
hipBLASLt Documentation

Release 0.12.1

Advanced Micro Devices, Inc.

Sep 10, 2025

INSTALL

1	Installation	3
1.1	Prerequisites	3
1.2	Installing prebuilt packages	3
1.3	Build using script	3
1.4	Manual build (all supported platforms)	4
2	hipBLASLt datatypes reference	7
2.1	hipblasLtEpilogue_t	7
2.2	hipblasLtPointerMode_t	8
2.3	hipblasLtHandle_t	8
2.4	hipblasLtMatmulAlgo_t	8
2.5	hipblasLtMatmulDesc_t	9
2.6	hipblasLtMatmulDescAttributes_t	9
2.7	hipblasLtMatmulHeuristicResult_t	11
2.8	hipblasLtMatmulPreference_t	11
2.9	hipblasLtMatmulPreferenceAttributes_t	11
2.10	hipblasLtMatrixLayout_t	11
2.11	hipblasLtMatrixLayoutAttribute_t	12
2.12	hipblasLtMatrixTransformDesc_t	13
3	hipBLASLt API reference	15
3.1	hipblasLtCreate()	15
3.2	hipblasLtDestroy()	15
3.3	hipblasLtMatrixLayoutCreate()	16
3.4	hipblasLtMatrixLayoutDestroy()	16
3.5	hipblasLtMatrixLayoutSetAttribute()	16
3.6	hipblasLtMatrixLayoutGetAttribute()	17
3.7	hipblasLtMatmulDescCreate()	17
3.8	hipblasLtMatmulDescDestroy()	18
3.9	hipblasLtMatmulDescSetAttribute()	18
3.10	hipblasLtMatmulDescGetAttribute()	18
3.11	hipblasLtMatmulPreferenceCreate()	19
3.12	hipblasLtMatmulPreferenceDestroy()	19
3.13	hipblasLtMatmulPreferenceSetAttribute()	20
3.14	hipblasLtMatmulPreferenceGetAttribute()	20
3.15	hipblasLtMatmulAlgoGetHeuristic()	21
3.16	hipblasLtMatmul()	21
3.17	hipblasLtMatrixTransformDescCreate()	23
3.18	hipblasLtMatrixTransformDescDestroy()	24
3.19	hipblasLtMatrixTransformDescSetAttribute()	24

3.20	hipblasLtMatrixTransformDescGetAttribute()	24
3.21	hipblasLtMatrixTransform()	25
4	hipBLASLtExt API reference	27
4.1	hipBLASLtExt datatypes reference	27
4.2	hipBLASLtExt Gemm Class Reference	33
4.3	hipBLASLtExt API Reference	46
4.4	hipblasLtExt usage	48
4.5	Sample codes	57
5	hipBLASLtExt operation API reference	65
5.1	hipblasLtExtSoftmax()	65
5.2	hipblasLtExtLayerNorm()	66
5.3	hipblasLtExtAMax()	66
5.4	hipblasLtExtAMaxWithScale()	67
6	Logging and heuristics	69
6.1	Logging	69
6.2	Heuristics cache	70
7	Using the hipBLASLt Tuning Utility	71
7.1	The template.yaml file	71
7.2	Running the tuning utility	71
8	User Offline Tuning	73
8.1	Using hipblaslt-bench to run the tuning with the best GEMM kernel	73
9	Using Stream-K with hipBLASLt	75
9.1	Configuring the kernel selection strategy	75
9.2	Configuring the kernel selection strategy	75
9.3	Recommendations for using Stream-K	76
10	Library source code organization	77
10.1	library directory	77
10.2	clients directory	77
10.3	Infrastructure	77
11	Clients	79
11.1	hipblaslt-test	79
11.2	hipblaslt-bench	79
12	License	83
	Index	85

hipBLASLt is a library that provides General Matrix-Matrix (GEMM) operations with flexible APIs and extends functionalities beyond the traditional BLAS library. To learn more, see [what-is-hipblaslt](#)

You can access hipBLASLt code on our [GitHub repository](#).

Our documentation is structured as follows:

Install

- [*Installation*](#)

API reference

- [*hipBLASLt datatypes reference*](#)
- [*hipBLASLt API reference*](#)
- [*hipBLASLtExt API reference*](#)
- [*hipBLASLtExt operation API reference*](#)
- [*Logging and heuristics*](#)

How to

- [*Using the hipBLASLt Tuning Utility*](#)
- [*User Offline Tuning*](#)
- [*Use Stream-K with hipBLASLt*](#)

Contribution

- [*Library source code organization*](#)
- [*Clients*](#)

To contribute to the documentation, refer to [Contributing to ROCm](#).

You can find licensing information on the [Licensing](#) page.

INSTALLATION

This document provides all the information required to build and install hipBLASLt on Linux systems.

1.1 Prerequisites

- A ROCm enabled platform. For more information refer to [ROCm Documentation](#).
- A compatible version of hipBLAS.

1.2 Installing prebuilt packages

Download prebuilt packages from [ROCm's native package manager](#) .

```
sudo apt update && sudo apt install hipblaslt
```

1.3 Build using script

You can use `install.sh` script to build and install hipBLASLt and dependencies. The following sections explain the usage of `install.sh` script along with various options.

1.3.1 Build library dependencies and library

The root of this repository has a helper bash script `install.sh` to build and install hipBLASLt with a single command. It takes a lot of options and hard-coded configuration that can be specified through invoking `cmake` directly, but it's a great way to get started quickly and can serve as an example of how to build/install. A few commands in the script need `sudo` access so it may prompt you for a password.

Typical uses of `install.sh` to build (library dependencies and library) are listed below:

Command	Description
<code>./install.sh -h</code>	Help information.
<code>./install.sh -d</code>	Build library dependencies and library in your local directory. Use <code>-d</code> flag only once. For subsequent invocations of <code>install.sh</code> , it is not necessary to rebuild the dependencies.
<code>./install.sh</code>	Build library in your local directory. The dependencies are assumed to be already built.
<code>./install.sh -i</code>	Build library, then build and install hipBLASLt package in <code>/opt/rocm/hipblaslt</code> . This prompts for <code>sudo</code> access and installs for all users. If you want to keep hipBLASLt in your local directory, don't use <code>-i</code> flag.

