed to launch on the GPU
<code>CUSPARSE_STATUS_INTERNAL_ERROR</code>	an internal operation failed.
<code>CUSPARSE_STATUS_MATRIX_TYPE_NOT_SUPPORTED</code>	the matrix type is not supported.

7.2. `cusparse<t>bsrxmv`

```

cusparseStatus_t
cusparseSbsrxmv(cusparseHandle_t handle, cusparseDirection_t dir,
    cusparseOperation_t trans, int sizeofMask,
    int mb, int nb, int nnzb,
    const float *alpha, const cusparseMatDescr_t descr,
    const float *bsrVal, const int *bsrMaskPtr,
    const int *bsrRowPtr, const int *bsrEndPtr, const int *bsrColInd,
    int blockDim, const float *x,
    const float *beta, float *y)
cusparseStatus_t
cusparseDbsrxmv(cusparseHandle_t handle, cusparseDirection_t dir,
    cusparseOperation_t trans, int sizeofMask,
    int mb, int nb, int nnzb,
    const double *alpha, const cusparseMatDescr_t descr,
    const double *bsrVal, const int *bsrMaskPtr,
    const int *bsrRowPtr, const int *bsrEndPtr, const int *bsrColInd,
    int blockDim, const double *x,
    const double *beta, double *y)
cusparseStatus_t
cusparseCbsrxmv(cusparseHandle_t handle, cusparseDirection_t dir,
    cusparseOperation_t trans, int sizeofMask,
    int mb, int nb, int nnzb,
    const cuComplex *alpha, const cusparseMatDescr_t descr,
    const cuComplex *bsrVal, const int *bsrMaskPtr,
    const int *bsrRowPtr, const int *bsrEndPtr, const int *bsrColInd,
    int blockDim, const cuComplex *x,
    const cuComplex *beta, cuComplex *y)
cusparseStatus_t
cusparseZbsrxmv(cusparseHandle_t handle, cusparseDirection_t dir,
    cusparseOperation_t trans, int sizeofMask,
    int mb, int nb, int nnzb,
    const cuDoubleComplex *alpha, const cusparseMatDescr_t descr,
    const cuDoubleComplex *bsrVal, const int *bsrMaskPtr,
    const int *bsrRowPtr, const int *bsrEndPtr, const int *bsrColInd,
    int blockDim, const cuDoubleComplex *x,
    const cuDoubleComplex *beta, cuDoubleComplex *y)

```

This function performs a bsr_mmv and a mask operation

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

where A is $(mb * \text{blockDim}) \times (nb * \text{blockDim})$ sparse matrix (that is defined in BSRX storage format by the four arrays **bsrVal**, **bsrRowPtr**, **bsrEndPtr**, and **bsrColInd**), x and y are vectors, α and β are scalars, and

$$\text{op}(A) = \begin{cases} A & \text{if trans} == \text{CUSPARSE_OPERATION_NON_TRANPOSE} \\ A^T & \text{if trans} == \text{CUSPARSE_OPERATION_TRANPOSE} \\ A^H & \text{if trans} == \text{CUSPARSE_OPERATION_CONJUGATE_TRANPOSE} \end{cases}$$

The mask operation is defined by array **bsrMaskPtr** which contains updated row indices of y . If row i is not specified in **bsrMaskPtr**, then **bsrxmv** does not touch row block i of A and $y[i]$.

For example, consider the 2×3 block matrix A :

$$A = \begin{bmatrix} A_{11} & A_{12} & O \\ A_{21} & A_{22} & A_{23} \end{bmatrix}$$

and its one-based BSR format (three vector form) is

$$\begin{aligned} \text{bsrVal} &= [A_{11} \ A_{12} \ A_{21} \ A_{22} \ A_{23}] \\ \text{bsrRowPtr} &= [1 \ 3 \ 6] \\ \text{bsrColInd} &= [1 \ 2 \ 1 \ 2 \ 3] \end{aligned}$$

Suppose we want to do the following **bsrmv** operation on a matrix \bar{A} which is slightly different from A .

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} := \alpha * (\tilde{A} = \begin{bmatrix} O & O & O \\ O & A_{22} & O \end{bmatrix}) * \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} y_1 \\ \beta * y_2 \end{bmatrix}$$

We don't need to create another BSR format for the new matrix \bar{A} , all that we should do is to keep **bsrVal** and **bsrColInd** unchanged, but modify **bsrRowPtr** and add additional array **bsrEndPtr** which points to last nonzero elements per row of \bar{A} plus 1.

For example, the following **bsrRowPtr** and **bsrEndPtr** can represent matrix \bar{A} :

$$\begin{aligned} \text{bsrRowPtr} &= [1 \ 4] \\ \text{bsrEndPtr} &= [1 \ 5] \end{aligned}$$

Further we can use mask operator (specified by array **bsrMaskPtr**) to update particular row indices of y only because y_1 is never changed. In this case, **bsrMaskPtr** = [2]

The mask operator is equivalent to the following operation (? stands for don't care)

$$\begin{bmatrix} ? \\ y_2 \end{bmatrix} := \alpha * \begin{bmatrix} ? & ? & ? \\ O & A_{22} & O \end{bmatrix} * \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \beta * \begin{bmatrix} ? \\ y_2 \end{bmatrix}$$

In other words, **bsrRowPtr**[0] and **bsrEndPtr**[0] are don't care.

$$\begin{aligned} \text{bsrRowPtr} &= [? \ 4] \\ \text{bsrEndPtr} &= [? \ 5] \end{aligned}$$

Several comments on `bsrxmv`:

Only `CUSPARSE_OPERATION_NON_TRANSPOSE` and `CUSPARSE_MATRIX_TYPE_GENERAL` are supported.

`bsrMaskPtr`, `bsrRowPtr`, `bsrEndPtr` and `bsrColInd` are consistent with base index, either one-based or zero-based. Above example is one-based.

Input

<code>handle</code>	handle to the CUSPARSE library context.
<code>dir</code>	storage format of blocks, either <code>CUSPARSE_DIRECTION_ROW</code> or <code>CUSPARSE_DIRECTION_COLUMN</code> .
<code>trans</code>	the operation $op(A)$. Only <code>CUSPARSE_OPERATION_NON_TRANSPOSE</code> is supported.
<code>sizeOfMask</code>	number of updated rows of y .
<code>mb</code>	number of block rows of matrix A .
<code>nb</code>	number of block columns of matrix A .
<code>nnzb</code>	number of nonz-zero blocks of matrix A .
<code>alpha</code>	<type> scalar used for multiplication.
<code>descr</code>	the descriptor of matrix A . The supported matrix type is <code>CUSPARSE_MATRIX_TYPE_GENERAL</code> . Also, the supported index bases are <code>CUSPARSE_INDEX_BASE_ZERO</code> and <code>CUSPARSE_INDEX_BASE_ONE</code> .
<code>bsrVal</code>	<type> array of <code>nnz</code> non-zero blocks of matrix A .
<code>bsrRowPtr</code>	integer array of <code>mb</code> elements that contains the start of every block row and the end of the last block row plus one.
<code>bsrEndPtr</code>	integer array of <code>mb</code> elements that contains the end of the every block row plus one.
<code>bsrColInd</code>	integer array of <code>nnzb</code> column indices of the non-zero blocks of matrix A .
<code>blockDim</code>	block dimension of sparse matrix A , larger than zero.
<code>x</code>	<type> vector of $nb * blockDim$ elements.
<code>beta</code>	<type> scalar used for multiplication. If <code>beta</code> is zero, y does not have to be a valid input.
<code>y</code>	<type> vector of $mb * blockDim$ elements.

Status Returned

<code>CUSPARSE_STATUS_SUCCESS</code>	the operation completed successfully.
<code>CUSPARSE_STATUS_NOT_INITIALIZED</code>	the library was not initialized.
<code>CUSPARSE_STATUS_INVALID_VALUE</code>	invalid parameters were passed ($m, n, nnz < 0$, <code>trans !=</code>

	CUSPARSE_OPERATION_NON_TRANSPOSE, <i>blockDim</i> < 1, <i>dir</i> is not row-major or column-major, or <i>IndexBase</i> of <i>descr</i> is not base-0 or base-1).
CUSPARSE_STATUS_ARCH_MISMATCH	the device does not support double precision.
CUSPARSE_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU
CUSPARSE_STATUS_INTERNAL_ERROR	an internal operation failed.
CUSPARSE_STATUS_MATRIX_TYPE_NOT_SUPPORTED	the matrix type is not supported.

7.3. `cusparse<t>csrcmv`

```

cusparseStatus_t
cusparseScsrcmv(cusparseHandle_t handle, cusparseOperation_t transA,
               int m, int n, int nnz, const float *alpha,
               const cusparseMatDescr_t descrA,
               const float *csrValA,
               const int *csrRowPtrA, const int *csrColIndA,
               const float *x, const float *y, const float *beta,
               float *y)
cusparseStatus_t
cusparseDcsrcmv(cusparseHandle_t handle, cusparseOperation_t transA,
               int m, int n, int nnz, const double *alpha,
               const cusparseMatDescr_t descrA,
               const double *csrValA,
               const int *csrRowPtrA, const int *csrColIndA,
               const double *x, const double *y, const double *beta,
               double *y)
cusparseStatus_t
cusparseCcsrcmv(cusparseHandle_t handle, cusparseOperation_t transA,
               int m, int n, int nnz, const cuComplex *alpha,
               const cusparseMatDescr_t descrA,
               const cuComplex *csrValA,
               const int *csrRowPtrA, const int *csrColIndA,
               const cuComplex *x, const cuComplex *y, const cuComplex *beta,
               cuComplex *y)
cusparseStatus_t
cusparseZcsrcmv(cusparseHandle_t handle, cusparseOperation_t transA,
               int m, int n, int nnz, const cuDoubleComplex *alpha,
               const cusparseMatDescr_t descrA,
               const cuDoubleComplex *csrValA,
               const int *csrRowPtrA, const int *csrColIndA,
               const cuDoubleComplex *x, const cuDoubleComplex *y, const cuDoubleComplex *beta,
               cuDoubleComplex *y)

```

This function performs the matrix-vector operation

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

where A is $m \times n$ sparse matrix (that is defined in CSR storage format by the three arrays **csrValA**, **csrRowPtrA**, and **csrColIndA**), x and y are vectors, α and β are scalars, and

$$\text{op}(A) = \begin{cases} A & \text{if trans == CUSPARSE_OPERATION_NON_TRANSPOSE} \\ A^T & \text{if trans == CUSPARSE_OPERATION_TRANSPOSE} \\ A^H & \text{if trans == CUSPARSE_OPERATION_CONJUGATE_TRANSPOSE} \end{cases}$$

When using the (conjugate) transpose of a general matrix or a Hermitian/symmetric matrix, this routine may produce slightly different results during different runs of this function with the same input parameters. For these matrix types it uses atomic operations to compute the final result, consequently many threads may be adding floating point numbers to the same memory location without any specific ordering, which may produce slightly different results for each run.

If exactly the same output is required for any input when multiplying by the transpose of a general matrix, the following procedure can be used:

1. Convert the matrix from CSR to CSC format using one of the **csr2csc()** functions. Notice that by interchanging the rows and columns of the result you are implicitly transposing the matrix.
2. Call the **csrmmv()** function with the **cusparseOperation_t** parameter set to **CUSPARSE_OPERATION_NON_TRANSPOSE** and with the interchanged rows and columns of the matrix stored in CSC format. This (implicitly) multiplies the vector by the transpose of the matrix in the original CSR format.

This function requires no extra storage for the general matrices when operation **CUSPARSE_OPERATION_NON_TRANSPOSE** is selected. It requires some extra storage for Hermitian/symmetric matrices and for the general matrices when operation different than **CUSPARSE_OPERATION_NON_TRANSPOSE** is selected. It is executed asynchronously with respect to the host and it may return control to the application on the host before the result is ready.

Input

handle	handle to the CUSPARSE library context.
trans	the operation $op(A)$
m	number of rows of matrix A .
n	number of columns of matrix A .
nnz	number of nonz-zero elements of matrix A .
alpha	<type> scalar used for multiplication.
descrA	the descriptor of matrix A . The supported matrix type is CUSPARSE_MATRIX_TYPE_GENERAL , CUSPARSE_MATRIX_TYPE_SYMMETRIC , and CUSPARSE_MATRIX_TYPE_HERMITIAN . Also, the supported index bases are CUSPARSE_INDEX_BASE_ZERO and CUSPARSE_INDEX_BASE_ONE .
csrValA	<type> array of nnz ($= \text{csrRowPtrA}(m) - \text{csrRowPtrA}(0)$) non-zero elements of matrix A .
csrRowPtrA	integer array of $m+1$ elements that contains the start of every row and the end of the last row plus one.
csrColIndA	integer array of nnz ($= \text{csrRowPtrA}(m) - \text{csrRowPtrA}(0)$) column indices of the non-zero elements of matrix A .

x	<type> vector of n elements if $\text{op}(A) = A$, and m elements if $\text{op}(A) = A^T$ or $\text{op}(A) = A^H$
beta	<type> scalar used for multiplication. If beta is zero, y does not have to be a valid input.
y	<type> vector of m elements if $\text{op}(A) = A$, and n elements if $\text{op}(A) = A^T$ or $\text{op}(A) = A^H$

Output

y	<type> updated vector.
----------	------------------------

Status Returned

CUSPARSE_STATUS_SUCCESS	the operation completed successfully.
CUSPARSE_STATUS_NOT_INITIALIZED	the library was not initialized.
CUSPARSE_STATUS_ALLOC_FAILED	the resources could not be allocated.
CUSPARSE_STATUS_INVALID_VALUE	invalid parameters were passed (m , n , nnz < 0).
CUSPARSE_STATUS_ARCH_MISMATCH	the device does not support double precision. (compute capability (c.c.) >= 1.3), symmetric/ Hermitian matrix (c.c. >= 1.2) or transpose operation (c.c. >= 1.1).
CUSPARSE_STATUS_INTERNAL_ERROR	an internal operation failed.
CUSPARSE_STATUS_MATRIX_TYPE_NOT_SUPPORTED	the matrix type is not supported.

7.4. cusparse<t>csrsv_analysis

```

cusparseStatus_t
cusparseScsrsv_analysis(cusparseHandle_t handle,
                        cusparseOperation_t transA,
                        int m, int nnz, const cusparseMatDescr_t descrA,
                        const float *csrValA,
                        const int *csrRowPtrA, const int *csrColIndA,
                        cusparseSolveAnalysisInfo_t info)

cusparseStatus_t
cusparseDcsrsv_analysis(cusparseHandle_t handle,
                        cusparseOperation_t transA,
                        int m, int nnz, const cusparseMatDescr_t descrA,
                        const double *csrValA,
                        const int *csrRowPtrA, const int *csrColIndA,
                        cusparseSolveAnalysisInfo_t info)

cusparseStatus_t
cusparseCcsrsv_analysis(cusparseHandle_t handle,
                        cusparseOperation_t transA,
                        int m, int nnz, const cusparseMatDescr_t descrA,
                        const cuComplex *csrValA,
                        const int *csrRowPtrA, const int *csrColIndA,
                        cusparseSolveAnalysisInfo_t info)

cusparseStatus_t
cusparseZcsrsv_analysis(cusparseHandle_t handle,
                        cusparseOperation_t transA,
```



```
int m, int nnz, const cusparseMatDescr_t descrA,
const cuDoubleComplex *csrValA,
const int *csrRowPtrA, const int *csrColIndA,
cusparseSolveAnalysisInfo_t info)
```

This function performs the analysis phase of the solution of a sparse triangular linear system

$$\text{op}(A) * y = \alpha * x$$

where A is $m \times m$ sparse matrix (that is defined in CSR storage format by the three arrays **csrValA**, **csrRowPtrA**, and **csrColIndA**), x and y are the right-hand-side and the solution vectors, α is a scalar, and

$$\text{op}(A) = \begin{cases} A & \text{if trans == CUSPARSE_OPERATION_NON_TRANSPOSE} \\ A^T & \text{if trans == CUSPARSE_OPERATION_TRANSPOSE} \\ A^H & \text{if trans == CUSPARSE_OPERATION_CONJUGATE_TRANSPOSE} \end{cases}$$

It is expected that this function will be executed only once for a given matrix and a particular operation type.

This function requires significant amount of extra storage that is proportional to the matrix size. It is executed asynchronously with respect to the host and it may return control to the application on the host before the result is ready.

Input

handle	handle to the CUSPARSE library context.
trans	the operation $\text{op}(A)$
m	number of rows of matrix A .
nnz	number of nonz-zero elements of matrix A .
descrA	the descriptor of matrix A . The supported matrix types are CUSPARSE_MATRIX_TYPE_TRIANGULAR and CUSPARSE_MATRIX_TYPE_GENERAL , while the supported diagonal types are CUSPARSE_DIAG_TYPE_UNIT and CUSPARSE_DIAG_TYPE_NON_UNIT .
csrValA	<type> array of $\text{nnz} (= \text{csrRowPtrA}(m) - \text{csrRowPtrA}(0))$ non-zero elements of matrix A .
csrRowPtrA	integer array of $m + 1$ elements that contains the start of every row and the end of the last row plus one.
csrColIndA	integer array of $\text{nnz} (= \text{csrRowPtrA}(m) - \text{csrRowPtrA}(0))$ column indices of the non-zero elements of matrix A .
info	structure initialized using cusparseCreateSolveAnalysisInfo .

Output

info	structure filled with information collected during the analysis phase (that should be passed to the solve phase unchanged).
-------------	---

Status Returned

CUSPARSE_STATUS_SUCCESS	the operation completed successfully.
CUSPARSE_STATUS_NOT_INITIALIZED	the library was not initialized.
CUSPARSE_STATUS_ALLOC_FAILED	the resources could not be allocated.
CUSPARSE_STATUS_INVALID_VALUE	invalid parameters were passed (<code>m, nnz < 0</code>).
CUSPARSE_STATUS_ARCH_MISMATCH	the device does not support double precision.
CUSPARSE_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU
CUSPARSE_STATUS_INTERNAL_ERROR	an internal operation failed.
CUSPARSE_STATUS_MATRIX_TYPE_NOT_SUPPORTED	the matrix type is not supported.

7.5. `cusparse<t>csrsv_solve`

```

cusparseStatus_t
cusparseScsrsv_solve(cusparseHandle_t handle,
                    cusparseOperation_t transA,
                    int m, const float *alpha,
                    const cusparseMatDescr_t descrA,
                    const float *csrValA,
                    const int *csrRowPtrA, const int *csrColIndA,
                    cusparseSolveAnalysisInfo_t info,
                    const float *x, float *y)

cusparseStatus_t
cusparseDcsrsv_solve(cusparseHandle_t handle,
                    cusparseOperation_t transA,
                    int m, const double *alpha,
                    const cusparseMatDescr_t descrA,
                    const double *csrValA,
                    const int *csrRowPtrA, const int *csrColIndA,
                    cusparseSolveAnalysisInfo_t info,
                    const double *x, double *y)

cusparseStatus_t
cusparseCcsrsv_solve(cusparseHandle_t handle,
                    cusparseOperation_t transA,
                    int m, const cuComplex *alpha,
                    const cusparseMatDescr_t descrA,
                    const cuComplex *csrValA,
                    const int *csrRowPtrA, const int *csrColIndA,
                    cusparseSolveAnalysisInfo_t info,
                    const cuComplex *x, cuComplex *y)

cusparseStatus_t
cusparseZcsrsv_solve(cusparseHandle_t handle,
                    cusparseOperation_t transA,
                    int m, const cuDoubleComplex *alpha,
                    const cusparseMatDescr_t descrA,
                    const cuDoubleComplex *csrValA,
                    const int *csrRowPtrA, const int *csrColIndA,
                    cusparseSolveAnalysisInfo_t info,
                    const cuDoubleComplex *x, cuDoubleComplex *y)

```

This function performs the solve phase of the solution of a sparse triangular linear system

$$\text{op}(A) * y = \alpha * x$$

where A is $m \times m$ sparse matrix (that is defined in CSR storage format by the three arrays **csrValA**, **csrRowPtrA**, and **csrColIndA**), x and y are the right-hand-side and the solution vectors, α is a scalar, and

$$\text{op}(A) = \begin{cases} A & \text{if trans} == \text{CUSPARSE_OPERATION_NON_TRANPOSE} \\ A^T & \text{if trans} == \text{CUSPARSE_OPERATION_TRANPOSE} \\ A^H & \text{if trans} == \text{CUSPARSE_OPERATION_CONJUGATE_TRANPOSE} \end{cases}$$

This function may be executed multiple times for a given matrix and a particular operation type.

This function requires no extra storage. It is executed asynchronously with respect to the host and it may return control to the application on the host before the result is ready.

Input

handle	handle to the CUSPARSE library context.
trans	the operation $\text{op}(A)$
m	number of rows and columns of matrix A .
alpha	<type> scalar used for multiplication.
descrA	the descriptor of matrix A . The supported matrix types are CUSPARSE_MATRIX_TYPE_TRIANGULAR and CUSPARSE_MATRIX_TYPE_GENERAL , while the supported diagonal types are CUSPARSE_DIAG_TYPE_UNIT and CUSPARSE_DIAG_TYPE_NON_UNIT .
csrValA	<type> array of $\text{nnz} (= \text{csrRowPtrA}(m) - \text{csrRowPtrA}(0))$ non-zero elements of matrix A .
csrRowPtrA	integer array of $m + 1$ elements that contains the start of every row and the end of the last row plus one.
csrColIndA	integer array of $\text{nnz} (= \text{csrRowPtrA}(m) - \text{csrRowPtrA}(0))$ column indices of the non-zero elements of matrix A .
info	structure with information collected during the analysis phase (that should have been passed to the solve phase unchanged).
x	<type> right-hand-side vector of size m .

Output

y	<type> solution vector of size m .
----------	--------------------------------------

Status Returned

CUSPARSE_STATUS_SUCCESS	the operation completed successfully.
CUSPARSE_STATUS_NOT_INITIALIZED	the library was not initialized.
CUSPARSE_STATUS_INVALID_VALUE	invalid parameters were passed ($m < 0$).

<code>CUSPARSE_STATUS_ARCH_MISMATCH</code>	the device does not support double precision.
<code>CUSPARSE_STATUS_MAPPING_ERROR</code>	the texture binding failed.
<code>CUSPARSE_STATUS_EXECUTION_FAILED</code>	the function failed to launch on the GPU
<code>CUSPARSE_STATUS_INTERNAL_ERROR</code>	an internal operation failed.
<code>CUSPARSE_STATUS_MATRIX_TYPE_NOT_SUPPORTED</code>	the matrix type is not supported.