1.3.2 Build library dependencies, client dependencies, library, and client

The client contains executables as listed below:

Executable Name	Description
hipblaslt-test	Runs Google tests to test the library
hipblaslt-bench	Executable to benchmark or test individual functions

Common uses of `install.sh` to build (dependencies, library, and client) are listed below:

Command	Description
<code>./install.sh -h</code>	Help information.
<code>./install.sh -dc</code>	Build library dependencies, client dependencies, library, and client in your local directory. Use <code>-d</code> flag only once. For subsequent invocations of <code>install.sh</code> , it is not necessary to rebuild the dependencies.
<code>./install.sh -c</code>	Build library and client in your local directory. The dependencies are assumed to be already built.
<code>./install.sh -idc</code>	Build library dependencies, client dependencies, library, client, then build and install the hipBLASLt package. This prompts for sudo access. To install for all users, use <code>-i</code> flag. To keep hipBLASLt in your local directory, don't use <code>-i</code> flag.
<code>./install.sh -ic</code>	Build and install hipBLASLt package, and build the client. This prompts for sudo access and installs for all users. To keep hipBLASLt in your local directory, don't use <code>-i</code> flag.

1.3.3 Static Library

To build static libraries with `install.sh`, use `--static`. This produces a non-standard static library build. This means it has an additional runtime dependency consisting of the entire `hipblaslt/` subdirectory, which is located in the `/opt/rocm/lib` folder. You can relocate this folder, but you must set the environment variable `HIPBLASLT_TENSILE_LIBPATH` accordingly.

1.3.4 Dependencies

Dependencies are listed in the `install.sh` script. Use `install.sh` with `-d` option to install dependencies. CMake has a minimum version requirement which is listed in `install.sh`. See `--cmake_install` flag in `install.sh` to upgrade automatically.

1.4 Manual build (all supported platforms)

This section provides information on how to configure cmake and build manually using individual commands.

1.4.1 Build library manually

Dependencies for hipBLASLt

1. The `hipBLAS-common` header files are required to be installed on your system.
2. Usually `roctx` library is pre-installed, if not please install it on your system.

Building hipBLASLt

```
mkdir -p [HIPBLASLT_BUILD_DIR]/release
cd [HIPBLASLT_BUILD_DIR]/release
# Default install location is in /opt/rocm, define -DCMAKE_INSTALL_PREFIX=<path> to
→specify other
# Default build config is 'Release', define -DCMAKE_BUILD_TYPE=<config> to specify other
CXX=/opt/rocm/bin/amdclang++ ccmake [HIPBLASLT_SOURCE]
make -j$(nproc)
sudo make install # sudo required if installing into system directory such as /opt/rocm
```

1.4.2 Build library, tests, benchmarks, and samples manually

The repository contains source for clients that serve as samples, tests, and benchmarks. You can find the clients source in the clients sub-directory.

Dependencies for hipBLASLt clients

The hipBLASLt samples have no external dependencies, but unit test and benchmarking applications do. These clients introduce the following dependencies:

- Lapack, Lapack itself brings a dependency on a fortran compiler
- googletest

Googletest and Lapack are not easy to install. Many distros don't provide a googletest package with precompiled libraries and the Lapack packages don't have the necessary cmake config files for cmake to configure linking the cblas library. hipBLASLt provides a cmake script that builds the above dependencies from source. This is an optional step; you can provide your own builds of these dependencies and help cmake find them by setting the CMAKE_PREFIX_PATH definition. The following is a sequence of steps to build dependencies and install them to the cmake default /usr/local.

One-time optional step

```
mkdir -p [HIPBLASLT_BUILD_DIR]/release/deps
cd [HIPBLASLT_BUILD_DIR]/release/deps
ccmake -DBUILD_BOOST=OFF [HIPBLASLT_SOURCE]/deps # assuming boost is installed through
→package manager as above
make -j$(nproc) install
```

Once dependencies are available on the system, it is possible to configure the clients to build. This requires a few extra cmake flags to the library cmake configure script. If the dependencies are not installed into system defaults (like /usr/local), pass the CMAKE_PREFIX_PATH to cmake to help find them.

```
-DCMAKE_PREFIX_PATH="<semicolon separated paths>"
# Default install location is in /opt/rocm, use -DCMAKE_INSTALL_PREFIX=<path> to specify
→other
CXX=/opt/rocm/bin/amdclang++ ccmake -DBUILD_CLIENTS_TESTS=ON -DBUILD_CLIENTS_
→BENCHMARKS=ON [HIPBLASLT_SOURCE]
make -j$(nproc)
sudo make install # sudo required if installing into system directory such as /opt/rocm
```


HIPBLASLT DATATYPES REFERENCE

2.1 hipblasLtEpilogue_t

enum **hipblasLtEpilogue_t**

Specify the enum type to set the postprocessing options for the epilogue.

Values:

enumerator **HIPBLASLT_EPILOGUE_DEFAULT**

No special postprocessing, just scale and quantize the results if necessary.

enumerator **HIPBLASLT_EPILOGUE_RELU**

Apply ReLU point-wise transform to the results:($x:=\max(x, 0)$)

enumerator **HIPBLASLT_EPILOGUE_BIAS**

Apply (broadcast) bias from the bias vector. Bias vector length must match matrix D rows, and it must be packed (such as stride between vector elements is 1). Bias vector is broadcast to all columns and added before applying the final postprocessing.

enumerator **HIPBLASLT_EPILOGUE_RELU_BIAS**

Apply bias and then ReLU transform.

enumerator **HIPBLASLT_EPILOGUE_GELU**

Apply GELU point-wise transform to the results ($x:=\text{GELU}(x)$).

enumerator **HIPBLASLT_EPILOGUE_GELU_BIAS**

Apply Bias and then GELU transform.

enumerator **HIPBLASLT_EPILOGUE_GELU_AUX**

Output GEMM results before applying GELU transform.

enumerator **HIPBLASLT_EPILOGUE_GELU_AUX_BIAS**

Output GEMM results after applying bias but before applying GELU transform.

enumerator **HIPBLASLT_EPILOGUE_DGELU**

Apply gradient GELU transform. Requires additional aux input.

enumerator **HIPBLASLT_EPILOGUE_DGELU_BGRAD**

Apply gradient GELU transform and bias gradient to the results. Requires additional aux input.

enumerator **HIPBLASLT_EPILOGUE_BGRADA**

Apply bias gradient to A and output gemm result.

enumerator **HIPBLASLT_EPILOGUE_BGRADB**

Apply bias gradient to B and output gemm result.

2.2 hipblasLtPointerMode_t

enum **hipblasLtPointerMode_t**

Pointer mode to use for alpha.

Values:

enumerator **HIPBLASLT_POINTER_MODE_HOST**

enumerator **HIPBLASLT_POINTER_MODE_DEVICE**

targets host memory

enumerator **HIPBLASLT_POINTER_MODE_ALPHA_DEVICE_VECTOR_BETA_HOST**

targets device memory

2.3 hipblasLtHandle_t

typedef void ***hipblasLtHandle_t**

Handle to the hipBLASLt library context queue.