7.6. `cusparse<t>hybm`

```

cusparseStatus_t
cusparseShybm(cusparseHandle_t handle, cusparseOperation_t transA,
              const float *alpha,
              const cusparseMatDescr_t descrA,
              const cusparseHybMat_t hybA, const float *x,
              const float *beta, float *y)

cusparseStatus_t
cusparseDhybm(cusparseHandle_t handle, cusparseOperation_t transA,
              const double *alpha,
              const cusparseMatDescr_t descrA,
              const cusparseHybMat_t hybA, const double *x,
              const double *beta, double *y)

cusparseStatus_t
cusparseChybm(cusparseHandle_t handle, cusparseOperation_t transA,
              const cuComplex *alpha,
              const cusparseMatDescr_t descrA,
              const cusparseHybMat_t hybA, const cuComplex *x,
              const cuComplex *beta, cuComplex *y)

cusparseStatus_t
cusparseZhybm(cusparseHandle_t handle, cusparseOperation_t transA,
              const cuDoubleComplex *alpha,
              const cusparseMatDescr_t descrA,
              const cusparseHybMat_t hybA, const cuDoubleComplex *x,
              const cuDoubleComplex *beta, cuDoubleComplex *y)

```

This function performs the matrix-vector operation

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

where A is an $m \times n$ sparse matrix (that is defined in the HYB storage format by an opaque data structure **hybA**), x and y are vectors, α and β are scalars, and

$\text{op}(A) = \{ A \quad \text{if } \text{transA} == \text{CUSPARSE_OPERATION_NON_TRANPOSE} \}$

Notice that currently only $\text{op}(A) = A$ is supported.

This function requires no extra storage. It is executed asynchronously with respect to the host and it may return control to the application on the host before the result is ready.

Input

handle	handle to the CUSPARSE library context.
transA	the operation $\text{op}(A)$ (currently only $\text{op}(A) = A$ is supported).
m	number of rows of matrix A .

n	number of columns of matrix <i>A</i> .
alpha	<type> scalar used for multiplication.
descrA	the descriptor of matrix <i>A</i> . The supported matrix type is <code>CUSPARSE_MATRIX_TYPE_GENERAL</code> .
hybA	the matrix <i>A</i> in HYB storage format.
x	<type> vector of <i>n</i> elements.
beta	<type> scalar used for multiplication. If beta is zero, y does not have to be a valid input.
y	<type> vector of <i>m</i> elements.

Output

y	<type> updated vector.
----------	------------------------

Status Returned

<code>CUSPARSE_STATUS_SUCCESS</code>	the operation completed successfully.
<code>CUSPARSE_STATUS_NOT_INITIALIZED</code>	the library was not initialized.
<code>CUSPARSE_STATUS_ALLOC_FAILED</code>	the resources could not be allocated.
<code>CUSPARSE_STATUS_INVALID_VALUE</code>	the internally stored hyb format parameters are invalid.
<code>CUSPARSE_STATUS_ARCH_MISMATCH</code>	the device does not support double precision.
<code>CUSPARSE_STATUS_EXECUTION_FAILED</code>	the function failed to launch on the GPU
<code>CUSPARSE_STATUS_INTERNAL_ERROR</code>	an internal operation failed.
<code>CUSPARSE_STATUS_MATRIX_TYPE_NOT_SUPPORTED</code>	the matrix type is not supported.

7.7. `cusparse<t>hybsv_analysis`

```

cusparseStatus_t
cusparseShybsv_analysis(cusparseHandle_t handle,
                        cusparseOperation_t transA,
                        const cusparseMatDescr_t descrA,
                        cusparseHybMat_t hybA,
                        cusparseSolveAnalysisInfo_t info)

cusparseStatus_t
cusparseDhybsv_analysis(cusparseHandle_t handle,
                        cusparseOperation_t transA,
                        const cusparseMatDescr_t descrA,
                        cusparseHybMat_t hybA,
                        cusparseSolveAnalysisInfo_t info)

cusparseStatus_t
cusparseChybsv_analysis(cusparseHandle_t handle,
                        cusparseOperation_t transA,
                        const cusparseMatDescr_t descrA,
                        cusparseHybMat_t hybA,
                        cusparseSolveAnalysisInfo_t info)

cusparseStatus_t
cusparseZhybsv_analysis(cusparseHandle_t handle,

```

```

cusparseOperation_t transA,
const cusparseMatDescr_t descrA,
cusparseHybMat_t hybA,
cusparseSolveAnalysisInfo_t info)

```

This function performs the analysis phase of the solution of a sparse triangular linear system

$$\text{op}(A) * y = \alpha * x$$

where A is $m \times m$ sparse matrix (that is defined in HYB storage format by an opaque data structure **hybA**), x and y are the right-hand-side and the solution vectors, α is a scalar, and

$\text{op}(A) = \{ A \quad \text{if } \text{transA} == \text{CUSPARSE_OPERATION_NON_TRANPOSE}$

Notice that currently only $\text{op}(A) = A$ is supported.

It is expected that this function will be executed only once for a given matrix and a particular operation type.

This function requires significant amount of extra storage that is proportional to the matrix size. It is executed asynchronously with respect to the host and it may return control to the application on the host before the result is ready.

Input

handle	handle to the CUSPARSE library context.
transA	the operation $\text{op}(A)$ (currently only $\text{op}(A) = A$ is supported).
descrA	the descriptor of matrix A . The supported matrix type is CUSPARSE_MATRIX_TYPE_TRIANGULAR and diagonal type USPARSE_DIAG_TYPE_NON_UNIT .
hybA	the matrix A in HYB storage format.
info	structure initialized using cusparseCreateSolveAnalysisInfo .

Output

info	structure filled with information collected during the analysis phase (that should be passed to the solve phase unchanged).
-------------	---

Status Returned

CUSPARSE_STATUS_SUCCESS	the operation completed successfully.
CUSPARSE_STATUS_NOT_INITIALIZED	the library was not initialized.
CUSPARSE_STATUS_ALLOC_FAILED	the resources could not be allocated.
CUSPARSE_STATUS_INVALID_VALUE	the internally stored hyb format parameters are invalid.
CUSPARSE_STATUS_ARCH_MISMATCH	the device does not support double precision.
CUSPARSE_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU

CUSPARSE_STATUS_INTERNAL_ERROR	an internal operation failed.
CUSPARSE_STATUS_MATRIX_TYPE_NOT_SUPPORTED	the matrix type is not supported.

7.8. cusparse<t>hybsv_solve

```

cusparseStatus_t
cusparseShybsv_solve(cusparseHandle_t handle,
                    cusparseOperation_t transA,
                    const float          *alpha,
                    const cusparseMatDescr_t descrA,
                    cusparseHybMat_t hybA,
                    cusparseSolveAnalysisInfo_t info,
                    const float          *x, float          *y)

cusparseStatus_t
cusparseDhybsv_solve(cusparseHandle_t handle,
                    cusparseOperation_t transA,
                    const double         *alpha,
                    const cusparseMatDescr_t descrA,
                    cusparseHybMat_t hybA,
                    cusparseSolveAnalysisInfo_t info,
                    const double         *x, double         *y)

cusparseStatus_t
cusparseChybsv_solve(cusparseHandle_t handle,
                    cusparseOperation_t transA,
                    const cuComplex      *alpha,
                    const cusparseMatDescr_t descrA,
                    cusparseHybMat_t hybA,
                    cusparseSolveAnalysisInfo_t info,
                    const cuComplex      *x, cuComplex      *y)

cusparseStatus_t
cusparseZhybsv_solve(cusparseHandle_t handle,
                    cusparseOperation_t transA,
                    const cuDoubleComplex *alpha,
                    const cusparseMatDescr_t descrA,
                    cusparseHybMat_t hybA,
                    cusparseSolveAnalysisInfo_t info,
                    const cuDoubleComplex *x, cuDoubleComplex *y)

```

This function performs the solve phase of the solution of a sparse triangular linear system

$$\text{op}(A) * y = \alpha * x$$

where A is $m \times m$ sparse matrix (that is defined in HYB storage format by an opaque data structure **hybA**), x and y are the right-hand-side and the solution vectors, α is a scalar, and

$\text{op}(A) = \{ A \quad \text{if } \text{transA} == \text{CUSPARSE_OPERATION_NON_TRANPOSE}$

Notice that currently only $\text{op}(A) = A$ is supported.

This function may be executed multiple times for a given matrix and a particular operation type.

This function requires no extra storage. It is executed asynchronously with respect to the host and it may return control to the application on the host before the result is ready.

Input

handle	handle to the CUSPARSE library context.
transA	the operation $\text{op}(A)$ (currently only $\text{op}(A) = A$ is supported).
alpha	<type> scalar used for multiplication.
descrA	the descriptor of matrix A . The supported matrix type is <code>CUSPARSE_MATRIX_TYPE_TRIANGULAR</code> and diagonal type <code>CUSPARSE_DIAG_TYPE_NON_UNIT</code> .
hybA	the matrix A in HYB storage format.
info	structure with information collected during the analysis phase (that should be passed to the solve phase unchanged).
x	<type> right-hand-side vector of size m .

Output

y	<type> solution vector of size m .
----------	--------------------------------------

Status Returned

<code>CUSPARSE_STATUS_SUCCESS</code>	the operation completed successfully.
<code>CUSPARSE_STATUS_NOT_INITIALIZED</code>	the library was not initialized.
<code>CUSPARSE_STATUS_INVALID_VALUE</code>	the internally stored hyb format parameters are invalid.
<code>CUSPARSE_STATUS_ARCH_MISMATCH</code>	the device does not support double precision.
<code>CUSPARSE_STATUS_MAPPING_ERROR</code>	the texture binding failed.
<code>CUSPARSE_STATUS_EXECUTION_FAILED</code>	the function failed to launch on the GPU
<code>CUSPARSE_STATUS_INTERNAL_ERROR</code>	an internal operation failed.
<code>CUSPARSE_STATUS_MATRIX_TYPE_NOT_SUPPORTED</code>	the matrix type is not supported.

Chapter 8.

CUSPARSE LEVEL 3 FUNCTION REFERENCE

This chapter describes sparse linear algebra functions that perform operations between sparse and (usually tall) dense matrices.

In particular, the solution of sparse triangular linear systems with multiple right-hand-sides is implemented in two phases. First, during the analysis phase, the sparse triangular matrix is analyzed to determine the dependencies between its elements by calling the appropriate `csrsm_analysis()` function. The analysis is specific to the sparsity pattern of the given matrix and to the selected `cusparseOperation_t` type. The information from the analysis phase is stored in the parameter of type `cusparseSolveAnalysisInfo_t` that has been initialized previously with a call to `cusparseCreateSolveAnalysisInfo()`.

Second, during the solve phase, the given sparse triangular linear system is solved using the information stored in the `cusparseSolveAnalysisInfo_t` parameter by calling the appropriate `csrsm_solve()` function. The solve phase may be performed multiple times with different multiple right-hand-sides, while the analysis phase needs to be performed only once. This is especially useful when a sparse triangular linear system must be solved for different sets of multiple right-hand-sides one at a time, while its coefficient matrix remains the same.

Finally, once all the solves have completed, the opaque data structure pointed to by the `cusparseSolveAnalysisInfo_t` parameter can be released by calling `cusparseDestroySolveAnalysisInfo()`. For more information please refer to [3].

8.1. `cusparse<t>csrmm`

```
cusparseStatus_t
cusparseScsrmm(cusparseHandle_t handle, cusparseOperation_t transA,
               int m, int n, int k, int nnz,
               const float *alpha,
               const cusparseMatDescr_t descrA,
               const float *csrValA,
               const int *csrRowPtrA, const int *csrColIndA,
               const float *B, int ldb,
               const float *beta, float *C, int ldc)
cusparseStatus_t
cusparseDcsrmm(cusparseHandle_t handle, cusparseOperation_t transA,
```

```

        int m, int n, int k, int nnz,
        const double *alpha,
        const cusparseMatDescr_t descrA,
        const double *csrValA,
        const int *csrRowPtrA, const int *csrColIndA,
        const double *B, int ldb,
        const double *beta, double *C, int ldc)
cusparseStatus_t
cusparseCcsrmm(cusparseHandle_t handle, cusparseOperation_t transA,
        int m, int n, int k, int nnz,
        const cuComplex *alpha,
        const cusparseMatDescr_t descrA,
        const cuComplex *csrValA,
        const int *csrRowPtrA, const int *csrColIndA,
        const cuComplex *B, int ldb,
        const cuComplex *beta, cuComplex *C, int ldc)
cusparseStatus_t
cusparseZcsrmm(cusparseHandle_t handle, cusparseOperation_t transA,
        int m, int n, int k, int nnz,
        const cuDoubleComplex *alpha,
        const cusparseMatDescr_t descrA,
        const cuDoubleComplex *csrValA,
        const int *csrRowPtrA, const int *csrColIndA,
        const cuDoubleComplex *B, int ldb,
        const cuDoubleComplex *beta, cuDoubleComplex *C, int ldc)

```

This function performs one of the following matrix-matrix operation

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

where A is $m \times k$ sparse matrix (that is defined in CSR storage format by the three arrays **csrValA**, **csrRowPtrA**, and **csrColIndA**), B and C are dense matrices, α and β are scalars, and

$$\text{op}(A) = \begin{cases} A & \text{if trans == CUSPARSE_OPERATION_NON_TRANSPOSE} \\ A^T & \text{if trans == CUSPARSE_OPERATION_TRANSPOSE} \\ A^H & \text{if trans == CUSPARSE_OPERATION_CONJUGATE_TRANSPOSE} \end{cases}$$

When using the (conjugate) transpose of a general matrix or a Hermitian/symmetric matrix, this routine may produce slightly different results during different runs of this function with the same input parameters. For these matrix types it uses atomic operations to compute the final result, consequently many threads may be adding floating point numbers to the same memory location without any specific ordering, which may produce slightly different results for each run.

If exactly the same output is required for any input when multiplying by the transpose of a general matrix, the following procedure can be used:

1. Convert the matrix from CSR to CSC format using one of the **csr2csc()** functions. Notice that by interchanging the rows and columns of the result you are implicitly transposing the matrix.
2. Call the **csrmm()** function with the **cusparseOperation_t** parameter set to **CUSPARSE_OPERATION_NON_TRANSPOSE** and with the interchanged rows and columns of the matrix stored in CSC format. This (implicitly) multiplies the vector by the transpose of the matrix in the original CSR format.

This function requires no extra storage for the general matrices when operation **CUSPARSE_OPERATION_NON_TRANSPOSE** is selected. It requires some extra storage for Hermitian/symmetric matrices and for the general matrices when operation different than **CUSPARSE_OPERATION_NON_TRANSPOSE** is selected. It is executed asynchronously with respect to the host and it may return control to the application on the host before the result is ready.

Input

handle	handle to the CUSPARSE library context.
transA	the operation $\text{op}(A)$
m	number of rows of sparse matrix A .
n	number of columns of dense matrix B and C .
k	number of columns of sparse matrix A .
nnz	number of nonz-zero elements of sparse matrix A .
alpha	<type> scalar used for multiplication.
descrA	the descriptor of matrix A . The supported matrix type is CUSPARSE_MATRIX_TYPE_GENERAL , CUSPARSE_MATRIX_TYPE_SYMMETRIC , and CUSPARSE_MATRIX_TYPE_HERMITIAN . Also, the supported index bases are CUSPARSE_INDEX_BASE_ZERO and CUSPARSE_INDEX_BASE_ONE .
csrValA	<type> array of nnz ($= \text{csrRowPtrA}(m) - \text{csrRowPtrA}(0)$) non-zero elements of matrix A .
csrRowPtrA	integer array of $m + 1$ elements that contains the start of every row and the end of the last row plus one.
csrColIndA	integer array of nnz ($= \text{csrRowPtrA}(m) - \text{csrRowPtrA}(0)$) column indices of the non-zero elements of matrix A .
B	array of dimensions (ldb , n).
ldb	leading dimension of B . It must be at least $\max(1, k)$ if $\text{op}(A) = A$ and at least $\max(1, m)$ otherwise.
beta	<type> scalar used for multiplication. If beta is zero, c does not have to be a valid input.
C	array of dimensions (ldc , n).
ldc	leading dimension of c . It must be at least $\max(1, k)$ if $\text{op}(A) = A$ and at least $\max(1, m)$ otherwise.

Output

c	<type> updated array of dimensions (ldc , n).
----------	---

Status Returned

<code>CUSPARSE_STATUS_SUCCESS</code>	the operation completed successfully.
<code>CUSPARSE_STATUS_NOT_INITIALIZED</code>	the library was not initialized.
<code>CUSPARSE_STATUS_ALLOC_FAILED</code>	the resources could not be allocated.
<code>CUSPARSE_STATUS_INVALID_VALUE</code>	invalid parameters were passed (<code>m,n,k,nnz<0</code> or <code>ldb</code> and <code>ldc</code> are incorrect).
<code>CUSPARSE_STATUS_ARCH_MISMATCH</code>	the device does not support double precision.
<code>CUSPARSE_STATUS_EXECUTION_FAILED</code>	the function failed to launch on the GPU
<code>CUSPARSE_STATUS_INTERNAL_ERROR</code>	an internal operation failed.
<code>CUSPARSE_STATUS_MATRIX_TYPE_NOT_SUPPORTED</code>	the matrix type is not supported.

8.2. `cusparse<t>csrmm2`

```

cusparseStatus_t
cusparseScsrmm2(cusparseHandle_t      handle,
                 cusparseOperation_t   transA,
                 cusparseOperation_t   transB,
                 int                    m,
                 int                    n,
                 int                    k,
                 int                    nnz,
                 const float            *alpha,
                 const cusparseMatDescr_t descrA,
                 const float            *csrValA,
                 const int              *csrRowPtrA,
                 const int              *csrColIndA,
                 const float            *B,
                 int                    ldb,
                 const float            *beta,
                 float                  *C,
                 int                    ldc)

cusparseStatus_t
cusparseDcsrmm2(cusparseHandle_t      handle,
                 cusparseOperation_t   transA,
                 cusparseOperation_t   transB,
                 int                    m,
                 int                    n,
                 int                    k,
                 int                    nnz,
                 const double           *alpha,
                 const cusparseMatDescr_t descrA,
                 const double           *csrValA,
                 const int              *csrRowPtrA,
                 const int              *csrColIndA,
                 const double           *B,
                 int                    ldb,
                 const double           *beta,
                 double                 *C,
                 int                    ldc)

cusparseStatus_t
cusparseCcsrmm2(cusparseHandle_t      handle,
                 cusparseOperation_t   transA,
                 cusparseOperation_t   transB,
                 int                    m,

```

```

        int n,
        int k,
        int nnz,
        const cuComplex *alpha,
        const cusparseMatDescr_t descrA,
        const cuComplex *csrValA,
        const int *csrRowPtrA,
        const int *csrColIndA,
        const cuComplex *B,
        int ldb,
        const cuComplex *beta,
        cuComplex *C,
        int ldc)
cusparseStatus_t
cusparseZcsrmm2(cusparseHandle_t handle,
               cusparseOperation_t transA,
               cusparseOperation_t transB,
               int m,
               int n,
               int k,
               int nnz,
               const cuDoubleComplex *alpha,
               const cusparseMatDescr_t descrA,
               const cuDoubleComplex *csrValA,
               const int *csrRowPtrA,
               const int *csrColIndA,
               const cuDoubleComplex *B,
               int ldb,
               const cuDoubleComplex *beta,
               cuDoubleComplex *C,
               int ldc)

```

This function performs one of the following matrix-matrix operation

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

where A is $m \times k$ sparse matrix (that is defined in CSR storage format by the three arrays **csrValA**, **csrRowPtrA**, and **csrColIndA**), B and C are dense matrices, α and β are scalars, and

$$\text{op}(A) = \begin{cases} A & \text{if transA} == \text{CUSPARSE_OPERATION_NON_TRANPOSE} \\ A^T & \text{if transA} == \text{CUSPARSE_OPERATION_TRANPOSE} \\ A^H & \text{if transA} == \text{CUSPARSE_OPERATION_CONJUGATE_TRANPOSE} \end{cases}$$

,

$$\text{op}(B) = \begin{cases} B & \text{if transB} == \text{CUSPARSE_OPERATION_NON_TRANPOSE} \\ B^T & \text{if transB} == \text{CUSPARSE_OPERATION_TRANPOSE} \\ B^H & \text{not supported} \end{cases}$$

If **op(B) = B**, **cusparse<t>csrmm2** is the same as **cusparse<t>csrmm**. Otherwise only **op(A) = A** is supported and matrix type must be **CUSPARSE_MATRIX_TYPE_GENERAL**.

The motivation of **transpose(B)** is to improve memory access of matrix B. The computational pattern of **A*transpose(B)** with matrix B in column-major order is equivalent to **A*B** with matrix B in row-major order.

In practice, no operation in iterative solver or eigenvalue solver uses **A*transpose(B)**. However we can perform **A*transpose(transpose(B))** which is the same as **A*B**. For

example, suppose A is $m \times k$, B is $k \times n$ and C is $m \times n$, the following code shows usage of `cusparseDcsrmm`.

```
// A is m*k, B is k*n and C is m*n
const int ldb_B = k ; // leading dimension of B
const int ldc = m ; // leading dimension of C
// perform C:=alpha*A*B + beta*C
cusparseSetMatType(descrA, CUSPARSE_MATRIX_TYPE_GENERAL ) ;
cusparseDcsrmm2(cusparse_handle,
                CUSPARSE_OPERATION_NON_TRANSPOSE,
                m, n, k, nnz, alpha,
                descrA, csrValA, csrRowPtrA, csrColIndA,
                B, ldb_B,
                beta, C, ldc);
```

Instead of using $A \times B$, our proposal is to transpose B to B^t first by calling `cublas<type>gemt`, then to perform $A \times \text{transpose}(B)$.

```
// step 1: Bt := transpose(B)
double *Bt;
const int ldb_Bt = n ; // leading dimension of Bt
cudaMalloc((void**)&Bt, sizeof(double)*ldb_Bt*k);
double one = 1.0;
double zero = 0.0;
cublasSetPointerMode(cublas_handle, CUBLAS_POINTER_MODE_HOST);
cublasDgemt(cublas_handle, CUBLAS_OP_T, CUBLAS_OP_T,
            n, k, &one, B, int ldb_B, &zero, Bt, int ldb_Bt, Bt, ldb_Bt);

// step 2: perform C:=alpha*A*transpose(Bt) + beta*C
cusparseDcsrmm2(cusparse_handle,
                CUSPARSE_OPERATION_NON_TRANSPOSE,
                CUSPARSE_OPERATION_TRANSPOSE,
                m, n, k, nnz, alpha,
                descrA, csrValA, csrRowPtrA, csrColIndA,
                Bt, ldb_Bt,
                beta, C, ldc);
```

Remark 1: `cublas<type>gemt` and `cusparse<type>csrmm2` are memory-bound. The complexity of `cublas<type>gemt` is $2 \times n \times k$ and the minimum complexity of `cusparse<type>csrmm2` is about $(nnz + nnz \times n + 2 \times m \times n)$. If nnz per column ($=nnz/k$) is large, it is worth paying extra cost on transposition because ' $A \times \text{transpose}(B)$ ' may be 2x faster than ' $A \times B$ ' if sparsity pattern of A is not good.

Remark 2: $A \times \text{transpose}(B)$ is only supported on compute capability 2.0 and above.