The `hipblasLtHandle_t` type is a pointer type to an opaque structure holding the hipBLASLt library context. Use the following functions to manipulate this library context:

hipblasLtCreate(): To initialize the hipBLASLt library context and return a handle to an opaque structure holding the hipBLASLt library context. *hipblasLtDestroy()*: To destroy a previously created hipBLASLt library context descriptor and release the resources.

2.4 hipblasLtMatmulAlgo_t

struct **hipblasLtMatmulAlgo_t**

Description of the matrix multiplication algorithm.

This is an opaque structure holding the description of the matrix multiplication algorithm. This structure can be trivially serialized and later restored for use with the same version of hipBLASLt library to save on selecting the right configuration again.

2.5 hipblasLtMatmulDesc_t

typedef hipblasLtMatmulDescOpaque_t ***hipblasLtMatmulDesc_t**

Descriptor of the matrix multiplication operation.

This is a pointer to an opaque structure holding the description of the matrix multiplication operation *hipblasLtMatmul()*. Use the following functions to manipulate this descriptor: *hipblasLtMatmulDescCreate()*: To create one instance of the descriptor. *hipblasLtMatmulDescDestroy()*: To destroy a previously created descriptor and release the resources.

2.6 hipblasLtMatmulDescAttributes_t

enum **hipblasLtMatmulDescAttributes_t**

Specify the attributes that define the specifics of the matrix multiply operation.

Values:

enumerator **HIPBLASLT_MATMUL_DESC_TRANSA**

Specifies the type of transformation operation that should be performed on matrix A. Default value is HIPBLAS_OP_N (for example, non-transpose operation). See hipblasOperation_t. Data Type:int32_t

enumerator **HIPBLASLT_MATMUL_DESC_TRANSB**

Specifies the type of transformation operation that should be performed on matrix B. Default value is HIPBLAS_OP_N (for example, non-transpose operation). See hipblasOperation_t. Data Type:int32_t

enumerator **HIPBLASLT_MATMUL_DESC_EPILOGUE**

Epilogue function. See hipblasLtEpilogue_t. Default value is: HIPBLASLT_EPILOGUE_DEFAULT. Data Type: uint32_t

enumerator **HIPBLASLT_MATMUL_DESC_BIAS_POINTER**

Bias or Bias gradient vector pointer in the device memory. Data Type:void* /const void*

enumerator **HIPBLASLT_MATMUL_DESC_BIAS_DATA_TYPE**

Type of the bias vector in the device memory. Can be set same as D matrix type or Scale type. Bias case: see HIPBLASLT_EPILOGUE_BIAS. Data Type:int32_t based on hipDataType

enumerator **HIPBLASLT_MATMUL_DESC_A_SCALE_POINTER**

Device pointer to the scale factor value that converts data in matrix A to the compute data type range. The scaling factor must have the same type as the compute type. If not specified, or set to NULL, the scaling factor is assumed to be 1. If set for an unsupported matrix data, scale, and compute type combination, calling *hipblasLtMatmul()* will return HIPBLAS_INVALID_VALUE. Default value: NULL Data Type: void* /const void*

enumerator **HIPBLASLT_MATMUL_DESC_B_SCALE_POINTER**

Equivalent to HIPBLASLT_MATMUL_DESC_A_SCALE_POINTER for matrix B. Default value: NULL Type: void* /const void*

enumerator **HIPBLASLT_MATMUL_DESC_C_SCALE_POINTER**

Equivalent to HIPBLASLT_MATMUL_DESC_A_SCALE_POINTER for matrix C. Default value: NULL
Type: void* /const void*

enumerator **HIPBLASLT_MATMUL_DESC_D_SCALE_POINTER**

Equivalent to HIPBLASLT_MATMUL_DESC_A_SCALE_POINTER for matrix D. Default value: NULL
Type: void* /const void*

enumerator **HIPBLASLT_MATMUL_DESC_EPILOGUE_AUX_SCALE_POINTER**

Equivalent to HIPBLASLT_MATMUL_DESC_A_SCALE_POINTER for matrix AUX. Default value: NULL
Type: void* /const void*

enumerator **HIPBLASLT_MATMUL_DESC_EPILOGUE_AUX_POINTER**

Epilogue auxiliary buffer pointer in the device memory. Data Type: void* /const void*

enumerator **HIPBLASLT_MATMUL_DESC_EPILOGUE_AUX_LD**

The leading dimension of the epilogue auxiliary buffer pointer in the device memory. Data Type: int64_t

enumerator **HIPBLASLT_MATMUL_DESC_EPILOGUE_AUX_BATCH_STRIDE**

The batch stride of the epilogue auxiliary buffer pointer in the device memory. Data Type: int64_t

enumerator **HIPBLASLT_MATMUL_DESC_POINTER_MODE**

Specifies alpha and beta are passed by reference, whether they are scalars on the host or on the device, or device vectors. Default value is: HIPBLASLT_POINTER_MODE_HOST (i.e., on the host). Data Type: int32_t based on hipblasLtPointerMode_t

enumerator **HIPBLASLT_MATMUL_DESC_AMAX_D_POINTER**

Device pointer to the memory location that on completion will be set to the maximum of absolute values in the output matrix. Data Type: void* /const void*

enumerator **HIPBLASLT_MATMUL_DESC_COMPUTE_INPUT_TYPE_A_EXT**

Compute input A types. Defines the data type used for the input A of matrix multiply.

enumerator **HIPBLASLT_MATMUL_DESC_COMPUTE_INPUT_TYPE_B_EXT**

Compute input B types. Defines the data type used for the input B of matrix multiply.

enumerator **HIPBLASLT_MATMUL_DESC_A_SCALE_POINTER_VEC_EXT**

Equivalent to HIPBLASLT_MATMUL_DESC_A_SCALE_POINTER but in vector. Default value: NULL
Type: void* /const void*

enumerator **HIPBLASLT_MATMUL_DESC_B_SCALE_POINTER_VEC_EXT**

Equivalent to HIPBLASLT_MATMUL_DESC_B_SCALE_POINTER but in vector. Default value: NULL
Type: void* /const void*

enumerator **HIPBLASLT_MATMUL_DESC_MAX**

2.7 hipblasLtMatmulHeuristicResult_t

struct **hipblasLtMatmulHeuristicResult_t**

Description of the matrix multiplication algorithm.

This is a descriptor that holds the configured matrix multiplication algorithm descriptor and its runtime properties. This structure can be trivially serialized and later restored for use with the same version of hipBLASLt library to save on selecting the right configuration again.

Param algo

hipblasLtMatmulAlgo_t struct

Param workspaceSize

Actual size of workspace memory required

Param state

Result status. Other fields are valid only if, after call to *hipblasLtMatmulAlgoGetHeuristic()*, this member is set to HIPBLAS_STATUS_SUCCESS

2.8 hipblasLtMatmulPreference_t

typedef hipblasLtMatmulPreferenceOpaque_t ***hipblasLtMatmulPreference_t**

Descriptor of the matrix multiplication preference.

This is a pointer to an opaque structure holding the description of the preferences for *hipblasLtMatmulAlgoGetHeuristic()* configuration. Use the following functions to manipulate this descriptor: *hipblasLtMatmulPreferenceCreate()*: To create one instance of the descriptor. *hipblasLtMatmulPreferenceDestroy()*: To destroy a previously created descriptor and release the resources.

2.9 hipblasLtMatmulPreferenceAttributes_t

enum **hipblasLtMatmulPreferenceAttributes_t**

It is an enumerated type used to apply algorithm search preferences while fine-tuning the heuristic function.

Values:

enumerator **HIPBLASLT_MATMUL_PREF_SEARCH_MODE**

Search mode. Data Type: uint32_t

enumerator **HIPBLASLT_MATMUL_PREF_MAX_WORKSPACE_BYTES**

Maximum allowed workspace memory. Default is 0 (no workspace memory allowed). Data Type: uint64_t

enumerator **HIPBLASLT_MATMUL_PREF_MAX**

2.10 hipblasLtMatrixLayout_t

typedef hipblasLtMatrixLayoutOpaque_t ***hipblasLtMatrixLayout_t**

Descriptor of the matrix layout.

This is a pointer to an opaque structure holding the description of a matrix layout. Use the following functions to manipulate this descriptor: *hipblasLtMatrixLayoutCreate()*: To create one instance of the descriptor. *hipblasLtMatrixLayoutDestroy()*: To destroy a previously created descriptor and release the resources.

2.11 hipblasLtMatrixLayoutAttribute_t

enum **hipblasLtMatrixLayoutAttribute_t**

Specify the attributes that define the details of the matrix.

Values:

enumerator **HIPBLASLT_MATRIX_LAYOUT_BATCH_COUNT**

Number of batch of this matrix. Default value is 1. Data Type: int32_t

enumerator **HIPBLASLT_MATRIX_LAYOUT_STRIDED_BATCH_OFFSET**

Stride (in elements) to the next matrix for the strided batch operation. Default value is 0. Data Type: int64_t

enumerator **HIPBLASLT_MATRIX_LAYOUT_TYPE**

Data type, see hipDataType.

uint32_t

enumerator **HIPBLASLT_MATRIX_LAYOUT_ORDER**

Memory order of the data, see hipblasLtOrder_t.

int32_t, default: HIPBLASLT_ORDER_COL

enumerator **HIPBLASLT_MATRIX_LAYOUT_ROWS**

Number of rows.

Usually only values that can be expressed as int32_t are supported.

uint64_t

enumerator **HIPBLASLT_MATRIX_LAYOUT_COLS**

Number of columns.

Usually only values that can be expressed as int32_t are supported.

uint64_t

enumerator **HIPBLASLT_MATRIX_LAYOUT_LD**

Matrix leading dimension.

For HIPBLASLT_ORDER_COL this is stride (in elements) of matrix column, for more details and documentation for other memory orders see documentation for hipblasLtOrder_t values.

Currently only non-negative values are supported, must be large enough so that matrix memory locations are not overlapping (e.g. greater or equal to HIPBLASLT_MATRIX_LAYOUT_ROWS in case of HIPBLASLT_ORDER_COL).

int64_t;

2.12 hipblasLtMatrixTransformDesc_t

typedef hipblasLtMatrixTransformDescOpaque_t ***hipblasLtMatrixTransformDesc_t**

Opaque descriptor for *hipblasLtMatrixTransform()* operation details.

The `hipblasLtMatrixTransformDesc_t` is a pointer to an opaque structure holding the description of a matrix transformation operation. *hipblasLtMatrixTransformDescCreate()*: To create one instance of the descriptor. *hipblasLtMatrixTransformDescDestroy()*: To destroy a previously created descriptor and release the resources.

HIPBLASLT API REFERENCE

3.1 hipblasLtCreate()

hipblasStatus_t **hipblasLtCreate**(*hipblasLtHandle_t* *handle)

Create a hipblaslt handle.

This function initializes the hipBLASLt library and creates a handle to an opaque structure holding the hipBLASLt library context. It allocates light hardware resources on the host and device, and must be called prior to making any other hipBLASLt library calls. The hipBLASLt library context is tied to the current ROCm device. To use the library on multiple devices, one hipBLASLt handle should be created for each device.

Parameters

handle – [out] Pointer to the allocated hipBLASLt handle for the created hipBLASLt context.

Return values

- **HIPBLAS_STATUS_SUCCESS** – The allocation completed successfully.
- **HIPBLAS_STATUS_INVALID_VALUE** – handle == NULL.

3.2 hipblasLtDestroy()

hipblasStatus_t **hipblasLtDestroy**(const *hipblasLtHandle_t* handle)

Destory a hipblaslt handle.

This function releases hardware resources used by the hipBLASLt library. This function is usually the last call with a particular handle to the hipBLASLt library. Because *hipblasLtCreate()* allocates some internal resources and the release of those resources by calling *hipblasLtDestroy()* will implicitly call *hipDeviceSynchronize()*, it is recommended to minimize the number of *hipblasLtCreate()/hipblasLtDestroy()* occurrences.

Parameters

handle – [in] Pointer to the hipBLASLt handle to be destroyed.

Return values

- **HIPBLAS_STATUS_SUCCESS** – The hipBLASLt context was successfully destroyed.
- **HIPBLAS_STATUS_NOT_INITIALIZED** – The hipBLASLt library was not initialized.
- **HIPBLAS_STATUS_INVALID_VALUE** – handle == NULL.

3.3 hipblasLtMatrixLayoutCreate()

hipblasStatus_t **hipblasLtMatrixLayoutCreate**(*hipblasLtMatrixLayout_t* *matLayout, hipDataType type, uint64_t rows, uint64_t cols, int64_t ld)

Create a matrix layout descriptor.

This function creates a matrix layout descriptor by allocating the memory needed to hold its opaque structure.

Parameters

- **matLayout** – [out] Pointer to the structure holding the matrix layout descriptor created by this function. see *hipblasLtMatrixLayout_t* .
- **type** – [in] Enumerant that specifies the data precision for the matrix layout descriptor this function creates. See hipDataType.
- **rows** – [in] Number of rows of the matrix.
- **cols** – [in] Number of columns of the matrix.
- **ld** – [in] The leading dimension of the matrix. In column major layout, this is the number of elements to jump to reach the next column. Thus ld >= m (number of rows).