Input

handle	handle to the CUSPARSE library context.
transA	the operation $op(A)$
transB	the operation $op(B)$
m	number of rows of sparse matrix A .
n	number of columns of dense matrix $op(B)$ and C .
k	number of columns of sparse matrix A .
nnz	number of nonz-zero elements of sparse matrix A .
alpha	<type> scalar used for multiplication.

descrA	the descriptor of matrix <i>A</i> . The supported matrix type is <code>CUSPARSE_MATRIX_TYPE_GENERAL</code> , <code>CUSPARSE_MATRIX_TYPE_SYMMETRIC</code> , and <code>CUSPARSE_MATRIX_TYPE_HERMITIAN</code> . Also, the supported index bases are <code>CUSPARSE_INDEX_BASE_ZERO</code> and <code>CUSPARSE_INDEX_BASE_ONE</code> .
csrValA	<type> array of <code>nnz (= csrRowPtrA(m) - csrRowPtrA(0))</code> non-zero elements of matrix <i>A</i> .
csrRowPtrA	integer array of <code>m + 1</code> elements that contains the start of every row and the end of the last row plus one.
csrColIndA	integer array of <code>nnz (= csrRowPtrA(m) - csrRowPtrA(0))</code> column indices of the non-zero elements of matrix <i>A</i> .
B	array of dimensions <code>(ldb, n)</code> if <code>op(B)=B</code> and <code>(ldb, k)</code> otherwise.
ldb	leading dimension of B . If <code>op(B)=B</code> , it must be at least <code>max(1, k)</code> if <code>op(A)=A</code> and at least <code>max(1, m)</code> otherwise. If <code>op(B) != B</code> , it must be at least <code>max(1, n)</code> .
beta	<type> scalar used for multiplication. If beta is zero, c does not have to be a valid input.
C	array of dimensions <code>(ldc, n)</code> .
ldc	leading dimension of C . It must be at least <code>max(1, k)</code> if <code>op(A)=A</code> and at least <code>max(1, m)</code> otherwise.

Output

C	<type> updated array of dimensions <code>(ldc, n)</code> .
----------	--

Status Returned

<code>CUSPARSE_STATUS_SUCCESS</code>	the operation completed successfully.
<code>CUSPARSE_STATUS_NOT_INITIALIZED</code>	the library was not initialized.
<code>CUSPARSE_STATUS_ALLOC_FAILED</code>	the resources could not be allocated.
<code>CUSPARSE_STATUS_INVALID_VALUE</code>	invalid parameters were passed (<code>m, n, k, nnz < 0</code> or <code>ldb</code> and <code>ldc</code> are incorrect).
<code>CUSPARSE_STATUS_ARCH_MISMATCH</code>	if <code>op(B)=B</code> the device does not support double precision or if <code>op(B)=transpose(B)</code> the device is below compute capability 2.0.
<code>CUSPARSE_STATUS_EXECUTION_FAILED</code>	the function failed to launch on the GPU
<code>CUSPARSE_STATUS_INTERNAL_ERROR</code>	an internal operation failed.
<code>CUSPARSE_STATUS_MATRIX_TYPE_NOT_SUPPORTED</code>	<code>CUSPARSE_MATRIX_TYPE_TRIANGULAR</code> is not supported if <code>op(B)=B</code> and only

CUSPARSE_MATRIX_TYPE_GENERAL is supported otherwise.

8.3. cusparse<t>csrsm_analysis

```

cusparseStatus_t
cusparseScsrsm_analysis(cusparseHandle_t handle,
                        cusparseOperation_t transA,
                        int m, int nnz,
                        const cusparseMatDescr_t descrA,
                        const float *csrValA,
                        const int *csrRowPtrA, const int *csrColIndA,
                        cusparseSolveAnalysisInfo_t info)

cusparseStatus_t
cusparseDcsrsm_analysis(cusparseHandle_t handle,
                        cusparseOperation_t transA,
                        int m, int nnz,
                        const cusparseMatDescr_t descrA,
                        const double *csrValA,
                        const int *csrRowPtrA, const int *csrColIndA,
                        cusparseSolveAnalysisInfo_t info)

cusparseStatus_t
cusparseCcsrsm_analysis(cusparseHandle_t handle,
                        cusparseOperation_t transA,
                        int m, int nnz,
                        const cusparseMatDescr_t descrA,
                        const cuComplex *csrValA,
                        const int *csrRowPtrA, const int *csrColIndA,
                        cusparseSolveAnalysisInfo_t info)

cusparseStatus_t
cusparseZcsrsm_analysis(cusparseHandle_t handle,
                        cusparseOperation_t transA,
                        int m, int nnz,
                        const cusparseMatDescr_t descrA,
                        const cuDoubleComplex *csrValA,
                        const int *csrRowPtrA, const int *csrColIndA,
                        cusparseSolveAnalysisInfo_t info)

```

This function performs the analysis phase of the solution of a sparse triangular linear system

$$\text{op}(A) * Y = \alpha * X$$

with multiple right-hand-sides, where A is $m \times m$ sparse matrix (that is defined in CSR storage format by the three arrays **csrValA**, **csrRowPtrA**, and **csrColIndA**), X and Y are the right-hand-side and the solution dense matrices, α is a scalar, and

$$\text{op}(A) = \begin{cases} A & \text{if trans == CUSPARSE_OPERATION_NON_TRANSPOSE} \\ A^T & \text{if trans == CUSPARSE_OPERATION_TRANSPOSE} \\ A^H & \text{if trans == CUSPARSE_OPERATION_CONJUGATE_TRANSPOSE} \end{cases}$$

It is expected that this function will be executed only once for a given matrix and a particular operation type.

This function requires significant amount of extra storage that is proportional to the matrix size. It is executed asynchronously with respect to the host and it may return control to the application on the host before the result is ready.

Input

handle	handle to the CUSPARSE library context.
transA	the operation $\text{op}(A)$
m	number of rows of matrix A .
nnz	number of nonz-zero elements of matrix A .
descrA	the descriptor of matrix A . The supported matrix types are <code>CUSPARSE_MATRIX_TYPE_TRIANGULAR</code> and <code>CUSPARSE_MATRIX_TYPE_GENERAL</code> , while the supported diagonal types are <code>CUSPARSE_DIAG_TYPE_UNIT</code> and <code>CUSPARSE_DIAG_TYPE_NON_UNIT</code> .
csrValA	<type> array of $\text{nnz} (= \text{csrRowPtrA}(m) - \text{csrRowPtrA}(0))$ non-zero elements of matrix A .
csrRowPtrA	integer array of $m + 1$ elements that contains the start of every row and the end of the last row plus one.
csrColIndA	integer array of $\text{nnz} (= \text{csrRowPtrA}(m) - \text{csrRowPtrA}(0))$ column indices of the non-zero elements of matrix A .
info	structure initialized using <code>cusparseCreateSolveAnalysisInfo</code> .

Output

info	structure filled with information collected during the analysis phase (that should be passed to the solve phase unchanged).
-------------	---

Status Returned

<code>CUSPARSE_STATUS_SUCCESS</code>	the operation completed successfully.
<code>CUSPARSE_STATUS_NOT_INITIALIZED</code>	the library was not initialized.
<code>CUSPARSE_STATUS_ALLOC_FAILED</code>	the resources could not be allocated.
<code>CUSPARSE_STATUS_INVALID_VALUE</code>	invalid parameters were passed ($m, \text{nnz} < 0$).
<code>CUSPARSE_STATUS_ARCH_MISMATCH</code>	the device does not support double precision.
<code>CUSPARSE_STATUS_EXECUTION_FAILED</code>	the function failed to launch on the GPU
<code>CUSPARSE_STATUS_INTERNAL_ERROR</code>	an internal operation failed.
<code>CUSPARSE_STATUS_MATRIX_TYPE_NOT_SUPPORTED</code>	the matrix type is not supported.

8.4. `cusparse<t>csrsm_solve`

```
cusparseStatus_t
cusparseScsrsm_solve(cusparseHandle_t handle,
                    cusparseOperation_t transA,
                    int m, int n, const float *alpha,
```

```

        const cusparseMatDescr_t descrA,
        const float               *csrValA,
        const int *csrRowPtrA, const int *csrColIndA,
        cusparseSolveAnalysisInfo_t info,
        const float               *X, int ldx,
        float                     *Y, int ldy)
cusparseStatus_t
cusparseDcsrsm_solve(cusparseHandle_t handle,
                    cusparseOperation_t transA,
                    int m, int n, const double *alpha,
                    const cusparseMatDescr_t descrA,
                    const double               *csrValA,
                    const int *csrRowPtrA, const int *csrColIndA,
                    cusparseSolveAnalysisInfo_t info,
                    const double               *X, int ldx,
                    double                     *Y, int ldy)
cusparseStatus_t
cusparseCcsrsm_solve(cusparseHandle_t handle,
                    cusparseOperation_t transA,
                    int m, int n, const cuComplex *alpha,
                    const cusparseMatDescr_t descrA,
                    const cuComplex           *csrValA,
                    const int *csrRowPtrA, const int *csrColIndA,
                    cusparseSolveAnalysisInfo_t info,
                    const cuComplex           *X, int ldx,
                    cuComplex                 *Y, int ldy)
cusparseStatus_t
cusparseZcsrsm_solve(cusparseHandle_t handle,
                    cusparseOperation_t transA,
                    int m, int n, const cuDoubleComplex *alpha,
                    const cusparseMatDescr_t descrA,
                    const cuDoubleComplex *csrValA,
                    const int *csrRowPtrA, const int *csrColIndA,
                    cusparseSolveAnalysisInfo_t info,
                    const cuDoubleComplex *X, int ldx,
                    cuDoubleComplex *Y, int ldy)

```

This function performs the solve phase of the solution of a sparse triangular linear system

$$\text{op}(A) * Y = \alpha * X$$

with multiple right-hand-sides, where A is $m \times n$ sparse matrix (that is defined in CSR storage format by the three arrays **csrValA**, **csrRowPtrA**, and **csrColIndA**), X and Y are the right-hand-side and the solution dense matrices, α is a scalar, and

$$\text{op}(A) = \begin{cases} A & \text{if trans == CUSPARSE_OPERATION_NON_TRANSPOSE} \\ A^T & \text{if trans == CUSPARSE_OPERATION_TRANSPOSE} \\ A^H & \text{if trans == CUSPARSE_OPERATION_CONJUGATE_TRANSPOSE} \end{cases}$$

This function may be executed multiple times for a given matrix and a particular operation type.

This function requires no extra storage. It is executed asynchronously with respect to the host and it may return control to the application on the host before the result is ready.

Input

handle	handle to the CUSPARSE library context.
---------------	---

transA	the operation $\text{op}(A)$
m	number of rows and columns of matrix A .
n	number of columns of matrix X and Y .
alpha	<type> scalar used for multiplication.
descrA	the descriptor of matrix A . The supported matrix types are <code>CUSPARSE_MATRIX_TYPE_TRIANGULAR</code> and <code>CUSPARSE_MATRIX_TYPE_GENERAL</code> , while the supported diagonal types are <code>CUSPARSE_DIAG_TYPE_UNIT</code> and <code>CUSPARSE_DIAG_TYPE_NON_UNIT</code> .
csrValA	<type> array of <code>nnz (= csrRowPtrA(m) - csrRowPtrA(0))</code> non-zero elements of matrix A .
csrRowPtrA	integer array of <code>m + 1</code> elements that contains the start of every row and the end of the last row plus one.
csrColIndA	integer array of <code>nnz (= csrRowPtrA(m) - csrRowPtrA(0))</code> column indices of the non-zero elements of matrix A .
info	structure with information collected during the analysis phase (that should be passed to the solve phase unchanged).
x	<type> right-hand-side array of dimensions <code>(ldx, n)</code> .
ldx	leading dimension of x . (that is $\geq \max(1, m)$).

Output

y	<type> solution array of dimensions <code>(ldy, n)</code> .
ldy	leading dimension of y . (that is $\geq \max(1, m)$).

Status Returned

<code>CUSPARSE_STATUS_SUCCESS</code>	the operation completed successfully.
<code>CUSPARSE_STATUS_NOT_INITIALIZED</code>	the library was not initialized.
<code>CUSPARSE_STATUS_INVALID_VALUE</code>	invalid parameters were passed (<code>m < 0</code>).
<code>CUSPARSE_STATUS_ARCH_MISMATCH</code>	the device does not support double precision.
<code>CUSPARSE_STATUS_MAPPING_ERROR</code>	the texture binding failed.
<code>CUSPARSE_STATUS_EXECUTION_FAILED</code>	the function failed to launch on the GPU
<code>CUSPARSE_STATUS_INTERNAL_ERROR</code>	an internal operation failed.
<code>CUSPARSE_STATUS_MATRIX_TYPE_NOT_SUPPORTED</code>	the matrix type is not supported.

Chapter 9.

CUSPARSE EXTRA FUNCTION REFERENCE

This chapter describes the extra routines used to manipulate sparse matrices.

9.1. `cusparse<t>csrgeam`

```
cusparseStatus_t
cusparseXcsrgeamNnz(cusparseHandle_t handle, int m, int n,
    const cusparseMatDescr_t descrA, int nnzA,
    const int *csrRowPtrA, const int *csrColIndA,
    const cusparseMatDescr_t descrB, int nnzB,
    const int *csrRowPtrB, const int *csrColIndB,
    const cusparseMatDescr_t descrC, int *csrRowPtrC,
    int *nnzTotalDevHostPtr)
cusparseStatus_t
cusparseScsrgeam(cusparseHandle_t handle, int m, int n,
    const float *alpha,
    const cusparseMatDescr_t descrA, int nnzA,
    const float *csrValA, const int *csrRowPtrA, const int *csrColIndA,
    const float *beta,
    const cusparseMatDescr_t descrB, int nnzB,
    const float *csrValB, const int *csrRowPtrB, const int *csrColIndB,
    const cusparseMatDescr_t descrC,
    float *csrValC, int *csrRowPtrC, int *csrColIndC)
cusparseStatus_t
cusparseDcsrgeam(cusparseHandle_t handle, int m, int n,
    const double *alpha,
    const cusparseMatDescr_t descrA, int nnzA,
    const double *csrValA, const int *csrRowPtrA, const int *csrColIndA,
    const double *beta,
    const cusparseMatDescr_t descrB, int nnzB,
    const double *csrValB, const int *csrRowPtrB, const int *csrColIndB,
    const cusparseMatDescr_t descrC,
    double *csrValC, int *csrRowPtrC, int *csrColIndC)
cusparseStatus_t
cusparseCcsrgeam(cusparseHandle_t handle, int m, int n,
    const cuComplex *alpha,
    const cusparseMatDescr_t descrA, int nnzA,
    const cuComplex *csrValA, const int *csrRowPtrA, const int
    *csrColIndA,
    const cuComplex *beta,
    const cusparseMatDescr_t descrB, int nnzB,
```

```

    const cuComplex *csrValB, const int *csrRowPtrB, const int
    *csrColIndB,
    const cusparseMatDescr_t descrC,
    cuComplex *csrValC, int *csrRowPtrC, int *csrColIndC)
cusparseStatus_t
cusparseZcsrgeam(cusparseHandle_t handle, int m, int n,
    const cuDoubleComplex *alpha,
    const cusparseMatDescr_t descrA, int nnzA,
    const cuDoubleComplex *csrValA, const int *csrRowPtrA,
    const int *csrColIndA,
    const cuDoubleComplex *beta,
    const cusparseMatDescr_t descrB, int nnzB,
    const cuDoubleComplex *csrValB, const int *csrRowPtrB,
    const int *csrColIndB,
    const cusparseMatDescr_t descrC,
    cuDoubleComplex *csrValC, int *csrRowPtrC, int *csrColIndC)

```

This function performs following matrix-matrix operation

$$C = \alpha * A + \beta * B$$

where A , B and C are $m \times n$ sparse matrices (defined in CSR storage format by the three arrays **csrValA**|**csrValB**|**csrValC**, **csrRowPtrA**|**csrRowPtrB**|**csrRowPtrC**, and **csrColIndA**|**csrColIndB**|**csrColIndC** respectively), and α and β are scalars. Since A and B have different sparsity patterns, CUSPARSE adopts two-step approach to complete sparse matrix C . In the first step, the user allocates **csrRowPtrC** of $m+1$ elements and uses function **cusparseXcsrgeamNnz** to determine **csrRowPtrC** and total number of nonzero elements. In the second step, the user gathers $nnzC$ (number of non-zero elements of matrix C) from either (**nnzC=*nnzTotalDevHostPtr**) or (**nnzC=csrRowPtrC(m)-csrRowPtrC(0)**) and allocates **csrValC**, **csrColIndC** of $nnzC$ elements respectively, then finally calls function **cusparse[S|D|C|Z]csrgeam** to complete matrix C .

The general procedure is as follows:

```

int baseC, nnzC;
// nnzTotalDevHostPtr points to host memory
int *nnzTotalDevHostPtr = &nnzC;
cusparseSetPointerNode(handle, CUSPARSE_POINTER_MODE_HOST);
cudaMalloc((void**)&csrRowPtrC, sizeof(int)*(m+1));
cusparseXcsrgeamNnz(handle, m, n,
    descrA, nnzA, csrRowPtrA, csrColIndA,
    descrB, nnzB, csrRowPtrB, csrColIndB,
    descrC, csrRowPtrC, nnzTotalDevHostPtr);
if (NULL != nnzTotalDevHostPtr){
    nnzC = *nnzTotalDevHostPtr;
}else{
    cudaMemcpy(&nnzC, csrRowPtrC+m, sizeof(int), cudaMemcpyDeviceToHost);
    cudaMemcpy(&baseC, csrRowPtrC, sizeof(int), cudaMemcpyDeviceToHost);
    nnzC -= baseC;
}
cudaMalloc((void**)&csrColIndC, sizeof(int)*nnzC);
cudaMalloc((void**)&csrValC, sizeof(float)*nnzC);
cusparseScsrgeam(handle, m, n,
    alpha,
    descrA, nnzA,
    csrValA, csrRowPtrA, csrColIndA,
    beta,
    descrB, nnzB,
    csrValB, csrRowPtrB, csrColIndB,
    descrC,

```

```
csrValC, csrRowPtrC, csrColIndC);
```

Several comments on csrgeam:

1. CUSPARSE does not support other three combinations, NT, TN and TT. In order to do any one of above three, the user should use the routine csr2csc to convert $A|B$ to $A^T|B^T$.
2. Only **CUSPARSE_MATRIX_TYPE_GENERAL** is supported, if either A or B is symmetric or hermitian, then the user must extend the matrix to a full one and reconfigure MatrixType field of descriptor to **CUSPARSE_MATRIX_TYPE_GENERAL**.
3. If the sparsity pattern of matrix C is known, then the user can skip the call to function cusparseXcsrgeamNnz. For example, suppose that the user has an iterative algorithm which would update A and B iteratively but keep sparsity patterns. The user can call function cusparseXcsrgeamNnz once to setup sparsity pattern of C , then call function cusparse[S|D|C|Z]geam only for each iteration.
4. The pointers, alpha and beta, must be valid.
5. CUSPARSE would not consider special case when alpha or beta is zero. The sparsity pattern of C is independent of value of alpha and beta. If the user want $C = 0 \times A + 1 \times B^T$, then csr2csc is better than csrgeam.

Input

handle	handle to the CUSPARSE library context.
m	number of rows of sparse matrix A, B, C .
n	number of columns of sparse matrix A, B, C .
alpha	<type> scalar used for multiplication.
descrA	the descriptor of matrix A . The supported matrix type is CUSPARSE_MATRIX_TYPE_GENERAL only.
nnzA	number of nonz-zero elements of sparse matrix A .
csrValA	<type> array of $nnzA (= csrRowPtrA(m) - csrRowPtrA(0))$ non-zero elements of matrix A .
csrRowPtrA	integer array of $m + 1$ elements that contains the start of every row and the end of the last row plus one.
csrColIndA	integer array of $nnzA (= csrRowPtrA(m) - csrRowPtrA(0))$ column indices of the non-zero elements of matrix A .
beta	<type> scalar used for multiplication. If beta is zero, y does not have to be a valid input.
descrB	the descriptor of matrix B . The supported matrix type is CUSPARSE_MATRIX_TYPE_GENERAL only.
nnzB	number of nonz-zero elements of sparse matrix B .
csrValB	<type> array of $nnzB (= csrRowPtrB(m) - csrRowPtrB(0))$ non-zero elements of matrix B .
csrRowPtrB	integer array of $m + 1$ elements that contains the start of every row and the end of the last row plus one.

csrColIndB	integer array of nnzB ($= \text{csrRowPtrB}(\text{m}) - \text{csrRowPtrB}(0)$) column indices of the non-zero elements of matrix <i>B</i> .
descrC	the descriptor of matrix <i>C</i> . The supported matrix type is CUSPARSE_MATRIX_TYPE_GENERAL only.

Output

csrValC	<type> array of nnzC ($= \text{csrRowPtrC}(\text{m}) - \text{csrRowPtrC}(0)$) non-zero elements of matrix <i>C</i> .
csrRowPtrC	integer array of $\text{m} + 1$ elements that contains the start of every row and the end of the last row plus one.
csrColIndC	integer array of nnzC ($= \text{csrRowPtrC}(\text{m}) - \text{csrRowPtrC}(0)$) column indices of the non-zero elements of matrix <i>C</i> .
nnzTotalDevHostPtr	total number of nonzero elements in device or host memory. It is equal to $(\text{csrRowPtrC}(\text{m}) - \text{csrRowPtrC}(0))$.

Status Returned

CUSPARSE_STATUS_SUCCESS	the operation completed successfully.
CUSPARSE_STATUS_NOT_INITIALIZED	the library was not initialized.
CUSPARSE_STATUS_ALLOC_FAILED	the resources could not be allocated.
CUSPARSE_STATUS_INVALID_VALUE	invalid parameters were passed ($\text{m}, \text{n}, \text{nnz} < 0$, IndexBase of descrA , descrB , descrC is not base-0 or base-1, or alpha or beta is nil).
CUSPARSE_STATUS_ARCH_MISMATCH	the device does not support double precision.
CUSPARSE_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU
CUSPARSE_STATUS_MATRIX_TYPE_NOT_SUPPORTED	the matrix type is not supported.
CUSPARSE_STATUS_INTERNAL_ERROR	an internal operation failed.