Return values

- **HIPBLAS_STATUS_SUCCESS** – If the descriptor was created successfully.
- **HIPBLAS_STATUS_ALLOC_FAILED** – If the memory could not be allocated.

3.4 hipblasLtMatrixLayoutDestroy()

hipblasStatus_t **hipblasLtMatrixLayoutDestroy**(const *hipblasLtMatrixLayout_t* matLayout)

Destory a matrix layout descriptor.

This function destroys a previously created matrix layout descriptor object.

Parameters

matLayout – [in] Pointer to the structure holding the matrix layout descriptor that should be destroyed by this function. see *hipblasLtMatrixLayout_t* .

Return values

HIPBLAS_STATUS_SUCCESS – If the operation was successful.

3.5 hipblasLtMatrixLayoutSetAttribute()

hipblasStatus_t **hipblasLtMatrixLayoutSetAttribute**(*hipblasLtMatrixLayout_t* matLayout, *hipblasLtMatrixLayoutAttribute_t* attr, const void *buf, size_t sizeInBytes)

Set attribute to a matrix descriptor.

This function sets the value of the specified attribute belonging to a previously created matrix descriptor.

Parameters

- **matLayout** – [in] Pointer to the previously created structure holding the matrix mdescriptor queried by this function. See *hipblasLtMatrixLayout_t*.
- **attr** – [in] The attribute that will be set by this function. See *hipblasLtMatrixLayoutAttribute_t*.

- **buf** – [in] The value to which the specified attribute should be set.
- **sizeInBytes** – [in] Size of buf buffer (in bytes) for verification.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If the attribute was set successfully..
- **HIPBLAS_STATUS_INVALID_VALUE** – If buf is NULL or sizeInBytes doesn't match the size of the internal storage for the selected attribute.

3.6 hipblasLtMatrixLayoutGetAttribute()

hipblasStatus_t **hipblasLtMatrixLayoutGetAttribute**(*hipblasLtMatrixLayout_t* matLayout, *hipblasLtMatrixLayoutAttribute_t* attr, void *buf, size_t sizeInBytes, size_t *sizeWritten)

Query attribute from a matrix descriptor.

This function returns the value of the queried attribute belonging to a previously created matrix descriptor.

Parameters

- **matLayout** – [in] Pointer to the previously created structure holding the matrix descriptor queried by this function. See *hipblasLtMatrixLayout_t*.
- **attr** – [in] The attribute that will be retrieved by this function. See *hipblasLtMatrixLayoutAttribute_t*.
- **buf** – [out] Memory address containing the attribute value retrieved by this function.
- **sizeInBytes** – [in] Size of buf buffer (in bytes) for verification.
- **sizeWritten** – [out] Valid only when the return value is HIPBLAS_STATUS_SUCCESS. If sizeInBytes is non-zero: then sizeWritten is the number of bytes actually written; if sizeInBytes is 0: then sizeWritten is the number of bytes needed to write full contents.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If attribute's value was successfully written to user memory.
- **HIPBLAS_STATUS_INVALID_VALUE** – If sizeInBytes is 0 and sizeWritten is NULL, or if sizeInBytes is non-zero and buf is NULL, or sizeInBytes doesn't match size of internal storage for the selected attribute.

3.7 hipblasLtMatmulDescCreate()

hipblasStatus_t **hipblasLtMatmulDescCreate**(*hipblasLtMatmulDesc_t* *matmulDesc, hipblasComputeType_t computeType, hipDataType scaleType)

Create a matrix multiply descriptor.

This function creates a matrix multiply descriptor by allocating the memory needed to hold its opaque structure.

Parameters

- **matmulDesc** – [out] Pointer to the structure holding the matrix multiply descriptor created by this function. See *hipblasLtMatmulDesc_t*.
- **computeType** – [in] Enumerant that specifies the data precision for the matrix multiply descriptor this function creates. See *hipblasComputeType_t*.

- **scaleType** – [in] Enumerant that specifies the data precision for the matrix transform descriptor this function creates. See `hipDataType`.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If the descriptor was created successfully.
- **HIPBLAS_STATUS_ALLOC_FAILED** – If the memory could not be allocated.

3.8 hipblasLtMatmulDescDestroy()

`hipblasStatus_t hipblasLtMatmulDescDestroy`(const *hipblasLtMatmulDesc_t* matmulDesc)

Destory a matrix multiply descriptor.

This function destroys a previously created matrix multiply descriptor object.

Parameters

matmulDesc – [in] Pointer to the structure holding the matrix multiply descriptor that should be destroyed by this function. See *hipblasLtMatmulDesc_t*.

Return values

HIPBLAS_STATUS_SUCCESS – If operation was successful.

3.9 hipblasLtMatmulDescSetAttribute()

`hipblasStatus_t hipblasLtMatmulDescSetAttribute`(*hipblasLtMatmulDesc_t* matmulDesc,
hipblasLtMatmulDescAttributes_t attr, const void *buf,
size_t sizeInBytes)

Set attribute to a matrix multiply descriptor.

This function sets the value of the specified attribute belonging to a previously created matrix multiply descriptor.

Parameters

- **matmulDesc** – [in] Pointer to the previously created structure holding the matrix multiply descriptor queried by this function. See *hipblasLtMatmulDesc_t*.
- **attr** – [in] The attribute that will be set by this function. See *hipblasLtMatmulDescAttributes_t*.
- **buf** – [in] The value to which the specified attribute should be set.
- **sizeInBytes** – [in] Size of buf buffer (in bytes) for verification.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If the attribute was set successfully..
- **HIPBLAS_STATUS_INVALID_VALUE** – If buf is NULL or sizeInBytes doesn't match the size of the internal storage for the selected attribute.

3.10 hipblasLtMatmulDescGetAttribute()

`hipblasStatus_t hipblasLtMatmulDescGetAttribute`(*hipblasLtMatmulDesc_t* matmulDesc,
hipblasLtMatmulDescAttributes_t attr, void *buf, size_t
sizeInBytes, size_t *sizeWritten)

Query attribute from a matrix multiply descriptor.

This function returns the value of the queried attribute belonging to a previously created matrix multiply descriptor.

Parameters

- **matmulDesc** – [in] Pointer to the previously created structure holding the matrix multiply descriptor queried by this function. See *hipblasLtMatmulDesc_t*.
- **attr** – [in] The attribute that will be retrieved by this function. See *hipblasLtMatmulDescAttributes_t*.
- **buf** – [out] Memory address containing the attribute value retrieved by this function.
- **sizeInBytes** – [in] Size of buf buffer (in bytes) for verification.
- **sizeWritten** – [out] Valid only when the return value is HIPBLAS_STATUS_SUCCESS. If sizeInBytes is non-zero: then sizeWritten is the number of bytes actually written; if sizeInBytes is 0: then sizeWritten is the number of bytes needed to write full contents.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If attribute’s value was successfully written to user memory.
- **HIPBLAS_STATUS_INVALID_VALUE** – If sizeInBytes is 0 and sizeWritten is NULL, or if sizeInBytes is non-zero and buf is NULL, or sizeInBytes doesn’t match size of internal storage for the selected attribute.