9.2. cusparse<t>csrgermm

```

cusparseStatus_t
cusparseXcsrgermmNnz(cusparseHandle_t handle,
    cusparseOperation_t transA, cusparseOperation_t transB,
    int m, int n, int k,
    const cusparseMatDescr_t descrA, const int nnzA,

    const int *csrRowPtrA, const int *csrColIndA,
    const cusparseMatDescr_t descrB, const int nnzB,

    const int *csrRowPtrB, const int *csrColIndB,
    const cusparseMatDescr_t descrC, int *csrRowPtrC,
    int *nnzTotalDevHostPtr )
cusparseStatus_t
cusparseScsrgermm(cusparseHandle_t handle,

```

```

    cuspOperation_t transA, cuspOperation_t transB,
    int m, int n, int k,
    const cuspMatDescr_t descrA, const int nnzA,
    const float *csrValA,
    const int *csrRowPtrA, const int *csrColIndA,
    const cuspMatDescr_t descrB, const int nnzB,

    const float *csrValB,
    const int *csrRowPtrB, const int *csrColIndB,
    const cuspMatDescr_t descrC,
    float *csrValC,
    const int *csrRowPtrC, int *csrColIndC )
cuspStatus_t
cuspDcsrgeMM(cuspHandle_t handle,
    cuspOperation_t transA, cuspOperation_t transB,
    int m, int n, int k,
    const cuspMatDescr_t descrA, const int nnzA,
    const double *csrValA,
    const int *csrRowPtrA, const int *csrColIndA,
    const cuspMatDescr_t descrB, const int nnzB,

    const double *csrValB,
    const int *csrRowPtrB, const int *csrColIndB,
    const cuspMatDescr_t descrC,
    double *csrValC,
    const int *csrRowPtrC, int *csrColIndC )
cuspStatus_t
cuspCcsrgeMM(cuspHandle_t handle,
    cuspOperation_t transA, cuspOperation_t transB,
    int m, int n, int k,
    const cuspMatDescr_t descrA, const int nnzA,
    const cuComplex *csrValA,
    const int *csrRowPtrA, const int *csrColIndA,
    const cuspMatDescr_t descrB, const int nnzB,

    const cuComplex *csrValB,
    const int *csrRowPtrB, const int *csrColIndB,
    const cuspMatDescr_t descrC,
    cuComplex *csrValC,
    const int *csrRowPtrC, int *csrColIndC )
cuspStatus_t
cuspZcsrgeMM(cuspHandle_t handle,
    cuspOperation_t transA, cuspOperation_t transB,
    int m, int n, int k,
    const cuspMatDescr_t descrA, const int nnzA,
    const cuDoubleComplex *csrValA,
    const int *csrRowPtrA, const int *csrColIndA,
    const cuspMatDescr_t descrB, const int nnzB,
    const cuDoubleComplex *csrValB,
    const int *csrRowPtrB, const int *csrColIndB,
    const cuspMatDescr_t descrC,
    cuDoubleComplex *csrValC,
    const int *csrRowPtrC, int *csrColIndC )

```

This function performs following matrix-matrix operation

$$C = \text{op}(A) * \text{op}(B)$$

where $\text{op}(A)$, $\text{op}(B)$ and C are $m \times k$, $k \times n$, and $m \times n$ sparse matrices (defined in CSR storage format by the three arrays **csrValA|csrValB|csrValC**, **csrRowPtrA|csrRowPtrB|**

csrRowPtrC, and **csrColIndA|csrColIndB|csrColIndC** respectively. The operation is defined by

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

There are four versions, NN, NT, TN and TT. NN stands for $C = A * B$, NT stands for $C = A * B^T$, TN stands for $C = A^T * B$ and TT stands for $C = A^T * A^T$.

CUSPARSE adopts two-step approach to complete sparse matrix. In the first step, the user allocates **csrRowPtrC** of **m+1** elements and uses function **cusparseXcsrsgemmNnz** to determine **csrRowPtrC** and total number of nonzero elements. In the second step, the user gathers **nnzC** (number of nonzero elements of matrix C) from either (**nnzC=*nnzTotalDevHostPtr**) or (**nnzC=csrRowPtrC(m)-csrRowPtrC(0)**) and allocates **csrValC**, **csrColIndC** of **nnzC** elements respectively, then finally calls function **cusparse[S|D|C|Z]csrsgemm** to complete matrix C.

The general procedure is as follows:

```
int baseC, nnzC;
// nnzTotalDevHostPtr points to host memory
int *nnzTotalDevHostPtr = &nnzC;
cusparseSetPointerMode(handle, CUSPARSE_POINTER_MODE_HOST);
cudaMalloc((void**)&csrRowPtrC, sizeof(int)*(m+1));
cusparseXcsrsgemmNnz(handle, m, n, k,
    descrA, nnzA, csrRowPtrA, csrColIndA,
    descrB, nnzB, csrRowPtrB, csrColIndB,
    descrC, csrRowPtrC, nnzTotalDevHostPtr );
if (NULL != nnzTotalDevHostPtr){
    nnzC = *nnzTotalDevHostPtr;
}else{
    cudaMemcpy(&nnzC, csrRowPtrC+m, sizeof(int), cudaMemcpyDeviceToHost);
    cudaMemcpy(&baseC, csrRowPtrC, sizeof(int), cudaMemcpyDeviceToHost);
    nnzC -= baseC;
}
cudaMalloc((void**)&csrColIndC, sizeof(int)*nnzC);
cudaMalloc((void**)&csrValC, sizeof(float)*nnzC);
cusparseScsrsgemm(handle, transA, transB, m, n, k,
    descrA, nnzA,
    csrValA, csrRowPtrA, csrColIndA,
    descrB, nnzB,
    csrValB, csrRowPtrB, csrColIndB,
    descrC,
    csrValC, csrRowPtrC, csrColIndC);
```

Several comments on **csrsgemm**:

1. Only NN version is implemented. For NT version, matrix B is converted to B^T by **csr2csc** and call NN version. The same technique applies to TN and TT. The **csr2csc** routine would allocate working space implicitly, if the user needs memory management, then NN version is better.
2. NN version needs working space of size **nnzA** integers at least.
3. Only **CUSPARSE_MATRIX_TYPE_GENERAL** is supported, if either A or B is symmetric or hermitian, then the user must extend the matrix to a full one and reconfigure **MatrixType** field of descriptor to **CUSPARSE_MATRIX_TYPE_GENERAL**.
4. Only support devices of compute capability 2.0 or above.

Input

handle	handle to the CUSPARSE library context.
transA	the operation $\text{op}(A)$
transB	the operation $\text{op}(B)$
m	number of rows of sparse matrix $\text{op}(A)$ and c .
n	number of columns of sparse matrix $\text{op}(B)$ and c .
k	number of columns/rows of sparse matrix $\text{op}(A) / \text{op}(B)$.
descrA	the descriptor of matrix A . The supported matrix type is <code>CUSPARSE_MATRIX_TYPE_GENERAL</code> only.
nnzA	number of nonz-zero elements of sparse matrix A .
csrValA	<type> array of $\text{nnzA} (= \text{csrRowPtrA}(m) - \text{csrRowPtrA}(0))$ non-zero elements of matrix A .
csrRowPtrA	integer array of $\tilde{m}+1$ elements that contains the start of every row and the end of the last row plus one. $\tilde{m} = m$ if <code>transA == CUSPARSE_OPERATION_NON_TRANSPOSE</code> , otherwise $\tilde{m} = k$.
csrColIndA	integer array of $\text{nnzA} (= \text{csrRowPtrA}(m) - \text{csrRowPtrA}(0))$ column indices of the non-zero elements of matrix A .
descrB	the descriptor of matrix B . The supported matrix type is <code>CUSPARSE_MATRIX_TYPE_GENERAL</code> only.
nnzB	number of nonz-zero elements of sparse matrix B .
csrValB	<type> array of nnzB non-zero elements of matrix B .
csrRowPtrB	integer array of $\tilde{k}+1$ elements that contains the start of every row and the end of the last row plus one. $\tilde{k} = k$ if <code>transB == CUSPARSE_OPERATION_NON_TRANSPOSE</code> , otherwise $\tilde{k} = n$.
csrColIndB	integer array of nnzB column indices of the non-zero elements of matrix B .
descrC	the descriptor of matrix C . The supported matrix type is <code>CUSPARSE_MATRIX_TYPE_GENERAL</code> only.

Output

csrValC	<type> array of $\text{nnzC} (= \text{csrRowPtrC}(m) - \text{csrRowPtrC}(0))$ non-zero elements of matrix C .
csrRowPtrC	integer array of $m+1$ elements that contains the start of every row and the end of the last row plus one.
csrColIndC	integer array of $\text{nnzC} (= \text{csrRowPtrC}(m) - \text{csrRowPtrC}(0))$ column indices of the non-zero elements of matrix C .

<code>nnzTotalDevHostPtr</code>	total number of nonzero elements in device or host memory. It is equal to <code>(csrRowPtrC(m) - csrRowPtrC(0))</code> .
---------------------------------	--

Status Returned

<code>CUSPARSE_STATUS_SUCCESS</code>	the operation completed successfully.
<code>CUSPARSE_STATUS_NOT_INITIALIZED</code>	the library was not initialized.
<code>CUSPARSE_STATUS_ALLOC_FAILED</code>	the resources could not be allocated.
<code>CUSPARSE_STATUS_INVALID_VALUE</code>	invalid parameters were passed (<code>m,n,k < 0</code> , <code>IndexBase</code> of <code>descrA</code> , <code>descrB</code> , <code>descrC</code> is not base-0 or base-1, or <code>alpha</code> or <code>beta</code> is nil).
<code>CUSPARSE_STATUS_ARCH_MISMATCH</code>	the device does not support double precision.
<code>CUSPARSE_STATUS_EXECUTION_FAILED</code>	the function failed to launch on the GPU
<code>CUSPARSE_STATUS_MATRIX_TYPE_NOT_SUPPORTED</code>	the matrix type is not supported.
<code>CUSPARSE_STATUS_INTERNAL_ERROR</code>	an internal operation failed.

Chapter 10.

CUSPARSE PRECONDITIONERS REFERENCE

This chapter describes the routines that implement different preconditioners.

In particular, the incomplete factorizations are implemented in two phases. First, during the analysis phase, the sparse triangular matrix is analyzed to determine the dependencies between its elements by calling the appropriate `csrsv_analysis()` function. The analysis is specific to the sparsity pattern of the given matrix and selected `cusparseOperation_t` type. The information from the analysis phase is stored in the parameter of type `cusparseSolveAnalysisInfo_t` that has been initialized previously with a call to `cusparseCreateSolveAnalysisInfo()`.

Second, during the numerical factorization phase, the given coefficient matrix is factorized using the information stored in the `cusparseSolveAnalysisInfo_t` parameter by calling the appropriate `csrilu0` or `csric0` function.

The analysis phase is shared across the sparse triangular solve and the incomplete factorization and must be performed only once. While the resulting information can be passed to the numerical factorization and the sparse triangular solve multiple times.

Finally, once the incomplete factorization and all the sparse triangular solves have completed, the opaque data structure pointed to by the `cusparseSolveAnalysisInfo_t` parameter can be released by calling `cusparseDestroySolveAnalysisInfo()`.

10.1. `cusparse<t>csric0`

```
cusparseStatus_t
cusparseScsric0(cusparseHandle_t handle, cusparseOperation_t trans,
               int m, const cusparseMatDescr_t descrA,
               float *csrValM,
               const int *csrRowPtrA, const int *csrColIndA,
               cusparseSolveAnalysisInfo_t info)
cusparseStatus_t
cusparseDcsric0(cusparseHandle_t handle, cusparseOperation_t trans,
               int m, const cusparseMatDescr_t descrA,
               double *csrValM,
               const int *csrRowPtrA, const int *csrColIndA,
               cusparseSolveAnalysisInfo_t info)
cusparseStatus_t
```

```

cusparseCcsric0(cusparseHandle_t handle, cusparseOperation_t trans,
               int m, const cusparseMatDescr_t descrA,
               cuComplex *csrValM,
               const int *csrRowPtrA, const int *csrColIndA,
               cusparseSolveAnalysisInfo_t info)
cusparseStatus_t
cusparseZcsric0(cusparseHandle_t handle, cusparseOperation_t trans,
               int m, const cusparseMatDescr_t descrA,
               cuDoubleComplex *csrValM,
               const int *csrRowPtrA, const int *csrColIndA,
               cusparseSolveAnalysisInfo_t info)

```

This function computes the incomplete-Cholesky factorization with 0 fill-in and no pivoting

$$op(A) \approx R^T R$$

where A is $m \times m$ Hermitian/symmetric positive definite sparse matrix (that is defined in CSR storage format by the three arrays **csrValM**, **csrRowPtrA** and **csrColIndA**) and

$$op(A) = \begin{cases} A & \text{if trans == CUSPARSE_OPERATION_NON_TRANSPOSE} \\ A^T & \text{if trans == CUSPARSE_OPERATION_TRANSPOSE} \\ A^H & \text{if trans == CUSPARSE_OPERATION_CONJUGATE_TRANSPOSE} \end{cases}$$

Notice that only a lower or upper Hermitian/symmetric part of the matrix A is actually stored. It is overwritten by the lower or upper triangular factor R^T or R , respectively.

A call to this routine must be preceded by a call to the **csrsv_analysis** routine.

This function requires some extra storage. It is executed asynchronously with respect to the host and it may return control to the application on the host before the result is ready.

Input

handle	handle to the CUSPARSE library context.
trans	the operation $op(A)$
m	number of rows and columns of matrix A .
descrA	the descriptor of matrix A . The supported matrix type is CUSPARSE_MATRIX_TYPE_SYMMETRIC and CUSPARSE_MATRIX_TYPE_HERMITIAN . Also, the supported index bases are CUSPARSE_INDEX_BASE_ZERO and CUSPARSE_INDEX_BASE_ONE .
csrValM	<type> array of $nnz (= \text{csrRowPtrA}(m) - \text{csrRowPtrA}(0))$ non-zero elements of matrix A .
csrRowPtrA	integer array of $m + 1$ elements that contains the start of every row and the end of the last row plus one.
csrColIndA	integer array of $nnz (= \text{csrRowPtrA}(m) - \text{csrRowPtrA}(0))$ column indices of the non-zero elements of matrix A .

info	structure with information collected during the analysis phase (that should have been passed to the solve phase unchanged).
-------------	---

Output

csrValM	<type> matrix containing the incomplete-Cholesky lower or upper triangular factor.
----------------	--

Status Returned

CUSPARSE_STATUS_SUCCESS	the operation completed successfully.
CUSPARSE_STATUS_NOT_INITIALIZED	the library was not initialized.
CUSPARSE_STATUS_ALLOC_FAILED	the resources could not be allocated.
CUSPARSE_STATUS_INVALID_VALUE	invalid parameters were passed ($m < 0$).
CUSPARSE_STATUS_ARCH_MISMATCH	the device does not support double precision.
CUSPARSE_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU
CUSPARSE_STATUS_INTERNAL_ERROR	an internal operation failed.
CUSPARSE_STATUS_MATRIX_TYPE_NOT_SUPPORTED	the matrix type is not supported.

10.2. `cusparse<t>csrilu0`

```

cusparseStatus_t
cusparseScsrilu0(cusparseHandle_t handle, cusparseOperation_t trans,
                int m, const cusparseMatDescr_t descrA,
                float *csrValM,
                const int *csrRowPtrA, const int *csrColIndA,
                cusparseSolveAnalysisInfo_t info)

cusparseStatus_t
cusparseDcsrilu0(cusparseHandle_t handle, cusparseOperation_t trans,
                int m, const cusparseMatDescr_t descrA,
                double *csrValM,
                const int *csrRowPtrA, const int *csrColIndA,
                cusparseSolveAnalysisInfo_t info)

cusparseStatus_t
cusparseCcsrilu0(cusparseHandle_t handle, cusparseOperation_t trans,
                int m, const cusparseMatDescr_t descrA,
                cuComplex *csrValM,
                const int *csrRowPtrA, const int *csrColIndA,
                cusparseSolveAnalysisInfo_t info)

cusparseStatus_t
cusparseZcsrilu0(cusparseHandle_t handle, cusparseOperation_t trans,
                int m, const cusparseMatDescr_t descrA,
                cuDoubleComplex *csrValM,
                const int *csrRowPtrA, const int *csrColIndA,
                cusparseSolveAnalysisInfo_t info)

```

This function computes the incomplete-LU factorization with 0 fill-in and no pivoting

$$op(A) \approx LU$$

where A is $m \times m$ sparse matrix (that is defined in CSR storage format by the three arrays **csrValM**, **csrRowPtrA** and **csrColIndA**) and

$$\text{op}(A) = \begin{cases} A & \text{if trans} == \text{CUSPARSE_OPERATION_NON_TRANPOSE} \\ A^T & \text{if trans} == \text{CUSPARSE_OPERATION_TRANPOSE} \\ A^H & \text{if trans} == \text{CUSPARSE_OPERATION_CONJUGATE_TRANPOSE} \end{cases}$$

Notice that the diagonal of lower triangular factor L is unitary and need not be stored. Therefore the input matrix is overwritten with the resulting lower and upper triangular factor L and U , respectively.

A call to this routine must be preceded by a call to the **csrsv_analysis** routine.

This function requires some extra storage. It is executed asynchronously with respect to the host and it may return control to the application on the host before the result is ready.

Input

handle	handle to the CUSPARSE library context.
trans	the operation $\text{op}(A)$
m	number of rows and columns of matrix A .
descrA	the descriptor of matrix A . The supported matrix type is CUSPARSE_MATRIX_TYPE_GENERAL . Also, the supported index bases are CUSPARSE_INDEX_BASE_ZERO and CUSPARSE_INDEX_BASE_ONE .
csrValM	<type> array of $\text{nnz} (= \text{csrRowPtrA}(m) - \text{csrRowPtrA}(0))$ non-zero elements of matrix A .
csrRowPtrA	integer array of $m + 1$ elements that contains the start of every row and the end of the last row plus one.
csrColIndA	integer array of $\text{nnz} (= \text{csrRowPtrA}(m) - \text{csrRowPtrA}(0))$ column indices of the non-zero elements of matrix A .
info	structure with information collected during the analysis phase (that should have been passed to the solve phase unchanged).

Output

csrValM	<type> matrix containing the incomplete-LU lower and upper triangular factors.
----------------	--

Status Returned

CUSPARSE_STATUS_SUCCESS	the operation completed successfully.
CUSPARSE_STATUS_NOT_INITIALIZED	the library was not initialized.
CUSPARSE_STATUS_ALLOC_FAILED	the resources could not be allocated.
CUSPARSE_STATUS_INVALID_VALUE	invalid parameters were passed ($m < 0$).
CUSPARSE_STATUS_ARCH_MISMATCH	the device does not support double precision.

<code>CUSPARSE_STATUS_EXECUTION_FAILED</code>	the function failed to launch on the GPU
<code>CUSPARSE_STATUS_INTERNAL_ERROR</code>	an internal operation failed.
<code>CUSPARSE_STATUS_MATRIX_TYPE_NOT_SUPPORTED</code>	the matrix type is not supported.

10.3. `cusparse<t>gtsv`

```

cusparseStatus_t
cusparseSgtsv(cusparseHandle_t handle, int m, int n,
              const float *dl, const float *d, const float *du, float *B, int ldb)

cusparseStatus_t
cusparseDgtsv(cusparseHandle_t handle, int m, int n,
              const double *dl, const double *d, const double *du, double *B, int ldb)

cusparseStatus_t
cusparseCgtsv(cusparseHandle_t handle, int m, int n,
              const cuComplex *dl, const cuComplex *d, const cuComplex *du, cuComplex *B, int ldb)

cusparseStatus_t
cusparseZgtsv(cusparseHandle_t handle, int m, int n,
              const cuDoubleComplex *dl, const cuDoubleComplex *d, const cuDoubleComplex *du, cuDoubleComplex *B, int ldb)

```

This function computes the solution of a tridiagonal linear system

$$A * Y = \alpha * X$$

with multiple right-hand-sides.

The coefficient matrix A of each of these tri-diagonal linear system is defined with three vectors corresponding to its lower (**ld**), main (**d**) and upper (**ud**) matrix diagonals, while the right-hand-sides are stored in the dense matrix X . Notice that the solutions Y overwrite the right-hand-sides X on exit.

The routine does perform pivoting, which usually results in more accurate and more stable results than `cusparse<t>gtsv_nopivot` at the expense of some execution time

This routine requires significant amount of temporary extra storage (**min(m, 8) * (3+n) * sizeof(<type>)**). It is executed asynchronously with respect to the host and it may return control to the application on the host before the result is ready.

Input

handle	handle to the CUSPARSE library context.
m	the size of the linear system (must be ≥ 3).
n	number of right-hand-sides, columns of matrix B .
dl	<type> dense array containing the lower diagonal of the tri-diagonal linear system. The first element of each lower diagonal must be zero.
d	<type> dense array containing the main diagonal of the tri-diagonal linear system.

du	<type> dense array containing the upper diagonal of the tri-diagonal linear system. The last element of each upper diagonal must be zero.
B	<type> dense right-hand-side array of dimensions (ldb, n).
ldb	leading dimension of B . (that is $\geq \max(1, m)$)

Output

B	<type> dense solution array of dimensions (ldb, m).
----------	---

Status Returned

CUSPARSE_STATUS_SUCCESS	the operation completed successfully.
CUSPARSE_STATUS_NOT_INITIALIZED	the library was not initialized.
CUSPARSE_STATUS_ALLOC_FAILED	the resources could not be allocated.
CUSPARSE_STATUS_INVALID_VALUE	invalid parameters were passed ($m < 3$, $n < 0$).
CUSPARSE_STATUS_ARCH_MISMATCH	the device does not support double precision.
CUSPARSE_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU
CUSPARSE_STATUS_INTERNAL_ERROR	an internal operation failed.
CUSPARSE_STATUS_MATRIX_TYPE_NOT_SUPPORTED	the matrix type is not supported.

10.4. cusparsesgtstsv_nopivot

```

cusparsesgtstsv_nopivot(cusparsesgtstsv_handle_t handle, int m, int n,
                        const float *dl, const float *du, float *B, int ldb)
                        *d,

cusparsesgtstsv_nopivot(cusparsesgtstsv_handle_t handle, int m, int n,
                        const double *dl, const double *du, double *B, int ldb)
                        *d,

cusparsesgtstsv_nopivot(cusparsesgtstsv_handle_t handle, int m, int n,
                        const cuComplex *dl, const cuComplex *du, cuComplex *B, int ldb)
                        *d,

cusparsesgtstsv_nopivot(cusparsesgtstsv_handle_t handle, int m, int n,
                        const cuDoubleComplex *dl, const cuDoubleComplex *du, cuDoubleComplex *B, int ldb)
                        *d,

```

This function computes the solution of a tridiagonal linear system

$$A * Y = \alpha * X$$

with multiple right-hand-sides.

The coefficient matrix A of each of these tri-diagonal linear system is defined with three vectors corresponding to its lower (**ld**), main (**d**) and upper (**ud**) matrix diagonals, while the right-hand-sides are stored in the dense matrix X . Notice that the solutions Y overwrite the right-hand-sides X on exit.

The routine does not perform any pivoting and uses a combination of the Cyclic Reduction (CR) and Parallel Cyclic Reduction (PCR) algorithms to find the solution. It achieves better performance when **m** is a power of 2.

This routine requires significant amount of temporary extra storage (**m** × (3+**n**) × **sizeof**(**<type>**)). It is executed asynchronously with respect to the host and it may return control to the application on the host before the result is ready.

Input

handle	handle to the CUSPARSE library context.
m	the size of the linear system (must be ≥ 3).
n	number of right-hand-sides, columns of matrix B .
dl	<type> dense array containing the lower diagonal of the tri-diagonal linear system. The first element of each lower diagonal must be zero.
d	<type> dense array containing the main diagonal of the tri-diagonal linear system.
du	<type> dense array containing the upper diagonal of the tri-diagonal linear system. The last element of each upper diagonal must be zero.
B	<type> dense right-hand-side array of dimensions (ldb , n).
ldb	leading dimension of B . (that is $\geq \max(1, m)$)

Output

B	<type> dense solution array of dimensions (ldb , m).
----------	---

Status Returned

CUSPARSE_STATUS_SUCCESS	the operation completed successfully.
CUSPARSE_STATUS_NOT_INITIALIZED	the library was not initialized.
CUSPARSE_STATUS_ALLOC_FAILED	the resources could not be allocated.
CUSPARSE_STATUS_INVALID_VALUE	invalid parameters were passed (m <3, n <0).
CUSPARSE_STATUS_ARCH_MISMATCH	the device does not support double precision.
CUSPARSE_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU
CUSPARSE_STATUS_INTERNAL_ERROR	an internal operation failed.
CUSPARSE_STATUS_MATRIX_TYPE_NOT_SUPPORTED	the matrix type is not supported.

10.5. `cusparse<t>gtsvStridedBatch`

```

cusparseStatus_t
cusparseSgtsvStridedBatch(cusparseHandle_t handle, int m,
                          const float          *dl,
                          const float          *d,
                          const float          *du, float          *x,
                          int batchCount, int batchStride)

cusparseStatus_t
cusparseDgtsvStridedBatch(cusparseHandle_t handle, int m,
                          const double         *dl,
                          const double         *d,
                          const double         *du, double        *x,
                          int batchCount, int batchStride)

cusparseStatus_t
cusparseCgtsvStridedBatch(cusparseHandle_t handle, int m,
                          const cuComplex      *dl,
                          const cuComplex      *d,
                          const cuComplex      *du, cuComplex      *x,
                          int batchCount, int batchStride)

cusparseStatus_t
cusparseZgtsvStridedBatch(cusparseHandle_t handle, int m,
                          const cuDoubleComplex *dl,
                          const cuDoubleComplex *d,
                          const cuDoubleComplex *du, cuDoubleComplex *x,
                          int batchCount, int batchStride)

```

This function computes the solution of multiple tridiagonal linear systems

$$A^{(i)} * y^{(i)} = \alpha * x^{(i)}$$

for $i=0, \dots, \text{batchCount}$.

The coefficient matrix A of each of these tri-diagonal linear system is defined with three vectors corresponding to its lower (**ld**), main (**d**) and upper (**ud**) matrix diagonals, while the right-hand-side is stored in the vector **x**. Notice that the solution **y** overwrites the right-hand-side **x** on exit. The different matrices are assumed to be of the same size and are stored with a fixed **batchStride** in memory.

The routine does not perform any pivoting and uses a combination of the Cyclic Reduction (CR) and Parallel Cyclic Reduction (PCR) algorithms to find the solution. It achieves better performance when **m** is a power of 2.

This routine requires significant amount of temporary extra storage ($(\text{batchCount} \times (4 \times m + 2048) \times \text{sizeof}(\text{<type>}))$). It is executed asynchronously with respect to the host and it may return control to the application on the host before the result is ready.

Input

handle	handle to the CUSPARSE library context.
---------------	---

m	the size of the linear system (must be ≥ 3).
dl	<type> dense array containing the lower diagonal of the tri-diagonal linear system. The lower diagonal $dl^{(i)}$ that corresponds to the i th linear system starts at location $dl + \text{batchStride} \times i$ in memory. Also, the first element of each lower diagonal must be zero.
d	<type> dense array containing the main diagonal of the tri-diagonal linear system. The main diagonal $d^{(i)}$ that corresponds to the i th linear system starts at location $d + \text{batchStride} \times i$ in memory.
du	<type> dense array containing the upper diagonal of the tri-diagonal linear system. The upper diagonal $du^{(i)}$ that corresponds to the i th linear system starts at location $du + \text{batchStride} \times i$ in memory. Also, the last element of each upper diagonal must be zero.
x	<type> dense array that contains the right-hand-side of the tri-diagonal linear system. The right-hand-side $x^{(i)}$ that corresponds to the i th linear system starts at location $x + \text{batchStride} \times i$ in memory.
batchCount	Number of systems to solve.
batchStride	stride (number of elements) that separates the vectors of every system (must be at least m).

Output

x	<type> dense array that contains the solution of the tri-diagonal linear system. The solution $x^{(i)}$ that corresponds to the i th linear system starts at location $x + \text{batchStride} \times i$ in memory.
----------	--

Status Returned

CUSPARSE_STATUS_SUCCESS	the operation completed successfully.
CUSPARSE_STATUS_NOT_INITIALIZED	the library was not initialized.
CUSPARSE_STATUS_ALLOC_FAILED	the resources could not be allocated.
CUSPARSE_STATUS_INVALID_VALUE	invalid parameters were passed ($m < 3$, $\text{batchCount} \leq 0$, $\text{batchStride} < m$).
CUSPARSE_STATUS_ARCH_MISMATCH	the device does not support double precision.
CUSPARSE_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU
CUSPARSE_STATUS_INTERNAL_ERROR	an internal operation failed.

Chapter 11.

CUSPARSE FORMAT CONVERSION

REFERENCE

This chapter describes the conversion routines between different sparse and dense storage formats.

11.1. `cusparse<t>bsr2csr`

```
cusparseStatus_t
cusparseSbsr2csr(cusparseHandle_t handle, cusparseDirection_t dirA,
    int mb, int nb,
    const cusparseMatDescr_t descrA, const float *bsrValA,
    const int *bsrRowPtrA, const int *bsrColIndA,
    int blockDim,
    const cusparseMatDescr_t descrC,
    float *csrValC, int *csrRowPtrC, int *csrColIndC)
cusparseStatus_t
cusparseDbsr2csr(cusparseHandle_t handle, cusparseDirection_t dirA,
    int mb, int nb,
    const cusparseMatDescr_t descrA, const double *bsrValA,
    const int *bsrRowPtrA, const int *bsrColIndA,
    int blockDim,
    const cusparseMatDescr_t descrC,
    double *csrValC, int *csrRowPtrC, int *csrColIndC)
cusparseStatus_t
cusparseCbsr2csr(cusparseHandle_t handle, cusparseDirection_t dirA,
    int mb, int nb,
    const cusparseMatDescr_t descrA, const cuComplex *bsrValA,
    const int *bsrRowPtrA, const int *bsrColIndA,
    int blockDim,
    const cusparseMatDescr_t descrC,
    cuComplex *csrValC, int *csrRowPtrC, int *csrColIndC)
cusparseStatus_t
cusparseZbsr2csr(cusparseHandle_t handle, cusparseDirection_t dirA,
    int mb, int nb,
    const cusparseMatDescr_t descrA, const cuDoubleComplex *bsrValA,
    const int *bsrRowPtrA, const int *bsrColIndA,
    int blockDim,
    const cusparseMatDescr_t descrC,
    cuDoubleComplex *csrValC, int *csrRowPtrC, int *csrColIndC)
```

This function converts a sparse matrix in BSR format (that is defined by the three arrays **bsrValA**, **bsrRowPtrA**, and **bsrColIndA**) into a sparse matrix in CSR format (that is defined by arrays **csrValC**, **csrRowPtrC**, and **csrColIndC**).

Let $m(=mb * blockDim)$ be number of rows of A and $n(=nb * blockDim)$ be number of columns of A , then A and C are $m \times n$ sparse matrices. BSR format of A contains $nnzb(= csrRowPtrC(mb) - csrRowPtrC(0))$ non-zero blocks whereas sparse matrix A contains $nnz(= nnzb * blockDim^2)$ elements. The user must allocate enough space for arrays **csrRowPtrC**, **csrColIndC** and **csrValC**. The requirements are

csrRowPtrC of $m+1$ elements,

csrValC of nnz elements, and

csrColIndC of nnz elements.

The general procedure is as follows:

```
// Given BSR format (bsrRowPtrA, bsrColIndA, bsrValA) and
// blocks of BSR format are stored in column-major order.
cusparsedirection_t dirA = CUSPARSE_DIRECTION_COLUMN;
int m = mb*blockDim;
int nnzb = bsrRowPtrA[mb] - bsrRowPtrA[0]; // number of blocks
int nnz = nnzb * blockDim * blockDim; // number of elements
cudaMalloc((void**)&csrRowPtrC, sizeof(int)*(m+1));
cudaMalloc((void**)&csrColIndC, sizeof(int)*nnz);
cudaMalloc((void**)&csrValC, sizeof(float)*nnz);
cusparsesbsr2csr(handle, dirA, mb, nb,
    descrA,
    bsrValA, bsrRowPtrA, bsrColIndA,
    blockDim,
    descrC,
    csrValC, csrRowPtrC, csrColIndC);
```

Input

handle	handle to the CUSPARSE library context.
dirA	storage format of blocks, either CUSPARSE_DIRECTION_ROW or CUSPARSE_DIRECTION_COLUMN .
mb	number of block rows of sparse matrix A . The number of rows of sparse matrix C is $m(= mb * blockDim)$
nb	number of block columns of sparse matrix A . The number of columns of sparse matrix C is $n(= nb * blockDim)$
descrA	the descriptor of matrix A .
bsrValA	<type> array of $nnzb * blockDim^2$ non-zero elements of matrix A .
bsrRowPtrA	integer array of $mb+1$ elements that contains the start of every block row and the end of the last block row plus one.
bsrColIndA	integer array of $nnzb$ column indices of the non-zero blocks of matrix A .
blockDim	block dimension of sparse matrix A , larger than zero.

descrC	the descriptor of matrix <i>C</i> .
---------------	-------------------------------------

Output

csrValC	<type> array of nnz (= csrRowPtrA (m) - csrRowPtrA (0)) non-zero elements of matrix <i>C</i> .
csrRowPtrC	integer array of m +1 elements that contains the start of every row and the end of the last row plus one.
csrColIndC	integer array of nnz column indices of the non-zero elements of matrix <i>C</i> .

Status Returned

CUSPARSE_STATUS_SUCCESS	the operation completed successfully.
CUSPARSE_STATUS_NOT_INITIALIZED	the library was not initialized.
CUSPARSE_STATUS_ALLOC_FAILED	the resources could not be allocated.
CUSPARSE_STATUS_INVALID_VALUE	invalid parameters were passed (mb , nb <0, IndexBase of descrA , descrC is not base-0 or base-1, dirA is not row-major or column-major, or blockDim <1).
CUSPARSE_STATUS_ARCH_MISMATCH	the device does not support double precision.
CUSPARSE_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU
CUSPARSE_STATUS_MATRIX_TYPE_NOT_SUPPORTED	the matrix type is not supported.
CUSPARSE_STATUS_INTERNAL_ERROR	an internal operation failed.

11.2. `cusparse<t>coo2csr`

```
cusparseStatus_t
cusparseXcoo2csr(cusparseHandle_t handle, const int *cooRowInd,
                 int nnz, int m, int *csrRowPtr, cusparseIndexBase_t
                 idxBase)
```

This function converts the array containing the uncompressed row indices (corresponding to COO format) into an array of compressed row pointers (corresponding to CSR format).

It can also be used to convert the array containing the uncompressed column indices (corresponding to COO format) into an array of column pointers (corresponding to CSC format).

This function requires no extra storage. It is executed asynchronously with respect to the host and it may return control to the application on the host before the result is ready.

Input

handle	handle to the CUSPARSE library context.
cooRowInd	integer array of nnz uncompressed row indices.

nnz	number of non-zeros of the sparse matrix (that is also the length of array <code>cooRowInd</code>).
m	number of rows of matrix <i>A</i> .
idxBase	<code>CUSPARSE_INDEX_BASE_ZERO</code> or <code>CUSPARSE_INDEX_BASE_ONE</code> .

Output

csrRowPtr	integer array of <code>m+1</code> elements that contains the start of every row and the end of the last row plus one.
------------------	---

Status Returned

<code>CUSPARSE_STATUS_SUCCESS</code>	the operation completed successfully.
<code>CUSPARSE_STATUS_NOT_INITIALIZED</code>	the library was not initialized.
<code>CUSPARSE_STATUS_INVALID_VALUE</code>	<code>IdxBase</code> is neither <code>CUSPARSE_INDEX_BASE_ZERO</code> nor <code>CUSPARSE_INDEX_BASE_ONE</code> .
<code>CUSPARSE_STATUS_EXECUTION_FAILED</code>	the function failed to launch on the GPU

11.3. `cusparse<t>csc2dense`

```

cusparseStatus_t
cusparseScsc2dense(cusparseHandle_t handle, int m, int n,
                  const cusparseMatDescr_t descrA,
                  const float *cscValA,
                  const int *cscRowIndA, const int *cscColPtrA,
                  float *A, int lda)

cusparseStatus_t
cusparseDcsc2dense(cusparseHandle_t handle, int m, int n,
                  const cusparseMatDescr_t descrA,
                  const double *cscValA,
                  const int *cscRowIndA, const int *cscColPtrA,
                  double *A, int lda)

cusparseStatus_t
cusparseCcsc2dense(cusparseHandle_t handle, int m, int n,
                  const cusparseMatDescr_t descrA,
                  const cuComplex *cscValA,
                  const int *cscRowIndA, const int *cscColPtrA,
                  cuComplex *A, int lda)

cusparseStatus_t
cusparseZcsc2dense(cusparseHandle_t handle, int m, int n,
                  const cusparseMatDescr_t descrA,
                  const cuDoubleComplex *cscValA,
                  const int *cscRowIndA, const int *cscColPtrA,
                  cuDoubleComplex *A, int lda)

```

This function converts the sparse matrix in CSC format (that is defined by the three arrays **`cscValA`**, **`cscColPtrA`** and **`cscRowIndA`**) into the matrix ***A*** in dense format. The dense matrix ***A*** is filled in with the values of the sparse matrix and with zeros elsewhere.

This function requires no extra storage. It is executed asynchronously with respect to the host and it may return control to the application on the host before the result is ready.

Input

handle	handle to the CUSPARSE library context.
m	number of rows of matrix <i>A</i> .
n	number of columns of matrix <i>A</i> .
descrA	the descriptor of matrix <i>A</i> . The supported matrix type is <code>CUSPARSE_MATRIX_TYPE_GENERAL</code> . Also, the supported index bases are <code>CUSPARSE_INDEX_BASE_ZERO</code> and <code>CUSPARSE_INDEX_BASE_ONE</code> .
cscValA	<type> array of <code>nnz (= cscColPtrA(m) - cscColPtrA(0))</code> non-zero elements of matrix <i>A</i> .
cscRowIndA	integer array of <code>nnz (= cscColPtrA(m) - cscColPtrA(0))</code> row indices of the non-zero elements of matrix <i>A</i> .
cscColPtrA	integer array of <code>n+1</code> elements that contains the start of every row and the end of the last column plus one.
lda	leading dimension of dense array <i>A</i> .

Output

A	array of dimensions (<code>lda</code> , <code>n</code>) that is filled in with the values of the sparse matrix.
----------	---

Status Returned

<code>CUSPARSE_STATUS_SUCCESS</code>	the operation completed successfully.
<code>CUSPARSE_STATUS_NOT_INITIALIZED</code>	the library was not initialized.
<code>CUSPARSE_STATUS_INVALID_VALUE</code>	invalid parameters were passed (<code>m, n < 0</code>).
<code>CUSPARSE_STATUS_ARCH_MISMATCH</code>	the device does not support double precision.
<code>CUSPARSE_STATUS_EXECUTION_FAILED</code>	the function failed to launch on the GPU
<code>CUSPARSE_STATUS_MATRIX_TYPE_NOT_SUPPORTED</code>	the matrix type is not supported.

11.4. `cusparse<t>csr2bsr`

```

cusparseStatus_t
cusparseXcsr2bsrNnz(cusparseHandle_t handle, cusparseDirection_t dirA,
    int m, int n,
    const cusparseMatDescr_t descrA,
    const int *csrRowPtrA, const int *csrColIndA,
    int blockDim,
    const cusparseMatDescr_t descrC,
    int *bsrRowPtrC,
    int *nnzTotalDevHostPtr)
cusparseStatus_t
cusparseScsr2bsr(cusparseHandle_t handle, cusparseDirection_t dirA,
    int m, int n,

```

```

    const cusparseMatDescr_t descrA, const float *csrValA,
    const int *csrRowPtrA, const int *csrColIndA,
    int blockDim,
    const cusparseMatDescr_t descrC,
    float *bsrValC, int *bsrRowPtrC, int *bsrColIndC)
cusparseStatus_t
cusparseDcsr2bsr(cusparseHandle_t handle, cusparseDirection_t dirA,
    int m, int n,
    const cusparseMatDescr_t descrA, const double *csrValA,
    const int *csrRowPtrA, const int *csrColIndA,
    int blockDim,
    const cusparseMatDescr_t descrC,
    double *bsrValC, int *bsrRowPtrC, int *bsrColIndC)
cusparseStatus_t
cusparseCcsr2bsr(cusparseHandle_t handle, cusparseDirection_t dirA,
    int m, int n,
    const cusparseMatDescr_t descrA, const cuComplex *csrValA,
    const int *csrRowPtrA, const int *csrColIndA,
    int blockDim,
    const cusparseMatDescr_t descrC,
    cuComplex *bsrValC, int *bsrRowPtrC, int *bsrColIndC)
cusparseStatus_t
cusparseZcsr2bsr(cusparseHandle_t handle, cusparseDirection_t dirA,
    int m, int n,
    const cusparseMatDescr_t descrA, const cuDoubleComplex *csrValA,
    const int *csrRowPtrA, const int *csrColIndA,
    int blockDim,
    const cusparseMatDescr_t descrC,
    cuDoubleComplex *bsrValC, int *bsrRowPtrC, int *bsrColIndC)

```

This function converts a sparse matrix in CSR format (that is defined by the three arrays **csrValA**, **csrRowPtrA** and **csrColIndA**) into a sparse matrix in BSR format (that is defined by arrays **bsrValC**, **bsrRowPtrC** and **bsrColIndC**).

A is $m \times n$ sparse matrix and is $(mb \times \text{blockDim}) \times (nb \times \text{blockDim})$ sparse matrix.

where $mb = \frac{m + \text{blockDim} - 1}{\text{blockDim}}$ is number of block rows of A and

$nb = \frac{n + \text{blockDim} - 1}{\text{blockDim}}$ is number of block columns of A. m and n need not be multiple of blockDim . If so, then zeros are filled in.

CUSPARSE adopts two-step approach to do the conversion. First, the user allocates **bsrRowPtrC** of $mb+1$ elements and uses function **cusparseXcsr2bsrNnz** to determine number of non-zero block columns per block row. Second, the user gathers **nnzb** (number of non-zero block columns of matrix A) from either (**nnzb=*nnzTotalDevHostPtr**) or (**nnzb=bsrRowPtrC(mb)-bsrRowPtrC(0)**) and allocates **bsrValC** of $nnzb \times \text{blockDim}^2$ elements and **bsrColIndC** of $nnzb$ elements. Finally function **cusparse[S|D|C|Z]csr2bsr** is called to complete the conversion.

The general procedure is as follows:

```

// Given CSR format (csrRowPtrA, csrColIndA, csrValA) and
// blocks of BSR format are stored in column-major order.
cusparseDirection_t dirA = CUSPARSE_DIRECTION_COLUMN;
int base, nnzb;
int mb = (m + blockDim-1)/blockDim;
// nnzTotalDevHostPtr points to host memory
int *nnzTotalDevHostPtr = &nnzb;
cudaMalloc((void**)&bsrRowPtrC, sizeof(int) * (mb+1));
cusparseXcsr2bsrNnz(handle, dirA, m, n,

```

```

        descrA, csrRowPtrA, csrColIndA,
        blockDim,
        descrC, bsrRowPtrC,
        nnzTotalDevHostPtr);
if (NULL != nnzTotalDevHostPtr){
    nnzb = *nnzTotalDevHostPtr;
}else{
    cudaMemcpy(&nnzb, bsrRowPtrC+mb, sizeof(int), cudaMemcpyDeviceToHost);
    cudaMemcpy(&base, bsrRowPtrC, sizeof(int), cudaMemcpyDeviceToHost);
    nnzb -= base;
}
cudaMalloc((void**)&bsrColIndC, sizeof(int)*nnzb);
cudaMalloc((void**)&bsrValC, sizeof(float)*(blockDim*blockDim)*nnzb);
cusparseScsr2bsr(handle, dirA, m, n,
        descrA,
        csrValA, csrRowPtrA, csrColIndA,
        blockDim,
        descrC,
        bsrValC, bsrRowPtrC, bsrColIndC);

```

If *blockDim* is large (typically a block cannot fit into shared memory), then `csr2bsr` routines will allocate temporary integer array of size $mb * blockDim$. If device memory is not available, then `CUSPARSE_STATUS_ALLOC_FAILED` is returned.

Input

handle	handle to the CUSPARSE library context.
dirA	storage format of blocks, either <code>CUSPARSE_DIRECTION_ROW</code> or <code>CUSPARSE_DIRECTION_COLUMN</code> .
m	number of rows of sparse matrix <i>A</i> .
n	number of columns of sparse matrix <i>A</i> .
descrA	the descriptor of matrix <i>A</i> .
nnzA	number of nonz-zero elements of sparse matrix <i>A</i> .
csrValA	<type> array of <code>nnz (= csrRowPtrA(m) - csrRowPtrA(0))</code> non-zero elements of matrix <i>A</i> .
csrRowPtrA	integer array of <code>m+1</code> elements that contains the start of every row and the end of the last row plus one.
csrColIndA	integer array of <code>nnz (= csrRowPtrA(m) - csrRowPtrA(0))</code> column indices of the non-zero elements of matrix <i>A</i> .
blockDim	block dimension of sparse matrix <i>A</i> . The range of <i>blockDim</i> is between 1 and $\min(m, n)$.
descrC	the descriptor of matrix <i>C</i> .

Output

bsrValC	<type> array of <code>nnzb * blockDim²</code> non-zero elements of matrix <i>C</i> .
bsrRowPtrC	integer array of <code>mb+1</code> elements that contains the start of every block row and the end of the last block row plus one.

bsrColIndC	integer array of nnzb column indices of the non-zero blocks of matrix <i>C</i> .
nnzTotalDevHostPtr	total number of nonzero elements in device or host memory. It is equal to (bsrRowPtrC(mb)–bsrRowPtrC(0)).

Status Returned

CUSPARSE_STATUS_SUCCESS	the operation completed successfully.
CUSPARSE_STATUS_NOT_INITIALIZED	the library was not initialized.
CUSPARSE_STATUS_ALLOC_FAILED	the resources could not be allocated.
CUSPARSE_STATUS_INVALID_VALUE	invalid parameters were passed (m , n <0). IndexBase field of descrA , descrC is not base-0 or base-1, dirA is not row-major or column-major, or blockDim is not between 1 and min(m , n)).
CUSPARSE_STATUS_ARCH_MISMATCH	the device does not support double precision.
CUSPARSE_STATUS_INTERNAL_ERROR	an internal operation failed.
CUSPARSE_STATUS_MATRIX_TYPE_NOT_SUPPORTED	the matrix type is not supported.
CUSPARSE_STATUS_INTERNAL_ERROR	an internal operation failed.