3.11 hipblasLtMatmulPreferenceCreate()

hipblasStatus_t **hipblasLtMatmulPreferenceCreate**(*hipblasLtMatmulPreference_t* *pref)

Create a preference descriptor.

This function creates a matrix multiply heuristic search preferences descriptor by allocating the memory needed to hold its opaque structure.

Parameters

- **pref** – [out] Pointer to the structure holding the matrix multiply preferences descriptor created by this function. see *hipblasLtMatmulPreference_t*.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If the descriptor was created successfully.
- **HIPBLAS_STATUS_ALLOC_FAILED** – If memory could not be allocated.

3.12 hipblasLtMatmulPreferenceDestroy()

hipblasStatus_t **hipblasLtMatmulPreferenceDestroy**(const *hipblasLtMatmulPreference_t* pref)

Destory a preferences descriptor.

This function destroys a previously created matrix multiply preferences descriptor object.

Parameters

- **pref** – [in] Pointer to the structure holding the matrix multiply preferences descriptor that should be destroyed by this function. See *hipblasLtMatmulPreference_t*.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If operation was successful.

3.13 hipblasLtMatmulPreferenceSetAttribute()

hipblasStatus_t **hipblasLtMatmulPreferenceSetAttribute**(*hipblasLtMatmulPreference_t* pref, *hipblasLtMatmulPreferenceAttributes_t* attr, const void *buf, size_t sizeInBytes)

Set attribute to a preference descriptor.

This function sets the value of the specified attribute belonging to a previously created matrix multiply preferences descriptor.

Parameters

- **pref** – [in] Pointer to the previously created structure holding the matrix multiply preferences descriptor queried by this function. See *hipblasLtMatmulPreference_t*.
- **attr** – [in] The attribute that will be set by this function. See *hipblasLtMatmulPreferenceAttributes_t*.
- **buf** – [in] The value to which the specified attribute should be set.
- **sizeInBytes** – [in] Size of buf buffer (in bytes) for verification.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If the attribute was set successfully..
- **HIPBLAS_STATUS_INVALID_VALUE** – If buf is NULL or sizeInBytes doesn't match the size of the internal storage for the selected attribute.

3.14 hipblasLtMatmulPreferenceGetAttribute()

hipblasStatus_t **hipblasLtMatmulPreferenceGetAttribute**(*hipblasLtMatmulPreference_t* pref, *hipblasLtMatmulPreferenceAttributes_t* attr, void *buf, size_t sizeInBytes, size_t *sizeWritten)

Query attribute from a preference descriptor.

This function returns the value of the queried attribute belonging to a previously created matrix multiply heuristic search preferences descriptor.

Parameters

- **pref** – [in] Pointer to the previously created structure holding the matrix multiply heuristic search preferences descriptor queried by this function. See *hipblasLtMatmulPreference_t*.
- **attr** – [in] The attribute that will be retrieved by this function. See *hipblasLtMatmulPreferenceAttributes_t*.
- **buf** – [out] Memory address containing the attribute value retrieved by this function.
- **sizeInBytes** – [in] Size of buf buffer (in bytes) for verification.
- **sizeWritten** – [out] Valid only when the return value is HIPBLAS_STATUS_SUCCESS. If sizeInBytes is non-zero: then sizeWritten is the number of bytes actually written; if sizeInBytes is 0: then sizeWritten is the number of bytes needed to write full contents.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If attribute's value was successfully written to user memory.
- **HIPBLAS_STATUS_INVALID_VALUE** – If sizeInBytes is 0 and sizeWritten is NULL, or if sizeInBytes is non-zero and buf is NULL, or sizeInBytes doesn't match size of internal storage for the selected attribute.

3.15 hipblasLtMatmulAlgoGetHeuristic()

```
hipblasStatus_t hipblasLtMatmulAlgoGetHeuristic(hipblasLtHandle_t handle, hipblasLtMatmulDesc_t
matmulDesc, hipblasLtMatrixLayout_t Adesc,
hipblasLtMatrixLayout_t Bdesc, hipblasLtMatrixLayout_t
Cdesc, hipblasLtMatrixLayout_t Ddesc,
hipblasLtMatmulPreference_t pref, int
requestedAlgoCount, hipblasLtMatmulHeuristicResult_t
heuristicResultsArray[], int *returnAlgoCount)
```

Retrieve the possible algorithms.

This function retrieves the possible algorithms for the matrix multiply operation `hipblasLtMatmul()` function with the given input matrices A, B and C, and the output matrix D. The output is placed in `heuristicResultsArray[]` in the order of increasing estimated compute time. Note that the wall duration increases if the `requestedAlgoCount` increases.

Parameters

- **handle** – [in] Pointer to the allocated hipBLASLt handle for the hipBLASLt context. See `hipblasLtHandle_t`.
- **matmulDesc** – [in] Handle to a previously created matrix multiplication descriptor of type `hipblasLtMatmulDesc_t`.
- **Adesc, Bdesc, Cdesc, Ddesc** – [in] Handles to the previously created matrix layout descriptors of the type `hipblasLtMatrixLayout_t`.
- **pref** – [in] Pointer to the structure holding the heuristic search preferences descriptor. See `hipblasLtMatmulPreference_t`.
- **requestedAlgoCount** – [in] Size of the `heuristicResultsArray` (in elements). This is the requested maximum number of algorithms to return.
- **heuristicResultsArray[]** – [out] Array containing the algorithm heuristics and associated runtime characteristics, returned by this function, in the order of increasing estimated compute time.
- **returnAlgoCount** – [out] Number of algorithms returned by this function. This is the number of `heuristicResultsArray` elements written.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If query was successful. Inspect `heuristicResultsArray[0 to (returnAlgoCount - 1)].state` for the status of the results.
- **HIPBLAS_STATUS_NOT_SUPPORTED** – If no heuristic function available for current configuration.
- **HIPBLAS_STATUS_INVALID_VALUE** – If `requestedAlgoCount` is less or equal to zero.