11.5. `cusparse<t>csr2coo`

```
cusparseStatus_t
cusparseXcsr2coo(cusparseHandle_t handle, const int *csrRowPtr,
                 int nnz, int m, int *cooRowInd,
                 cusparseIndexBase_t idxBase)
```

This function converts the array containing the compressed row pointers (corresponding to CSR format) into an array of uncompressed row indices (corresponding to COO format).

It can also be used to convert the array containing the compressed column indices (corresponding to CSC format) into an array of uncompressed column indices (corresponding to COO format).

This function requires no extra storage. It is executed asynchronously with respect to the host and it may return control to the application on the host before the result is ready.

Input

handle	handle to the CUSPARSE library context.
csrRowPtr	integer array of m+1 elements that contains the start of every row and the end of the last row plus one.
nnz	number of non-zeros of the sparse matrix (that is also the length of array cooRowInd).
m	number of rows of matrix <i>A</i> .

<code>idxBase</code>	<code>CUSPARSE_INDEX_BASE_ZERO</code> or <code>CUSPARSE_INDEX_BASE_ONE</code> .
----------------------	--

Output

<code>cooRowInd</code>	integer array of <code>nnz</code> uncompressed row indices.
------------------------	---

Status Returned

<code>CUSPARSE_STATUS_SUCCESS</code>	the operation completed successfully.
<code>CUSPARSE_STATUS_NOT_INITIALIZED</code>	the library was not initialized.
<code>CUSPARSE_STATUS_INVALID_VALUE</code>	<code>IndexBase</code> is neither <code>CUSPARSE_INDEX_BASE_ZERO</code> nor <code>CUSPARSE_INDEX_BASE_ONE</code> .
<code>CUSPARSE_STATUS_EXECUTION_FAILED</code>	the function failed to launch on the GPU

11.6. `cusparse<t>csr2csc`

```

cusparseStatus_t
cusparseScsr2csc(cusparseHandle_t handle, int m, int n, int nnz,
                 const float *csrVal, const int *csrRowPtr,
                 const int *csrColInd, float *cscVal,
                 int *cscRowInd, int *cscColPtr,
                 cusparseAction_t copyValues,
                 cusparseIndexBase_t idxBase)

cusparseStatus_t
cusparseDcsr2csc(cusparseHandle_t handle, int m, int n, int nnz,
                 const double *csrVal, const int *csrRowPtr,
                 const int *csrColInd, double *cscVal,
                 int *cscRowInd, int *cscColPtr,
                 cusparseAction_t copyValues,
                 cusparseIndexBase_t idxBase)

cusparseStatus_t
cusparseCcsr2csc(cusparseHandle_t handle, int m, int n, int nnz,
                 const cuComplex *csrVal, const int *csrRowPtr,
                 const int *csrColInd, cuComplex *cscVal,
                 int *cscRowInd, int *cscColPtr,
                 cusparseAction_t copyValues,
                 cusparseIndexBase_t idxBase)

cusparseStatus_t
cusparseZcsr2csc(cusparseHandle_t handle, int m, int n, int nnz,
                 const cuDoubleComplex *csrVal, const int *csrRowPtr,
                 const int *csrColInd, cuDoubleComplex *cscVal,
                 int *cscRowInd, int *cscColPtr,
                 cusparseAction_t copyValues,
                 cusparseIndexBase_t idxBase)

```

This function converts a sparse matrix in CSR format (that is defined by the three arrays **csrVal**, **csrRowPtr** and **csrColInd**) into a sparse matrix in CSC format (that is defined by arrays **cscVal**, **cscRowInd**, and **cscColPtr**). The resulting matrix can also be seen as the transpose of the original sparse matrix. Notice that this routine can also be used to convert a matrix in CSC format into a matrix in CSR format.

This function requires significant amount of extra storage that is proportional to the matrix size. It is executed asynchronously with respect to the host and it may return control to the application on the host before the result is ready.

Input

handle	handle to the CUSPARSE library context.
m	number of rows of matrix <i>A</i> .
n	number of columns of matrix <i>A</i> .
nnz	number of nonz-zero elements of matrix <i>A</i> .
csrVal	<type> array of <code>nnz (= csrRowPtr(m) - csrRowPtr(0))</code> non-zero elements of matrix <i>A</i> .
csrRowPtr	integer array of <code>m+1</code> elements that contains the start of every row and the end of the last row plus one.
csrColInd	integer array of <code>nnz (= csrRowPtr(m) - csrRowPtr(0))</code> column indices of the non-zero elements of matrix <i>A</i> .
copyValues	CUSPARSE_ACTION_SYMBOLIC or CUSPARSE_ACTION_NUMERIC.
idxBase	CUSPARSE_INDEX_BASE_ZERO or CUSPARSE_INDEX_BASE_ONE.

Output

cscVal	<type> array of <code>nnz (= cscColPtr(n) - cscColPtr(0))</code> non-zero elements of matrix <i>A</i> . It is only filled-in if <code>copyValues</code> is set to CUSPARSE_ACTION_NUMERIC.
cscRowInd	integer array of <code>nnz (= cscColPtr(n) - cscColPtr(0))</code> column indices of the non-zero elements of matrix <i>A</i> .
cscColPtr	integer array of <code>n+1</code> elements that contains the start of every column and the end of the last column plus one.

Status Returned

CUSPARSE_STATUS_SUCCESS	the operation completed successfully.
CUSPARSE_STATUS_NOT_INITIALIZED	the library was not initialized.
CUSPARSE_STATUS_ALLOC_FAILED	the resources could not be allocated.
CUSPARSE_STATUS_INVALID_VALUE	invalid parameters were passed (<code>m, n, nnz < 0</code>).
CUSPARSE_STATUS_ARCH_MISMATCH	the device does not support double precision.
CUSPARSE_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU
CUSPARSE_STATUS_INTERNAL_ERROR	an internal operation failed.

11.7. `cusparse<t>csr2dense`

```

cusparseStatus_t
cusparseScsr2csc(cusparseHandle_t handle, int m, int n, int nnz,
                 const float *csrVal, const int *csrRowPtr,
                 const int *csrColInd, float *cscVal,
                 int *cscRowInd, int *cscColPtr,
                 cusparseAction_t copyValues,
                 cusparseIndexBase_t idxBase)

cusparseStatus_t
cusparseDcsr2csc(cusparseHandle_t handle, int m, int n, int nnz,
                 const double *csrVal, const int *csrRowPtr,
                 const int *csrColInd, double *cscVal,
                 int *cscRowInd, int *cscColPtr,
                 cusparseAction_t copyValues,
                 cusparseIndexBase_t idxBase)

cusparseStatus_t
cusparseCcsr2csc(cusparseHandle_t handle, int m, int n, int nnz,
                 const cuComplex *csrVal, const int *csrRowPtr,
                 const int *csrColInd, cuComplex *cscVal,
                 int *cscRowInd, int *cscColPtr,
                 cusparseAction_t copyValues,
                 cusparseIndexBase_t idxBase)

cusparseStatus_t
cusparseZcsr2csc(cusparseHandle_t handle, int m, int n, int nnz,
                 const cuDoubleComplex *csrVal, const int *csrRowPtr,
                 const int *csrColInd, cuDoubleComplex *cscVal,
                 int *cscRowInd, int *cscColPtr,
                 cusparseAction_t copyValues,
                 cusparseIndexBase_t idxBase)

```

This function converts the sparse matrix in CSR format (that is defined by the three arrays **csrValA**, **csrRowPtrA** and **csrColIndA**) into the matrix **A** in dense format. The dense matrix **A** is filled in with the values of the sparse matrix and with zeros elsewhere.

This function requires no extra storage. It is executed asynchronously with respect to the host and it may return control to the application on the host before the result is ready.

Input

handle	handle to the CUSPARSE library context.
m	number of rows of matrix <i>A</i> .
n	number of columns of matrix <i>A</i> .
descrA	the descriptor of matrix <i>A</i> . The supported matrix type is <code>CUSPARSE_MATRIX_TYPE_GENERAL</code> . Also, the supported index bases are <code>CUSPARSE_INDEX_BASE_ZERO</code> and <code>CUSPARSE_INDEX_BASE_ONE</code> .
csrValA	<type> array of <code>nnz (= csrRowPtrA(m) - csrRowPtrA(0))</code> non-zero elements of matrix <i>A</i> .
csrRowPtrA	integer array of <code>m+1</code> elements that contains the start of every row and the end of the last row plus one.

csrColIndA	integer array of nnz (= csrRowPtrA (m) - csrRowPtrA (0)) column indices of the non-zero elements of matrix <i>A</i> .
lda	leading dimension of array matrix <i>A</i> .

Output

A	array of dimensions (lda , n) that is filled in with the values of the sparse matrix.
----------	---

Status Returned

CUSPARSE_STATUS_SUCCESS	the operation completed successfully.
CUSPARSE_STATUS_NOT_INITIALIZED	the library was not initialized.
CUSPARSE_STATUS_INVALID_VALUE	invalid parameters were passed (m , n <0).
CUSPARSE_STATUS_ARCH_MISMATCH	the device does not support double precision.
CUSPARSE_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU
CUSPARSE_STATUS_MATRIX_TYPE_NOT_SUPPORTED	the matrix type is not supported.

11.8. `cusparse<t>csr2hyb`

```

cusparseStatus_t
cusparseScsr2hyb(cusparseHandle_t handle, int m, int n,
                 const cusparseMatDescr_t descrA,
                 const float *csrValA,
                 const int *csrRowPtrA, const int *csrColIndA,
                 cusparseHybMat_t hybA, int userEllWidth,
                 cusparseHybPartition_t partitionType)
cusparseStatus_t
cusparseDcsr2hyb(cusparseHandle_t handle, int m, int n,
                 const cusparseMatDescr_t descrA,
                 const double *csrValA,
                 const int *csrRowPtrA, const int *csrColIndA,
                 cusparseHybMat_t hybA, int userEllWidth,
                 cusparseHybPartition_t partitionType)
cusparseStatus_t
cusparseCcsr2hyb(cusparseHandle_t handle, int m, int n,
                 const cusparseMatDescr_t descrA,
                 const cuComplex *csrValA,
                 const int *csrRowPtrA, const int *csrColIndA,
                 cusparseHybMat_t hybA, int userEllWidth,
                 cusparseHybPartition_t partitionType)
cusparseStatus_t
cusparseZcsr2hyb(cusparseHandle_t handle, int m, int n,
                 const cusparseMatDescr_t descrA,
                 const cuDoubleComplex *csrValA,
                 const int *csrRowPtrA, const int *csrColIndA,
                 cusparseHybMat_t hybA, int userEllWidth,
                 cusparseHybPartition_t partitionType)

```


This function converts a sparse matrix in CSR format into a sparse matrix in HYB format. It assumes that the **hybA** parameter has been initialized with **cusparseCreateHybMat** routine before calling this function.

This function requires some amount of temporary storage and a significant amount of storage for the matrix in HYB format. It is executed asynchronously with respect to the host and it may return control to the application on the host before the result is ready.

Input

handle	handle to the CUSPARSE library context.
m	number of rows of matrix <i>A</i> .
n	number of columns of matrix <i>A</i> .
descrA	the descriptor of matrix <i>A</i> . The supported matrix type is CUSPARSE_MATRIX_TYPE_GENERAL . Also, the supported index bases are CUSPARSE_INDEX_BASE_ZERO and CUSPARSE_INDEX_BASE_ONE .
csrValA	<type> array of nnz (= csrRowPtrA (m) - csrRowPtrA (0)) non-zero elements of matrix <i>A</i> .
csrRowPtrA	integer array of m+1 elements that contains the start of every row and the end of the last row plus one.
csrColIndA	integer array of nnz (= csrRowPtrA (m) - csrRowPtrA (0)) column indices of the non-zero elements of matrix <i>A</i> .
userEllWidth	width of the regular (ELL) part of the matrix in HYB format, which should be less than maximum number of non-zeros per row and is only required if partitionType == CUSPARSE_HYB_PARTITION_USER .
partitionType	partitioning method to be used in the conversion (please refer to cusparseHybPartition_t on page ?? for details).

Output

hybA	the matrix A in HYB storage format.
-------------	--

Status Returned

CUSPARSE_STATUS_SUCCESS	the operation completed successfully.
CUSPARSE_STATUS_NOT_INITIALIZED	the library was not initialized.
CUSPARSE_STATUS_ALLOC_FAILED	the resources could not be allocated.
CUSPARSE_STATUS_INVALID_VALUE	invalid parameters were passed (m , n < 0).
CUSPARSE_STATUS_ARCH_MISMATCH	the device does not support double precision.
CUSPARSE_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU
CUSPARSE_STATUS_INTERNAL_ERROR	an internal operation failed.
CUSPARSE_STATUS_MATRIX_TYPE_NOT_SUPPORTED	

the matrix type is not supported.

11.9. `cusparse<t>dense2csc`

```
cusparseStatus_t
cusparseSdense2csc(cusparseHandle_t handle, int m, int n,
                  const cusparseMatDescr_t descrA,
                  const float *A,
                  int lda, const int *nnzPerCol,
                  float *cscValA,
                  int *cscRowIndA, int *cscColPtrA)

cusparseStatus_t
cusparseDdense2csc(cusparseHandle_t handle, int m, int n,
                  const cusparseMatDescr_t descrA,
                  const double *A,
                  int lda, const int *nnzPerCol,
                  double *cscValA,
                  int *cscRowIndA, int *cscColPtrA)

cusparseStatus_t
cusparseCdense2csc(cusparseHandle_t handle, int m, int n,
                  const cusparseMatDescr_t descrA,
                  const cuComplex *A,
                  int lda, const int *nnzPerCol,
                  cuComplex *cscValA,
                  int *cscRowIndA, int *cscColPtrA)

cusparseStatus_t
cusparseZdense2csc(cusparseHandle_t handle, int m, int n,
                  const cusparseMatDescr_t descrA,
                  const cuDoubleComplex *A,
                  int lda, const int *nnzPerCol,
                  cuDoubleComplex *cscValA,
                  int *cscRowIndA, int *cscColPtrA)
```

This function converts the matrix **A** in dense format into a sparse matrix in CSC format. All the parameters are assumed to have been pre-allocated by the user and the arrays are filled in based on **nnzPerCol**, which can be pre-computed with **`cusparse<t>nnz()`**.

This function requires no extra storage. It is executed asynchronously with respect to the host and it may return control to the application on the host before the result is ready.

Input

handle	handle to the CUSPARSE library context.
m	number of rows of matrix <i>A</i> .
n	number of columns of matrix <i>A</i> .
descrA	the descriptor of matrix <i>A</i> . The supported matrix type is CUSPARSE_MATRIX_TYPE_GENERAL . Also, the supported index bases are CUSPARSE_INDEX_BASE_ZERO and CUSPARSE_INDEX_BASE_ONE .
A	array of dimensions (<i>lda</i> , <i>n</i>).
lda	leading dimension of dense array <i>A</i> .
nnzPerCol	array of size <i>n</i> containing the number of non-zero elements per column.

Output

cscValA	<type> array of nnz (= cscRowPtrA (m) - cscRowPtrA (0)) non-zero elements of matrix <i>A</i> . It is only filled-in if copyValues is set to CUSPARSE_ACTION_NUMERIC .
cscRowIndA	integer array of nnz (= cscRowPtrA (m) - cscRowPtrA (0)) row indices of the non-zero elements of matrix <i>A</i> .
cscColPtrA	integer array of n+1 elements that contains the start of every column and the end of the last column plus one.

Status Returned

CUSPARSE_STATUS_SUCCESS	the operation completed successfully.
CUSPARSE_STATUS_NOT_INITIALIZED	the library was not initialized.
CUSPARSE_STATUS_INVALID_VALUE	invalid parameters were passed (m , n <0).
CUSPARSE_STATUS_ARCH_MISMATCH	the device does not support double precision.
CUSPARSE_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU
CUSPARSE_STATUS_MATRIX_TYPE_NOT_SUPPORTED	the matrix type is not supported.

11.10. cusparse<t>dense2csr

```

cusparseStatus_t
cusparseSdense2csr(cusparseHandle_t handle, int m, int n,
                  const cusparseMatDescr_t descrA,
                  const float *A,
                  int lda, const int *nnzPerRow,
                  float *csrValA,
                  int *csrRowPtrA, int *csrColIndA)

cusparseStatus_t
cusparseDdense2csr(cusparseHandle_t handle, int m, int n,
                  const cusparseMatDescr_t descrA,
                  const double *A,
                  int lda, const int *nnzPerRow,
                  double *csrValA,
                  int *csrRowPtrA, int *csrColIndA)

cusparseStatus_t
cusparseCdense2csr(cusparseHandle_t handle, int m, int n,
                  const cusparseMatDescr_t descrA,
                  const cuComplex *A,
                  int lda, const int *nnzPerRow,
                  cuComplex *csrValA,
                  int *csrRowPtrA, int *csrColIndA)

cusparseStatus_t
cusparseZdense2csr(cusparseHandle_t handle, int m, int n,
                  const cusparseMatDescr_t descrA,
                  const cuDoubleComplex *A,
                  int lda, const int *nnzPerRow,
                  cuDoubleComplex *csrValA,
                  int *csrRowPtrA, int *csrColIndA)

```

This function converts the matrix **A** in dense format into a sparse matrix in CSR format. All the parameters are assumed to have been pre-allocated by the user and the arrays are filled in based on the **nnzPerRow**, which can be pre-computed with **cusparselt>nnz()**.

This function requires no extra storage. It is executed asynchronously with respect to the host and it may return control to the application on the host before the result is ready.

Input

handle	handle to the CUSPARSE library context.
m	number of rows of matrix <i>A</i> .
n	number of columns of matrix <i>A</i> .
descrA	the descriptor of matrix <i>A</i> . The supported matrix type is CUSPARSE_MATRIX_TYPE_GENERAL . Also, the supported index bases are CUSPARSE_INDEX_BASE_ZERO and CUSPARSE_INDEX_BASE_ONE .
A	array of dimensions (lda , n).
lda	leading dimension of dense array <i>A</i> .
nnzPerRow	array of size n containing the number of non-zero elements per row.

Output

csrValA	<type> array of nnz (= csrRowPtrA (m) - csrRowPtrA (0)) non-zero elements of matrix <i>A</i> .
csrRowPtrA	integer array of m+1 elements that contains the start of every column and the end of the last column plus one.
csrColIndA	integer array of nnz (= csrRowPtrA (m) - csrRowPtrA (0)) column indices of the non-zero elements of matrix <i>A</i> .

Status Returned

CUSPARSE_STATUS_SUCCESS	the operation completed successfully.
CUSPARSE_STATUS_NOT_INITIALIZED	the library was not initialized.
CUSPARSE_STATUS_ALLOC_FAILED	the resources could not be allocated.
CUSPARSE_STATUS_INVALID_VALUE	invalid parameters were passed (m , n < 0).
CUSPARSE_STATUS_ARCH_MISMATCH	the device does not support double precision.
CUSPARSE_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU
CUSPARSE_STATUS_MATRIX_TYPE_NOT_SUPPORTED	the matrix type is not supported.

11.11. `cusparse<t>dense2hyb`

```

cusparseStatus_t
cusparseSdense2hyb(cusparseHandle_t handle, int m, int n,
                  const cusparseMatDescr_t descrA,
                  const float *A,
                  int lda, const int *nnzPerRow, cusparseHybMat_t hybA,
                  int userEllWidth,
                  cusparseHybPartition_t partitionType)

cusparseStatus_t
cusparseDdense2hyb(cusparseHandle_t handle, int m, int n,
                  const cusparseMatDescr_t descrA,
                  const double *A,
                  int lda, const int *nnzPerRow, cusparseHybMat_t
                  hybA,
                  int userEllWidth,
                  cusparseHybPartition_t partitionType)

cusparseStatus_t
cusparseCdense2hyb(cusparseHandle_t handle, int m, int n,
                  const cusparseMatDescr_t descrA,
                  const cuComplex *A,
                  int lda, const int *nnzPerRow, cusparseHybMat_t hybA,
                  int userEllWidth,
                  cusparseHybPartition_t partitionType)

cusparseStatus_t
cusparseZdense2hyb(cusparseHandle_t handle, int m, int n,
                  const cusparseMatDescr_t descrA,
                  const cuDoubleComplex *A,
                  int lda, const int *nnzPerRow, cusparseHybMat_t hybA,
                  int userEllWidth,
                  cusparseHybPartition_t partitionType)

```

This function converts the matrix **A** in dense format into a sparse matrix in HYB format. It assumes that the routine **`cusparseCreateHybMat`** was used to initialize the opaque structure **`hybA`** and that the array **`nnzPerRow`** was pre-computed with **`cusparse<t>nnz()`**.

This function requires some amount of temporary storage and a significant amount of storage for the matrix in HYB format. It is executed asynchronously with respect to the host and it may return control to the application on the host before the result is ready.

Input

handle	handle to the CUSPARSE library context.
m	number of rows of matrix A .
n	number of columns of matrix A .
descrA	the descriptor of matrix A . The supported matrix type is <code>CUSPARSE_MATRIX_TYPE_GENERAL</code> .
A	array of dimensions (lda, n) .
lda	leading dimension of dense array A .
nnzPerRow	array of size m containing the number of non-zero elements per row.

userEllWidth	width of the regular (ELL) part of the matrix in HYB format, which should be less than maximum number of non-zeros per row and is only required if partitionType == CUSPARSE_HYB_PARTITION_USER .
partitionType	partitioning method to be used in the conversion (please refer to <code>cusparseHybPartition_t</code> on page ?? for details).

Output

hybA	the matrix A in HYB storage format.
-------------	--

Status Returned

CUSPARSE_STATUS_SUCCESS	the operation completed successfully.
CUSPARSE_STATUS_NOT_INITIALIZED	the library was not initialized.
CUSPARSE_STATUS_ALLOC_FAILED	the resources could not be allocated.
CUSPARSE_STATUS_INVALID_VALUE	invalid parameters were passed (m , n <0).
CUSPARSE_STATUS_ARCH_MISMATCH	the device does not support double precision.
CUSPARSE_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU
CUSPARSE_STATUS_INTERNAL_ERROR	an internal operation failed.
CUSPARSE_STATUS_MATRIX_TYPE_NOT_SUPPORTED	the matrix type is not supported.

11.12. `cusparse<t>hyb2csr`

```

cusparseStatus_t
cusparseShyb2csr(cusparseHandle_t handle,
                 const cusparseMatDescr_t descrA,
                 const cusparseHybMat_t hybA,
                 float *csrValA, int *csrRowPtrA, int *csrColIndA)

cusparseStatus_t
cusparseDhyb2csr(cusparseHandle_t handle,
                 const cusparseMatDescr_t descrA,
                 const cusparseHybMat_t hybA,
                 double *csrValA, int *csrRowPtrA, int *csrColIndA)

cusparseStatus_t
cusparseChyb2csr(cusparseHandle_t handle,
                 const cusparseMatDescr_t descrA,
                 const cusparseHybMat_t hybA,
                 cuComplex *csrValA, int *csrRowPtrA, int *csrColIndA)

cusparseStatus_t
cusparseZhyb2csr(cusparseHandle_t handle,
                 const cusparseMatDescr_t descrA,
                 const cusparseHybMat_t hybA,
                 cuDoubleComplex *csrValA, int *csrRowPtrA, int
                 *csrColIndA)

```

This function converts a sparse matrix in HYB format into a sparse matrix in CSR format.