3.16 hipblasLtMatmul()

```
hipblasStatus_t hipblasLtMatmul(hipblasLtHandle_t handle, hipblasLtMatmulDesc_t matmulDesc, const void
*alpha, const void *A, hipblasLtMatrixLayout_t Adesc, const void *B,
hipblasLtMatrixLayout_t Bdesc, const void *beta, const void *C,
hipblasLtMatrixLayout_t Cdesc, void *D, hipblasLtMatrixLayout_t Ddesc,
const hipblasLtMatmulAlgo_t *algo, void *workspace, size_t
workspaceSizeInBytes, hipStream_t stream)
```

Retrieve the possible algorithms.

This function computes the matrix multiplication of matrices A and B to produce the output matrix D, according to the following operation: $D = \alpha * (A * B) + \beta * (C)$, where A, B, and C are input matrices, and alpha and beta are input scalars. Note: This function supports both in-place matrix multiplication ($C == D$ and $Cdesc == Ddesc$) and out-of-place matrix multiplication ($C != D$, both matrices must have the same data type, number of rows, number of columns, batch size, and memory order). In the out-of-place case, the leading dimension of C can be different from the leading dimension of D. Specifically the leading dimension of C can be 0 to achieve row or column broadcast. If Cdesc is omitted, this function assumes it to be equal to Ddesc.

Parameters

- **handle** – [in] Pointer to the allocated hipBLASLt handle for the hipBLASLt context. See [hipblasLtHandle_t](#).
- **matmulDesc** – [in] Handle to a previously created matrix multiplication descriptor of type [hipblasLtMatmulDesc_t](#).
- **alpha, beta** – [in] Pointers to the scalars used in the multiplication.
- **Adesc, Bdesc, Cdesc, Ddesc** – [in] Handles to the previously created matrix layout descriptors of the type [hipblasLtMatrixLayout_t](#).
- **A, B, C** – [in] Pointers to the GPU memory associated with the corresponding descriptors Adesc, Bdesc and Cdesc.
- **D** – [out] Pointer to the GPU memory associated with the descriptor Ddesc.
- **algo** – [in] Handle for matrix multiplication algorithm to be used. See [hipblasLtMatmulAlgo_t](#). When NULL, an implicit heuristics query with default search preferences will be performed to determine actual algorithm to use.
- **workspace** – [in] Pointer to the workspace buffer allocated in the GPU memory. Pointer must be 16B aligned (that is, lowest 4 bits of address must be 0).
- **workspaceSizeInBytes** – [in] Size of the workspace.
- **stream** – [in] The HIP stream where all the GPU work will be submitted.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If the operation completed successfully.
- **HIPBLAS_STATUS_EXECUTION_FAILED** – If HIP reported an execution error from the device.
- **HIPBLAS_STATUS_ARCH_MISMATCH** – If the configured operation cannot be run using the selected device.
- **HIPBLAS_STATUS_NOT_SUPPORTED** – If the current implementation on the selected device doesn't support the configured operation.
- **HIPBLAS_STATUS_INVALID_VALUE** – If the parameters are unexpectedly NULL, in conflict or in an impossible configuration. For example, when workspaceSizeInBytes is less than workspace required by the configured algo.
- **HIPBLAS_STATUS_NOT_INITIALIZED** – If hipBLASLt handle has not been initialized.

Datatypes Supported:

hipblasLtMatmul supports the following computeType, scaleType, Atype/Btype, Ctype/Dtype and Bias Type:

computeType	scaleType/Bias Type	Atype/Btype	Ctype/Dtype
HIPBLAS_COMPUTE_32F	HIP_R_32F	HIP_R_32F	HIP_R_32F
HIPBLAS_COMPUTE_32F_FAST_TF32	HIP_R_32F	HIP_R_32F	HIP_R_32F
HIPBLAS_COMPUTE_32F	HIP_R_32F	HIP_R_16F	HIP_R_16F
HIPBLAS_COMPUTE_32F	HIP_R_32F	HIP_R_16F	HIP_R_32F
HIPBLAS_COMPUTE_32F	HIP_R_32F	HIP_R_16BF	HIP_R_16BF

For FP8 type Matmul, hipBLASLt supports the type combinations shown in the following table:

- This table uses simpler breviations:
 - **FP16** means **HIP_R_16F**
 - **BF16** means **HIP_R_16BF**
 - **FP32** means **HIP_R_32F**
 - **FP8** means **HIP_R_8F_E4M3_FNUZ** and
 - **BF8** means **HIP_R_8F_E5M2_FNUZ**
- This table applies to all transpose types (NN/NT/TT/TN)
- **Default Bias Type** means the type when users don't explicitly specify the bias type

Atype	Btype	Ctype	Dtype	compute-Type	scaleA,B	scaleC,D	Bias Type	Default Type	Bias	
FP8	FP8	FP16	FP16	FP32	Yes	No	FP32, FP16	FP16		
		BF16	BF16				FP32, BF16	BF16		
		FP32	FP32				FP32, BF16	BF16		
		FP8	FP8				Yes	FP32, FP16		FP16
		BF8	BF8				FP32, FP16	FP16		
		FP8	FP8				FP32, FP16	FP16		
	BF8	FP16	FP16	No	FP32, FP16	FP16				
		BF16	BF16	FP32, BF16	BF16					
		FP32	FP32	FP32, BF16	BF16					
		FP8	FP8	Yes	FP32, FP16	FP16				
		BF8	BF8	FP32, FP16	FP16					
		FP8	FP8	FP32, FP16	FP16					
BF8	FP8	FP16	FP16			No	FP32, FP16	FP16		
		BF16	BF16				FP32, BF16	BF16		
		FP32	FP32				FP32, BF16	BF16		
		FP8	FP8				Yes	FP32, FP16		FP16
		BF8	BF8				FP32, FP16	FP16		
		FP8	FP8				FP32, FP16	FP16		
	BF8	FP16	FP16	No	FP32, FP16	FP16				
		BF16	BF16	FP32, BF16	BF16					
		FP32	FP32	FP32, BF16	BF16					
		FP8	FP8	Yes	FP32, FP16	FP16				
		BF8	BF8	FP32, FP16	FP16					
		FP8	FP8	FP32, FP16	FP16					

3.17 hipblasLtMatrixTransformDescCreate()

hipblasStatus_t **hipblasLtMatrixTransformDescCreate**(hipblasLtMatrixTransformDesc_t *transformDesc, hipDataType scaleType)

Create new matrix transform operation descriptor.

Return values

- **HIPBLAS_STATUS_ALLOC_FAILED** – if memory could not be allocated
- **HIPBLAS_STATUS_SUCCESS** – if descriptor was created successfully

3.18 hipblasLtMatrixTransformDescDestroy()

hipblasStatus_t **hipblasLtMatrixTransformDescDestroy**(*hipblasLtMatrixTransformDesc_t* transformDesc)

Destroy matrix transform operation descriptor.