This function requires some amount of temporary storage. It is executed asynchronously with respect to the host and it may return control to the application on the host before the result is ready.

Input

handle	handle to the CUSPARSE library context.
descrA	the descriptor of matrix <i>A</i> in Hyb format. The supported matrix type is <code>CUSPARSE_MATRIX_TYPE_GENERAL</code> .
hybA	the matrix <i>A</i> in HYB storage format.

Output

csrValA	<type> array of <code>nnz (= csrRowPtrA(m) - csrRowPtrA(0))</code> non-zero elements of matrix <i>A</i> .
csrRowPtrA	integer array of <code>m+1</code> elements that contains the start of every column and the end of the last row plus one.
csrColIndA	integer array of <code>nnz (= csrRowPtrA(m) - csrRowPtrA(0))</code> column indices of the non-zero elements of matrix <i>A</i> .

Status Returned

<code>CUSPARSE_STATUS_SUCCESS</code>	the operation completed successfully.
<code>CUSPARSE_STATUS_NOT_INITIALIZED</code>	the library was not initialized.
<code>CUSPARSE_STATUS_ALLOC_FAILED</code>	the resources could not be allocated.
<code>CUSPARSE_STATUS_INVALID_VALUE</code>	invalid parameters were passed (<code>m, n < 0</code>).
<code>CUSPARSE_STATUS_ARCH_MISMATCH</code>	the device does not support double precision.
<code>CUSPARSE_STATUS_EXECUTION_FAILED</code>	the function failed to launch on the GPU
<code>CUSPARSE_STATUS_INTERNAL_ERROR</code>	an internal operation failed.
<code>CUSPARSE_STATUS_MATRIX_TYPE_NOT_SUPPORTED</code>	the matrix type is not supported.

11.13. `cusparse<t>hyb2dense`

```

cusparseStatus_t
cusparseShyb2csr(cusparseHandle_t handle,
                 const cusparseMatDescr_t descrA,
                 const cusparseHybMat_t hybA,
                 float *csrValA, int *csrRowPtrA, int *csrColIndA)

cusparseStatus_t
cusparseDhyb2csr(cusparseHandle_t handle,
                 const cusparseMatDescr_t descrA,
                 const cusparseHybMat_t hybA,
                 double *csrValA, int *csrRowPtrA, int *csrColIndA)

cusparseStatus_t
cusparseChyb2csr(cusparseHandle_t handle,
```

```

        const cusparseMatDescr_t descrA,
        const cusparseHybMat_t hybA,
        cuComplex *csrValA, int *csrRowPtrA, int *csrColIndA)
cusparseStatus_t
cusparseZhyb2csr(cusparseHandle_t handle,
        const cusparseMatDescr_t descrA,
        const cusparseHybMat_t hybA,
        cuDoubleComplex *csrValA, int *csrRowPtrA, int
        *csrColIndA)

```

This function converts a sparse matrix in HYB format (contained in the opaque structure) into a matrix **A** in dense format. The dense matrix **A** is filled in with the values of the sparse matrix and with zeros elsewhere.

This function requires no extra storage. It is executed asynchronously with respect to the host and it may return control to the application on the host before the result is ready.

Input

handle	handle to the CUSPARSE library context.
descrA	the descriptor of matrix A in Hyb format. The supported matrix type is CUSPARSE_MATRIX_TYPE_GENERAL .
hybA	the matrix A in HYB storage format.
lda	leading dimension of dense array A .

Output

A	array of dimensions (lda , n) that is filled in with the values of the sparse matrix.
----------	---

Status Returned

CUSPARSE_STATUS_SUCCESS	the operation completed successfully.
CUSPARSE_STATUS_NOT_INITIALIZED	the library was not initialized.
CUSPARSE_STATUS_INVALID_VALUE	the internally stored hyb format parameters are invalid.
CUSPARSE_STATUS_ARCH_MISMATCH	the device does not support double precision.
CUSPARSE_STATUS_EXECUTION_FAILED	the function failed to launch on the GPU
CUSPARSE_STATUS_MATRIX_TYPE_NOT_SUPPORTED	the matrix type is not supported.

11.14. cusparse<t>nnz

```

cusparseStatus_t
cusparseSnnz(cusparseHandle_t handle, cusparseDirection_t dirA, int m,
        int n, const cusparseMatDescr_t descrA,
        const float *A,
        int lda, int *nnzPerRowColumn, int *nnzTotalDevHostPtr)
cusparseStatus_t
cusparseDnnz(cusparseHandle_t handle, cusparseDirection_t dirA, int m,
        int n, const cusparseMatDescr_t descrA,

```



```

        const double          *A,
        int lda, int *nnzPerRowColumn, int *nnzTotalDevHostPtr)
cusparseStatus_t
cusparseCnnz(cusparseHandle_t handle, cusparseDirection_t dirA, int m,
             int n, const cusparseMatDescr_t descrA,
             const cuComplex   *A,
             int lda, int *nnzPerRowColumn, int *nnzTotalDevHostPtr)
cusparseStatus_t
cusparseZnnz(cusparseHandle_t handle, cusparseDirection_t dirA, int m,
             int n, const cusparseMatDescr_t descrA,
             const cuDoubleComplex *A,
             int lda, int *nnzPerRowColumn, int *nnzTotalDevHostPtr)

```

This function computes the number of non-zero elements per row or column and the total number of non-zero elements in a dense matrix.

This function requires no extra storage. It is executed asynchronously with respect to the host and it may return control to the application on the host before the result is ready.

Input

handle	handle to the CUSPARSE library context.
dirA	direction that specifies whether to count non-zero elements by <code>CUSPARSE_DIRECTION_ROW</code> or <code>CUSPARSE_DIRECTION_COLUMN</code> .
m	number of rows of matrix <i>A</i> .
n	number of columns of matrix <i>A</i> .
descrA	the descriptor of matrix <i>A</i> in Hyb format. The supported matrix type is <code>CUSPARSE_MATRIX_TYPE_GENERAL</code> . Also, the supported index bases are <code>CUSPARSE_INDEX_BASE_ZERO</code> and <code>CUSPARSE_INDEX_BASE_ONE</code> .
A	array of dimensions <code>(lda, n)</code> .
lda	leading dimension of dense array <i>A</i> .

Output

nnzPerRowColumn	array of size <i>m</i> or <i>n</i> containing the number of non-zero elements per row or column, respectively.
nnzTotalDevHostPtr	total number of non-zero elements in device or host memory.

Status Returned

<code>CUSPARSE_STATUS_SUCCESS</code>	the operation completed successfully.
<code>CUSPARSE_STATUS_NOT_INITIALIZED</code>	the library was not initialized.
<code>CUSPARSE_STATUS_ALLOC_FAILED</code>	the resources could not be allocated.
<code>CUSPARSE_STATUS_INVALID_VALUE</code>	invalid parameters were passed (<i>m</i> , <i>n</i> < 0).
<code>CUSPARSE_STATUS_ARCH_MISMATCH</code>	the device does not support double precision.
<code>CUSPARSE_STATUS_EXECUTION_FAILED</code>	the function failed to launch on the GPU

CUSPARSE_STATUS_INTERNAL_ERROR	an internal operation failed.
CUSPARSE_STATUS_MATRIX_TYPE_NOT_SUPPORTED	the matrix type is not supported.

Chapter 12.

APPENDIX A: CUSPARSE LIBRARY C++ EXAMPLE

For sample code reference please see the example code below. It shows an application written in C++ using the CUSPARSE library API. The code performs the following actions:

1. Creates a sparse test matrix in COO format.
2. Creates a sparse and dense vector.
3. Allocates GPU memory and copies the matrix and vectors into it.
4. Initializes the CUSPARSE library.
5. Creates and sets up the matrix descriptor.
6. Converts the matrix from COO to CSR format.
7. Exercises Level 1 routines.
8. Exercises Level 2 routines.
9. Exercises Level 3 routines.
10. Destroys the matrix descriptor.
11. Releases resources allocated for the CUSPARSE library.

```
//Example: Application using C++ and the CUSPARSE library
//-----
#include <stdio.h>
#include <stdlib.h>
#include <cuda_runtime.h>
#include "cusparse_v2.h"

#define CLEANUP(s) \
do { \
    printf ("%s\n", s); \
    if (yHostPtr) free(yHostPtr); \
    if (zHostPtr) free(zHostPtr); \
    if (xIndHostPtr) free(xIndHostPtr); \
    if (xValHostPtr) free(xValHostPtr); \
    if (cooRowIndexHostPtr) free(cooRowIndexHostPtr); \
    if (cooColIndexHostPtr) free(cooColIndexHostPtr); \
    if (cooValHostPtr) free(cooValHostPtr); \
}
```

```

    if (y)                cudaFree(y);                \
    if (z)                cudaFree(z);                \
    if (xInd)             cudaFree(xInd);             \
    if (xVal)             cudaFree(xVal);             \
    if (csrRowPtr)        cudaFree(csrRowPtr);        \
    if (cooRowIndex)     cudaFree(cooRowIndex);       \
    if (cooColIndex)     cudaFree(cooColIndex);       \
    if (cooVal)          cudaFree(cooVal);           \
    if (descr)           cusparseDestroyMatDescr(descr);\
    if (handle)          cusparseDestroy(handle);     \
    cudaDeviceReset();   \
    fflush(stdout);      \
} while (0)

int main(){
    cudaError_t cudaStat1,cudaStat2,cudaStat3,cudaStat4,cudaStat5,cudaStat6;
    cusparseStatus_t status;
    cusparseHandle_t handle=0;
    cusparseMatDescr_t descr=0;
    int *    cooRowIndexHostPtr=0;
    int *    cooColIndexHostPtr=0;
    double * cooValHostPtr=0;
    int *    cooRowIndex=0;
    int *    cooColIndex=0;
    double * cooVal=0;
    int *    xIndHostPtr=0;
    double * xValHostPtr=0;
    double * yHostPtr=0;
    int *    xInd=0;
    double * xVal=0;
    double * y=0;
    int *    csrRowPtr=0;
    double * zHostPtr=0;
    double * z=0;
    int      n, nnz, nnz_vector;
    double dzero =0.0;
    double dtwo  =2.0;
    double dthree=3.0;
    double dfive =5.0;

    printf("testing example\n");
    /* create the following sparse test matrix in COO format */
    /* |1.0    2.0 3.0|
       |    4.0    |
       |5.0    6.0 7.0|
       |    8.0    9.0| */
    n=4; nnz=9;
    cooRowIndexHostPtr = (int *)    malloc(nnz*sizeof(cooRowIndexHostPtr[0]));
    cooColIndexHostPtr = (int *)    malloc(nnz*sizeof(cooColIndexHostPtr[0]));
    cooValHostPtr      = (double *) malloc(nnz*sizeof(cooValHostPtr[0]));
    if ((!cooRowIndexHostPtr) || (!cooColIndexHostPtr) || (!cooValHostPtr)){
        CLEANUP("Host malloc failed (matrix)");
        return 1;
    }
    cooRowIndexHostPtr[0]=0; cooColIndexHostPtr[0]=0; cooValHostPtr[0]=1.0;
    cooRowIndexHostPtr[1]=0; cooColIndexHostPtr[1]=2; cooValHostPtr[1]=2.0;
    cooRowIndexHostPtr[2]=0; cooColIndexHostPtr[2]=3; cooValHostPtr[2]=3.0;
    cooRowIndexHostPtr[3]=1; cooColIndexHostPtr[3]=1; cooValHostPtr[3]=4.0;
    cooRowIndexHostPtr[4]=2; cooColIndexHostPtr[4]=0; cooValHostPtr[4]=5.0;
    cooRowIndexHostPtr[5]=2; cooColIndexHostPtr[5]=2; cooValHostPtr[5]=6.0;
    cooRowIndexHostPtr[6]=2; cooColIndexHostPtr[6]=3; cooValHostPtr[6]=7.0;
    cooRowIndexHostPtr[7]=3; cooColIndexHostPtr[7]=1; cooValHostPtr[7]=8.0;
    cooRowIndexHostPtr[8]=3; cooColIndexHostPtr[8]=3; cooValHostPtr[8]=9.0;
    /*
    //print the matrix
    printf("Input data:\n");
    for (int i=0; i<nnz; i++){

```

```

        printf("cooRowIndexHostPtr[%d]=%d  ", i, cooRowIndexHostPtr[i]);
        printf("cooColIndexHostPtr[%d]=%d  ", i, cooColIndexHostPtr[i]);
        printf("cooValHostPtr[%d]=%f      \n", i, cooValHostPtr[i]);
    }
    */

    /* create a sparse and dense vector */
    /* xVal= [100.0 200.0 400.0]      (sparse)
       xInd= [0      1      3      ]
       y   = [10.0 20.0 30.0 40.0 | 50.0 60.0 70.0 80.0] (dense) */
    nnz_vector = 3;
    xIndHostPtr = (int *) malloc(nnz_vector*sizeof(xIndHostPtr[0]));
    xValHostPtr = (double *) malloc(nnz_vector*sizeof(xValHostPtr[0]));
    yHostPtr = (double *) malloc(2*n*sizeof(yHostPtr[0]));
    zHostPtr = (double *) malloc(2*(n+1)*sizeof(zHostPtr[0]));
    if((!xIndHostPtr) || (!xValHostPtr) || (!yHostPtr) || (!zHostPtr)){
        CLEANUP("Host malloc failed (vectors)");
        return 1;
    }
    yHostPtr[0] = 10.0; xIndHostPtr[0]=0; xValHostPtr[0]=100.0;
    yHostPtr[1] = 20.0; xIndHostPtr[1]=1; xValHostPtr[1]=200.0;
    yHostPtr[2] = 30.0;
    yHostPtr[3] = 40.0; xIndHostPtr[2]=3; xValHostPtr[2]=400.0;
    yHostPtr[4] = 50.0;
    yHostPtr[5] = 60.0;
    yHostPtr[6] = 70.0;
    yHostPtr[7] = 80.0;
    /*
    //print the vectors
    for (int j=0; j<2; j++){
        for (int i=0; i<n; i++){
            printf("yHostPtr[%d,%d]=%f\n", i, j, yHostPtr[i+n*j]);
        }
    }
    for (int i=0; i<nnz_vector; i++){
        printf("xIndHostPtr[%d]=%d  ", i, xIndHostPtr[i]);
        printf("xValHostPtr[%d]=%f\n", i, xValHostPtr[i]);
    }
    */

    /* allocate GPU memory and copy the matrix and vectors into it */
    cudaStat1 = cudaMalloc((void**) &cooRowIndex, nnz*sizeof(cooRowIndex[0]));
    cudaStat2 = cudaMalloc((void**) &cooColIndex, nnz*sizeof(cooColIndex[0]));
    cudaStat3 = cudaMalloc((void**) &cooVal, nnz*sizeof(cooVal[0]));
    cudaStat4 = cudaMalloc((void**) &y, 2*n*sizeof(y[0]));
    cudaStat5 = cudaMalloc((void**) &xInd, nnz_vector*sizeof(xInd[0]));
    cudaStat6 = cudaMalloc((void**) &xVal, nnz_vector*sizeof(xVal[0]));
    if ((cudaStat1 != cudaSuccess) ||
        (cudaStat2 != cudaSuccess) ||
        (cudaStat3 != cudaSuccess) ||
        (cudaStat4 != cudaSuccess) ||
        (cudaStat5 != cudaSuccess) ||
        (cudaStat6 != cudaSuccess)) {
        CLEANUP("Device malloc failed");
        return 1;
    }
    cudaStat1 = cudaMemcpy(cooRowIndex, cooRowIndexHostPtr,
        (size_t) (nnz*sizeof(cooRowIndex[0])),
        cudaMemcpyHostToDevice);
    cudaStat2 = cudaMemcpy(cooColIndex, cooColIndexHostPtr,
        (size_t) (nnz*sizeof(cooColIndex[0])),
        cudaMemcpyHostToDevice);
    cudaStat3 = cudaMemcpy(cooVal, cooValHostPtr,
        (size_t) (nnz*sizeof(cooVal[0])),
        cudaMemcpyHostToDevice);
    cudaStat4 = cudaMemcpy(y, yHostPtr,
        (size_t) (2*n*sizeof(y[0])),

```

```

        cudaMemcpyHostToDevice);
cudaStat5 = cudaMemcpy(xInd, xIndHostPtr,
    (size_t)(nnz_vector*sizeof(xInd[0])),
    cudaMemcpyHostToDevice);
cudaStat6 = cudaMemcpy(xVal, xValHostPtr,
    (size_t)(nnz_vector*sizeof(xVal[0])),
    cudaMemcpyHostToDevice);
if ((cudaStat1 != cudaSuccess) ||
    (cudaStat2 != cudaSuccess) ||
    (cudaStat3 != cudaSuccess) ||
    (cudaStat4 != cudaSuccess) ||
    (cudaStat5 != cudaSuccess) ||
    (cudaStat6 != cudaSuccess)) {
    CLEANUP("Memcpy from Host to Device failed");
    return 1;
}

/* initialize cusparse library */
status= cusparseCreate(&handle);
if (status != CUSPARSE_STATUS_SUCCESS) {
    CLEANUP("CUSPARSE Library initialization failed");
    return 1;
}

/* create and setup matrix descriptor */
status= cusparseCreateMatDescr(&descr);
if (status != CUSPARSE_STATUS_SUCCESS) {
    CLEANUP("Matrix descriptor initialization failed");
    return 1;
}
cusparseSetMatType(descr,CUSPARSE_MATRIX_TYPE_GENERAL);
cusparseSetMatIndexBase(descr,CUSPARSE_INDEX_BASE_ZERO);

/* exercise conversion routines (convert matrix from COO 2 CSR format) */
cudaStat1 = cudaMalloc((void**) &csrRowPtr, (n+1)*sizeof(csrRowPtr[0]));
if (cudaStat1 != cudaSuccess) {
    CLEANUP("Device malloc failed (csrRowPtr)");
    return 1;
}
status= cusparseXcoo2csr(handle,cooRowIndex,nnz,n,
    csrRowPtr,CUSPARSE_INDEX_BASE_ZERO);
if (status != CUSPARSE_STATUS_SUCCESS) {
    CLEANUP("Conversion from COO to CSR format failed");
    return 1;
}
//csrRowPtr = [0 3 4 7 9]

/* exercise Level 1 routines (scatter vector elements) */
status= cusparseDsctr(handle, nnz_vector, xVal, xInd,
    &y[n], CUSPARSE_INDEX_BASE_ZERO);
if (status != CUSPARSE_STATUS_SUCCESS) {
    CLEANUP("Scatter from sparse to dense vector failed");
    return 1;
}
//y = [10 20 30 40 | 100 200 70 400]

/* exercise Level 2 routines (csrcmv) */
status= cusparseDcsrcmv(handle,CUSPARSE_OPERATION_NON_TRANSPOSE, n, n, nnz,
    &dtwo, descr, cooVal, csrRowPtr, cooColIndex,
    &y[0], &dthree, &y[n]);
if (status != CUSPARSE_STATUS_SUCCESS) {
    CLEANUP("Matrix-vector multiplication failed");
    return 1;
}
//y = [10 20 30 40 | 680 760 1230 2240]
cudaMemcpy(yHostPtr, y, (size_t)(2*n*sizeof(y[0])), cudaMemcpyDeviceToHost);
/*

```

```

printf("Intermediate results:\n");
for (int j=0; j<2; j++){
    for (int i=0; i<n; i++){
        printf("yHostPtr[%d,%d]=%f\n",i,j,yHostPtr[i+n*j]);
    }
}
*/

/* exercise Level 3 routines (csrmm) */
cudaStat1 = cudaMalloc((void**) &z, 2*(n+1)*sizeof(z[0]));
if (cudaStat1 != cudaSuccess) {
    CLEANUP("Device malloc failed (z)");
    return 1;
}
cudaStat1 = cudaMemset((void *)z,0, 2*(n+1)*sizeof(z[0]));
if (cudaStat1 != cudaSuccess) {
    CLEANUP("Memset on Device failed");
    return 1;
}
status= cusparseDcsrmm(handle, CUSPARSE_OPERATION_NON_TRANSPOSE, n, 2, n,
                        nnz, &dfive, descr, cooVal, csrRowPtr, cooColIndex,
                        y, n, &dzero, z, n+1);
if (status != CUSPARSE_STATUS_SUCCESS) {
    CLEANUP("Matrix-matrix multiplication failed");
    return 1;
}

/* print final results (z) */
cudaStat1 = cudaMemcpy(zHostPtr, z,
                        (size_t)(2*(n+1)*sizeof(z[0])),
                        cudaMemcpyDeviceToHost);
if (cudaStat1 != cudaSuccess) {
    CLEANUP("Memcpy from Device to Host failed");
    return 1;
}
//z = [950 400 2550 2600 0 | 49300 15200 132300 131200 0]
/*
printf("Final results:\n");
for (int j=0; j<2; j++){
    for (int i=0; i<n+1; i++){
        printf("z[%d,%d]=%f\n",i,j,zHostPtr[i+(n+1)*j]);
    }
}
*/

/* destroy matrix descriptor */
status = cusparseDestroyMatDescr(descr);
descr = 0;
if (status != CUSPARSE_STATUS_SUCCESS) {
    CLEANUP("Matrix descriptor destruction failed");
    return 1;
}

/* destroy handle */
status = cusparseDestroy(handle);
handle = 0;
if (status != CUSPARSE_STATUS_SUCCESS) {
    CLEANUP("CUSPARSE Library release of resources failed");
    return 1;
}

/* check the results */
/* Notice that CLEANUP() contains a call to cusparseDestroy(handle) */
if ((zHostPtr[0] != 950.0) ||
    (zHostPtr[1] != 400.0) ||
    (zHostPtr[2] != 2550.0) ||
    (zHostPtr[3] != 2600.0) ||

```

```

        (zHostPtr[4] != 0.0)      ||
        (zHostPtr[5] != 49300.0) ||
        (zHostPtr[6] != 15200.0) ||
        (zHostPtr[7] != 132300.0) ||
        (zHostPtr[8] != 131200.0) ||
        (zHostPtr[9] != 0.0)      ||
        (yHostPtr[0] != 10.0)     ||
        (yHostPtr[1] != 20.0)     ||
        (yHostPtr[2] != 30.0)     ||
        (yHostPtr[3] != 40.0)     ||
        (yHostPtr[4] != 680.0)    ||
        (yHostPtr[5] != 760.0)    ||
        (yHostPtr[6] != 1230.0)   ||
        (yHostPtr[7] != 2240.0)){
    CLEANUP("example test FAILED");
    return 1;
}
else{
    CLEANUP("example test PASSED");
    return 0;
}
}

```


Chapter 13.

APPENDIX B: CUSPARSE FORTRAN BINDINGS

The CUSPARSE library is implemented using the C-based CUDA toolchain, and it thus provides a C-style API that makes interfacing to applications written in C or C++ trivial. There are also many applications implemented in Fortran that would benefit from using CUSPARSE, and therefore a CUSPARSE Fortran interface has been developed.

Unfortunately, Fortran-to-C calling conventions are not standardized and differ by platform and toolchain. In particular, differences may exist in the following areas:

- Symbol names (capitalization, name decoration)

- Argument passing (by value or reference)

- Passing of pointer arguments (size of the pointer)

To provide maximum flexibility in addressing those differences, the CUSPARSE Fortran interface is provided in the form of wrapper functions, which are written in C and are located in the file `cusparses_fortran.c`. This file also contains a few additional wrapper functions (for `cudaMalloc()`, `cudaMemset`, and so on) that can be used to allocate memory on the GPU.

The CUSPARSE Fortran wrapper code is provided as an example only and needs to be compiled into an application for it to call the CUSPARSE API functions. Providing this source code allows users to make any changes necessary for a particular platform and toolchain.

The CUSPARSE Fortran wrapper code has been used to demonstrate interoperability with the compilers g95 0.91 (on 32-bit and 64-bit Linux) and g95 0.92 (on 32-bit and 64-bit Mac OS X). In order to use other compilers, users have to make any changes to the wrapper code that may be required.

The direct wrappers, intended for production code, substitute device pointers for vector and matrix arguments in all CUSPARSE functions. To use these interfaces, existing applications need to be modified slightly to allocate and deallocate data structures in GPU memory space (using `CUDA_MALLOC()` and `CUDA_FREE()`) and to copy data between GPU and CPU memory spaces (using the `CUDA_MEMCPY()` routines). The sample wrappers provided in `cusparses_fortran.c` map device pointers to the OS-

dependent type `size_t`, which is 32 bits wide on 32-bit platforms and 64 bits wide on a 64-bit platforms.

One approach to dealing with index arithmetic on device pointers in Fortran code is to use C-style macros and to use the C preprocessor to expand them. On Linux and Mac OS X, preprocessing can be done by using the option `'-cpp'` with `g95` or `gfortran`. The function `GET_SHIFTED_ADDRESS()`, provided with the CUSPARSE Fortran wrappers, can also be used, as shown in example B.

Example B shows the the C++ of example A implemented in Fortran 77 on the host. This example should be compiled with `ARCH_64` defined as 1 on a 64-bit OS system and as undefined on a 32-bit OS system. For example, on `g95` or `gfortran`, it can be done directly on the command line using the option `-cpp -DARCH_64=1`.

13.1. Example B, Fortran Application

```
c      #define ARCH_64 0
c      #define ARCH_64 1

      program cusparse_fortran_example
      implicit none
      integer cuda_malloc
      external cuda_free
      integer cuda_memcpy_c2fort_int
      integer cuda_memcpy_c2fort_real
      integer cuda_memcpy_fort2c_int
      integer cuda_memcpy_fort2c_real
      integer cuda_memset
      integer cusparse_create
      external cusparse_destroy
      integer cusparse_get_version
      integer cusparse_create_mat_descr
      external cusparse_destroy_mat_descr
      integer cusparse_set_mat_type
      integer cusparse_get_mat_type
      integer cusparse_get_mat_fill_mode
      integer cusparse_get_mat_diag_type
      integer cusparse_set_mat_index_base
      integer cusparse_get_mat_index_base
      integer cusparse_xcoo2csr
      integer cusparse_dsctr
      integer cusparse_dcsmv
      integer cusparse_dcsrmm
      external get_shifted_address
      #if ARCH_64
        integer*8 handle
        integer*8 descrA
        integer*8 cooRowIndex
        integer*8 cooColIndex
        integer*8 cooVal
        integer*8 xInd
        integer*8 xVal
        integer*8 y
        integer*8 z
        integer*8 csrRowPtr
        integer*8 ynp1
      #else
        integer*4 handle
        integer*4 descrA
        integer*4 cooRowIndex
        integer*4 cooColIndex
```

```

integer*4 cooVal
integer*4 xInd
integer*4 xVal
integer*4 y
integer*4 z
integer*4 csrRowPtr
integer*4 ynp1
#endif
integer status
integer cudaStat1,cudaStat2,cudaStat3
integer cudaStat4,cudaStat5,cudaStat6
integer n, nnz, nnz_vector
parameter (n=4, nnz=9, nnz_vector=3)
integer cooRowIndexHostPtr(nnz)
integer cooColIndexHostPtr(nnz)
real*8 cooValHostPtr(nnz)
integer xIndHostPtr(nnz_vector)
real*8 xValHostPtr(nnz_vector)
real*8 yHostPtr(2*n)
real*8 zHostPtr(2*(n+1))
integer i, j
integer version, mtype, fmode, dtype, ibase
real*8 dzero,dtwo,dthree,dfive
real*8 epsilon

write(*,*) "testing fortran example"
c predefined constants (need to be careful with them)
dzero = 0.0
dtwo = 2.0
dthree= 3.0
dfive = 5.0
c create the following sparse test matrix in COO format
c (notice one-based indexing)
c |1.0      2.0 3.0|
c |      4.0      |
c |5.0      6.0 7.0|
c |      8.0      9.0|
c cooRowIndexHostPtr(1)=1
c cooColIndexHostPtr(1)=1
c cooValHostPtr(1)      =1.0
c cooRowIndexHostPtr(2)=1
c cooColIndexHostPtr(2)=3
c cooValHostPtr(2)      =2.0
c cooRowIndexHostPtr(3)=1
c cooColIndexHostPtr(3)=4
c cooValHostPtr(3)      =3.0
c cooRowIndexHostPtr(4)=2
c cooColIndexHostPtr(4)=2
c cooValHostPtr(4)      =4.0
c cooRowIndexHostPtr(5)=3
c cooColIndexHostPtr(5)=1
c cooValHostPtr(5)      =5.0
c cooRowIndexHostPtr(6)=3
c cooColIndexHostPtr(6)=3
c cooValHostPtr(6)      =6.0
c cooRowIndexHostPtr(7)=3
c cooColIndexHostPtr(7)=4
c cooValHostPtr(7)      =7.0
c cooRowIndexHostPtr(8)=4
c cooColIndexHostPtr(8)=2
c cooValHostPtr(8)      =8.0
c cooRowIndexHostPtr(9)=4
c cooColIndexHostPtr(9)=4
c cooValHostPtr(9)      =9.0
c print the matrix
write(*,*) "Input data:"

```

```

do i=1,nnz
  write(*,*) "cooRowIndexHostPtr[" ,i,"]=",cooRowIndexHostPtr(i)
  write(*,*) "cooColIndexHostPtr[" ,i,"]=",cooColIndexHostPtr(i)
  write(*,*) "cooValHostPtr[" ,i,"]=",cooValHostPtr(i)
enddo

c   create a sparse and dense vector
c   xVal= [100.0 200.0 400.0]      (sparse)
c   xInd= [0      1      3      ]
c   y   = [10.0 20.0 30.0 40.0 | 50.0 60.0 70.0 80.0] (dense)
c   (notice one-based indexing)
yHostPtr(1) = 10.0
yHostPtr(2) = 20.0
yHostPtr(3) = 30.0
yHostPtr(4) = 40.0
yHostPtr(5) = 50.0
yHostPtr(6) = 60.0
yHostPtr(7) = 70.0
yHostPtr(8) = 80.0
xIndHostPtr(1)=1
xValHostPtr(1)=100.0
xIndHostPtr(2)=2
xValHostPtr(2)=200.0
xIndHostPtr(3)=4
xValHostPtr(3)=400.0
c   print the vectors
do j=1,2
  do i=1,n
    write(*,*) "yHostPtr[" ,i," ",j,"]=",yHostPtr(i+n*(j-1))
  enddo
enddo
do i=1,nnz_vector
  write(*,*) "xIndHostPtr[" ,i,"]=",xIndHostPtr(i)
  write(*,*) "xValHostPtr[" ,i,"]=",xValHostPtr(i)
enddo

c   allocate GPU memory and copy the matrix and vectors into it
c   cudaSuccess=0
c   cudaMemcpyHostToDevice=1
cudaStat1 = cuda_malloc(cooRowIndex,nnz*4)
cudaStat2 = cuda_malloc(cooColIndex,nnz*4)
cudaStat3 = cuda_malloc(cooVal,      nnz*8)
cudaStat4 = cuda_malloc(y,          2*n*8)
cudaStat5 = cuda_malloc(xInd,nnz_vector*4)
cudaStat6 = cuda_malloc(xVal,nnz_vector*8)
if ((cudaStat1 /= 0) .OR.
$   (cudaStat2 /= 0) .OR.
$   (cudaStat3 /= 0) .OR.
$   (cudaStat4 /= 0) .OR.
$   (cudaStat5 /= 0) .OR.
$   (cudaStat6 /= 0)) then
  write(*,*) "Device malloc failed"
  write(*,*) "cudaStat1=",cudaStat1
  write(*,*) "cudaStat2=",cudaStat2
  write(*,*) "cudaStat3=",cudaStat3
  write(*,*) "cudaStat4=",cudaStat4
  write(*,*) "cudaStat5=",cudaStat5
  write(*,*) "cudaStat6=",cudaStat6
  stop
endif
cudaStat1 = cuda_memcpy_fort2c_int(cooRowIndex,cooRowIndexHostPtr,
$   nnz*4,1)
cudaStat2 = cuda_memcpy_fort2c_int(cooColIndex,cooColIndexHostPtr,
$   nnz*4,1)
cudaStat3 = cuda_memcpy_fort2c_real(cooVal,      cooValHostPtr,
$   nnz*8,1)
cudaStat4 = cuda_memcpy_fort2c_real(y,          yHostPtr,

```

```

$                                2*n*8,1)
  cudaStat5 = cuda_memcpy_fort2c_int(xInd,      xIndHostPtr,
$                                nnz_vector*4,1)
  cudaStat6 = cuda_memcpy_fort2c_real(xVal,      xValHostPtr,
$                                nnz_vector*8,1)
  if ((cudaStat1 /= 0) .OR.
$    (cudaStat2 /= 0) .OR.
$    (cudaStat3 /= 0) .OR.
$    (cudaStat4 /= 0) .OR.
$    (cudaStat5 /= 0) .OR.
$    (cudaStat6 /= 0)) then
    write(*,*) "Memcpy from Host to Device failed"
    write(*,*) "cudaStat1=",cudaStat1
    write(*,*) "cudaStat2=",cudaStat2
    write(*,*) "cudaStat3=",cudaStat3
    write(*,*) "cudaStat4=",cudaStat4
    write(*,*) "cudaStat5=",cudaStat5
    write(*,*) "cudaStat6=",cudaStat6
    call cuda_free(cooRowIndex)
    call cuda_free(cooColIndex)
    call cuda_free(cooVal)
    call cuda_free(xInd)
    call cuda_free(xVal)
    call cuda_free(y)
    stop
  endif

c  initialize cusparse library
c  CUSPARSE_STATUS_SUCCESS=0
  status = cusparse_create(handle)
  if (status /= 0) then
    write(*,*) "CUSPARSE Library initialization failed"
    call cuda_free(cooRowIndex)
    call cuda_free(cooColIndex)
    call cuda_free(cooVal)
    call cuda_free(xInd)
    call cuda_free(xVal)
    call cuda_free(y)
    stop
  endif
c  get version
c  CUSPARSE_STATUS_SUCCESS=0
  status = cusparse_get_version(handle,version)
  if (status /= 0) then
    write(*,*) "CUSPARSE Library initialization failed"
    call cuda_free(cooRowIndex)
    call cuda_free(cooColIndex)
    call cuda_free(cooVal)
    call cuda_free(xInd)
    call cuda_free(xVal)
    call cuda_free(y)
    call cusparse_destroy(handle)
    stop
  endif
  write(*,*) "CUSPARSE Library version",version

c  create and setup the matrix descriptor
c  CUSPARSE_STATUS_SUCCESS=0
c  CUSPARSE_MATRIX_TYPE_GENERAL=0
c  CUSPARSE_INDEX_BASE_ONE=1
  status= cusparse_create_mat_descr(descrA)
  if (status /= 0) then
    write(*,*) "Creating matrix descriptor failed"
    call cuda_free(cooRowIndex)
    call cuda_free(cooColIndex)
    call cuda_free(cooVal)
    call cuda_free(xInd)

```

```

        call cuda_free(xVal)
        call cuda_free(y)
        call cusparse_destroy(handle)
        stop
    endif
    status = cusparse_set_mat_type(descrA,0)
    status = cusparse_set_mat_index_base(descrA,1)
c    print the matrix descriptor
    mtype = cusparse_get_mat_type(descrA)
    fmode = cusparse_get_mat_fill_mode(descrA)
    dtype = cusparse_get_mat_diag_type(descrA)
    ibase = cusparse_get_mat_index_base(descrA)
    write (*,*) "matrix descriptor:"
    write (*,*) "t=",mtype,"m=",fmode,"d=",dtype,"b=",ibase

c    exercise conversion routines (convert matrix from COO 2 CSR format)
c    cudaSuccess=0
c    CUSPARSE_STATUS_SUCCESS=0
c    CUSPARSE_INDEX_BASE_ONE=1
    cudaStat1 = cuda_malloc(csrRowPtr,(n+1)*4)
    if (cudaStat1 /= 0) then
        call cuda_free(cooRowIndex)
        call cuda_free(cooColIndex)
        call cuda_free(cooVal)
        call cuda_free(xInd)
        call cuda_free(xVal)
        call cuda_free(y)
        call cusparse_destroy_mat_descr(descrA)
        call cusparse_destroy(handle)
        write(*,*) "Device malloc failed (csrRowPtr)"
        stop
    endif
    status= cusparse_xcoo2csr(handle,cooRowIndex,nnz,n,
$                                csrRowPtr,1)
    if (status /= 0) then
        call cuda_free(cooRowIndex)
        call cuda_free(cooColIndex)
        call cuda_free(cooVal)
        call cuda_free(xInd)
        call cuda_free(xVal)
        call cuda_free(y)
        call cuda_free(csrRowPtr)
        call cusparse_destroy_mat_descr(descrA)
        call cusparse_destroy(handle)
        write(*,*) "Conversion from COO to CSR format failed"
        stop
    endif
c    csrRowPtr = [0 3 4 7 9]

c    exercise Level 1 routines (scatter vector elements)
c    CUSPARSE_STATUS_SUCCESS=0
c    CUSPARSE_INDEX_BASE_ONE=1
    call get_shifted_address(y,n*8,ynp1)
    status= cusparse_dsctr(handle, nnz_vector, xVal, xInd,
$                                ynp1, 1)
    if (status /= 0) then
        call cuda_free(cooRowIndex)
        call cuda_free(cooColIndex)
        call cuda_free(cooVal)
        call cuda_free(xInd)
        call cuda_free(xVal)
        call cuda_free(y)
        call cuda_free(csrRowPtr)
        call cusparse_destroy_mat_descr(descrA)
        call cusparse_destroy(handle)
        write(*,*) "Scatter from sparse to dense vector failed"
        stop
    endif

```

```

endif
c      y = [10 20 30 40 | 100 200 70 400]

c      exercise Level 2 routines (csrcmv)
c      CUSPARSE_STATUS_SUCCESS=0
c      CUSPARSE_OPERATION_NON_TRANSPOSE=0
      status= cusparse_dcsrmv(handle, 0, n, n, nnz, dtwo,
$          descrA, cooVal, csrRowPtr, cooColIndex,
$          y, dthree, ynpl)
      if (status /= 0) then
          call cuda_free(cooRowIndex)
          call cuda_free(cooColIndex)
          call cuda_free(cooVal)
          call cuda_free(xInd)
          call cuda_free(xVal)
          call cuda_free(y)
          call cuda_free(csrRowPtr)
          call cusparse_destroy_mat_descr(descrA)
          call cusparse_destroy(handle)
          write(*,*) "Matrix-vector multiplication failed"
          stop
      endif

c      print intermediate results (y)
c      y = [10 20 30 40 | 680 760 1230 2240]
c      cudaSuccess=0
c      cudaMemcpyDeviceToHost=2
      cudaStat1 = cuda_memcpy_c2fort_real(yHostPtr, y, 2*n*8, 2)
      if (cudaStat1 /= 0) then
          call cuda_free(cooRowIndex)
          call cuda_free(cooColIndex)
          call cuda_free(cooVal)
          call cuda_free(xInd)
          call cuda_free(xVal)
          call cuda_free(y)
          call cuda_free(csrRowPtr)
          call cusparse_destroy_mat_descr(descrA)
          call cusparse_destroy(handle)
          write(*,*) "Memcpy from Device to Host failed"
          stop
      endif
      write(*,*) "Intermediate results:"
      do j=1,2
          do i=1,n
              write(*,*) "yHostPtr[" ,i, ",",j,"]=",yHostPtr(i+n*(j-1))
          enddo
      enddo

c      exercise Level 3 routines (csrmm)
c      cudaSuccess=0
c      CUSPARSE_STATUS_SUCCESS=0
c      CUSPARSE_OPERATION_NON_TRANSPOSE=0
      cudaStat1 = cuda_malloc(z, 2*(n+1)*8)
      if (cudaStat1 /= 0) then
          call cuda_free(cooRowIndex)
          call cuda_free(cooColIndex)
          call cuda_free(cooVal)
          call cuda_free(xInd)
          call cuda_free(xVal)
          call cuda_free(y)
          call cuda_free(csrRowPtr)
          call cusparse_destroy_mat_descr(descrA)
          call cusparse_destroy(handle)
          write(*,*) "Device malloc failed (z)"
          stop
      endif
      cudaStat1 = cuda_memset(z, 0, 2*(n+1)*8)

```

```

    if (cudaStat1 /= 0) then
        call cuda_free(cooRowIndex)
        call cuda_free(cooColIndex)
        call cuda_free(cooVal)
        call cuda_free(xInd)
        call cuda_free(xVal)
        call cuda_free(y)
        call cuda_free(z)
        call cuda_free(csrRowPtr)
        call cusparse_destroy_mat_descr(descrA)
        call cusparse_destroy(handle)
        write(*,*) "Memset on Device failed"
        stop
    endif
    status= cusparse_dcsrmm(handle, 0, n, 2, n, nnz, dfive,
$                                descrA, cooVal, csrRowPtr, cooColIndex,
$                                y, n, dzero, z, n+1)
    if (status /= 0) then
        call cuda_free(cooRowIndex)
        call cuda_free(cooColIndex)
        call cuda_free(cooVal)
        call cuda_free(xInd)
        call cuda_free(xVal)
        call cuda_free(y)
        call cuda_free(z)
        call cuda_free(csrRowPtr)
        call cusparse_destroy_mat_descr(descrA)
        call cusparse_destroy(handle)
        write(*,*) "Matrix-matrix multiplication failed"
        stop
    endif

c    print final results (z)
c    cudaSuccess=0
c    cudaMemcpyDeviceToHost=2
    cudaStat1 = cuda_memcpy_c2fort_real(zHostPtr, z, 2*(n+1)*8, 2)
    if (cudaStat1 /= 0) then
        call cuda_free(cooRowIndex)
        call cuda_free(cooColIndex)
        call cuda_free(cooVal)
        call cuda_free(xInd)
        call cuda_free(xVal)
        call cuda_free(y)
        call cuda_free(z)
        call cuda_free(csrRowPtr)
        call cusparse_destroy_mat_descr(descrA)
        call cusparse_destroy(handle)
        write(*,*) "Memcpy from Device to Host failed"
        stop
    endif
c    z = [950 400 2550 2600 0 | 49300 15200 132300 131200 0]
    write(*,*) "Final results:"
    do j=1,2
        do i=1,n+1
            write(*,*) "z[" ,i ,",",j ,"]=",zHostPtr(i+(n+1)*(j-1))
        enddo
    enddo

c    check the results
    epsilon = 0.000000000000001
    if ((DABS(zHostPtr(1) - 950.0) .GT. epsilon) .OR.
$      (DABS(zHostPtr(2) - 400.0) .GT. epsilon) .OR.
$      (DABS(zHostPtr(3) - 2550.0) .GT. epsilon) .OR.
$      (DABS(zHostPtr(4) - 2600.0) .GT. epsilon) .OR.
$      (DABS(zHostPtr(5) - 0.0) .GT. epsilon) .OR.
$      (DABS(zHostPtr(6) - 49300.0) .GT. epsilon) .OR.
$      (DABS(zHostPtr(7) - 15200.0) .GT. epsilon) .OR.

```



```

$      (DABS(zHostPtr(8) - 132300.0).GT. epsilon) .OR.
$      (DABS(zHostPtr(9) - 131200.0).GT. epsilon) .OR.
$      (DABS(zHostPtr(10) - 0.0) .GT. epsilon) .OR.
$      (DABS(yHostPtr(1) - 10.0) .GT. epsilon) .OR.
$      (DABS(yHostPtr(2) - 20.0) .GT. epsilon) .OR.
$      (DABS(yHostPtr(3) - 30.0) .GT. epsilon) .OR.
$      (DABS(yHostPtr(4) - 40.0) .GT. epsilon) .OR.
$      (DABS(yHostPtr(5) - 680.0) .GT. epsilon) .OR.
$      (DABS(yHostPtr(6) - 760.0) .GT. epsilon) .OR.
$      (DABS(yHostPtr(7) - 1230.0) .GT. epsilon) .OR.
$      (DABS(yHostPtr(8) - 2240.0) .GT. epsilon) then
    write(*,*) "fortran example test FAILED"
  else
    write(*,*) "fortran example test PASSED"
  endif
c
  deallocate GPU memory and exit
  call cuda_free(cooRowIndex)
  call cuda_free(cooColIndex)
  call cuda_free(cooVal)
  call cuda_free(xInd)
  call cuda_free(xVal)
  call cuda_free(y)
  call cuda_free(z)
  call cuda_free(csrRowPtr)
  call cusparse_destroy_mat_descr(descrA)
  call cusparse_destroy(handle)

  stop
end

```

Chapter 14.

APPENDIX C: ACKNOWLEDGEMENTS

NVIDIA would like to thank the following individuals and institutions for their contributions:

- ▶ The cusparse\timesgtsv implementation is derived from a version developed by Li-Wen Chang from the University of Illinois.

Chapter 15.

BIBLIOGRAPHY

- [1] N. Bell and M. Garland, “Implementing Sparse Matrix-Vector Multiplication on Throughput-Oriented Processors”, Supercomputing, 2009.
- [2] R. Grimes, D. Kincaid, and D. Young, “ITPACK 2.0 User’s Guide”, Technical Report CNA-150, Center for Numerical Analysis, University of Texas, 1979.
- [3] M. Naumov, “Incomplete-LU and Cholesky Preconditioned Iterative Methods Using CUSPARSE and CUBLAS”, Technical Report and White Paper, 2011.

Notice

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

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

Trademarks

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

Copyright

© 2007-2013 NVIDIA Corporation. All rights reserved.