Return values

- **HIPBLAS_STATUS_SUCCESS** – if operation was successful

3.19 hipblasLtMatrixTransformDescSetAttribute()

hipblasStatus_t **hipblasLtMatrixTransformDescSetAttribute**(*hipblasLtMatrixTransformDesc_t* transformDesc,
hipblasLtMatrixTransformDescAttributes_t attr, const void *buf, size_t sizeInBytes)

Set matrix transform operation descriptor attribute.

Parameters

- **transformDesc** – [in] The descriptor
- **attr** – [in] The attribute
- **buf** – [in] memory address containing the new value
- **sizeInBytes** – [in] size of buf buffer for verification (in bytes)

Return values

- **HIPBLAS_STATUS_INVALID_VALUE** – if buf is NULL or sizeInBytes doesn't match size of internal storage for selected attribute
- **HIPBLAS_STATUS_SUCCESS** – if attribute was set successfully

3.20 hipblasLtMatrixTransformDescGetAttribute()

hipblasStatus_t **hipblasLtMatrixTransformDescGetAttribute**(*hipblasLtMatrixTransformDesc_t* transformDesc,
hipblasLtMatrixTransformDescAttributes_t attr, void *buf, size_t sizeInBytes, size_t *sizeWritten)

Matrix transform operation getter.

Get matrix transform operation descriptor attribute.

Parameters

- **transformDesc** – [in] The descriptor
- **attr** – [in] The attribute
- **buf** – [out] memory address containing the new value
- **sizeInBytes** – [in] size of buf buffer for verification (in bytes)

- **sizeWritten** – [out] only valid when return value is HIPBLAS_STATUS_SUCCESS. If sizeInBytes is non-zero: number of bytes actually written, if sizeInBytes is 0: number of bytes needed to write full contents

Return values

- **HIPBLAS_STATUS_INVALID_VALUE** – if sizeInBytes is 0 and sizeWritten is NULL, or if sizeInBytes is non-zero and buf is NULL or sizeInBytes doesn't match size of internal storage for selected attribute
- **HIPBLAS_STATUS_SUCCESS** – if attribute's value was successfully written to user memory

3.21 hipblasLtMatrixTransform()

hipblasStatus_t **hipblasLtMatrixTransform**(*hipblasLtHandle_t* lightHandle, *hipblasLtMatrixTransformDesc_t* transformDesc, const void *alpha, const void *A, *hipblasLtMatrixLayout_t* Adesc, const void *beta, const void *B, *hipblasLtMatrixLayout_t* Bdesc, void *C, *hipblasLtMatrixLayout_t* Cdesc, *hipStream_t* stream)

Matrix layout conversion helper.

Matrix layout conversion helper ($C = \alpha * \text{op}(A) + \beta * \text{op}(B)$), can be used to change memory order of data or to scale and shift the values.

Parameters

- **lightHandle** – [in] Pointer to the allocated hipBLASLt handle for the hipBLASLt context. See *hipblasLtHandle_t*.
- **transformDesc** – [in] Pointer to allocated matrix transform descriptor.
- **alpha** – [in] Pointer to scalar alpha, either pointer to host or device address.
- **A** – [in] Pointer to matrix A, must be pointer to device address.
- **Adesc** – [in] Pointer to layout for input matrix A.
- **beta** – [in] Pointer to scalar beta, either pointer to host or device address.
- **B** – [in] Pointer to layout for matrix B, must be pointer to device address
- **Bdesc** – [in] Pointer to layout for inputmatrix B.
- **C** – [in] Pointer to matrix C, must be pointer to device address
- **Cdesc** – [out] Pointer to layout for output matrix C.
- **stream** – [in] The HIP stream where all the GPU work will be submitted.

Return values

- **HIPBLAS_STATUS_NOT_INITIALIZED** – if hipBLASLt handle has not been initialized
- **HIPBLAS_STATUS_INVALID_VALUE** – if parameters are in conflict or in an impossible configuration; e.g. when A is not NULL, but Adesc is NULL
- **HIPBLAS_STATUS_NOT_SUPPORTED** – if current implementation on selected device doesn't support configured operation
- **HIPBLAS_STATUS_ARCH_MISMATCH** – if configured operation cannot be run using selected device
- **HIPBLAS_STATUS_EXECUTION_FAILED** – if HIP reported execution error from the device

- **HIPBLAS_STATUS_SUCCESS** – if the operation completed successfully

hipblasLtMatrixTransform supports the following Atype/Btype/Ctype and scaleType:

Atype/Btype/Ctype	scaleType
HIP_R_32F	HIP_R_32F
HIP_R_16F	HIP_R_32F/HIP_R_16F
HIP_R_16BF	HIP_R_32F
HIP_R_8I	HIP_R_32F
HIP_R_32I	HIP_R_32F

HIPBLASLTTEXT API REFERENCE

hipBLASLt has extension APIs with namespace `hipblaslt_ext`. It is only C++ compatible. The extensions support:

1. *Gemm*
2. *Grouped gemm*
3. *Get all algorithms*

4.1 hipBLASLtExt datatypes reference

4.1.1 GemmType

enum class `hipblaslt_ext::GemmType`

It is an enumerated type used to specific the type of the gemm problem in `hipblasLtExt` APIs.

Values:

enumerator `HIPBLASLT_GEMM`

enumerator `HIPBLASLT_GROUPED_GEMM`

4.1.2 GemmProblemType

struct `GemmProblemType`

`hipblasLt` extension `ProblemType` for gemm problems. (deprecated)

This structure sets the problem type of a gemm problem.

Public Members

`hipblasOperation_t op_a`

The A matrix transpose.

`hipblasOperation_t op_b`

The B matrix transpose.

`hipDataType type_a`

The A matrix datatype.

hipDataType **type_b**
The B matrix datatype.

hipDataType **type_c**
The C matrix datatype.

hipDataType **type_d**
The D matrix datatype.

hipblasComputeType_t **type_compute**
The compute datatype.

4.1.3 GemmProblemTypeV2

class **GemmProblemTypeV2**

hipblasLt extension ProblemType for gemm problems.

This structure sets the problem type of a gemm problem.

Public Functions

void **setOpA**(hipblasOperation_t op)
Set the A matrix transpose.

void **setOpB**(hipblasOperation_t op)
Set the B matrix transpose.

void **setTypeA**(hipDataType type)
Set the A matrix datatype.

void **setTypeB**(hipDataType type)
Set the B matrix datatype.

void **setTypeC**(hipDataType type)
Set the C matrix datatype.

void **setTypeD**(hipDataType type)
Set the D matrix datatype.

void **setTypeCompute**(hipblasComputeType_t type)
Set the compute datatype.

hipblasOperation_t **getOpA**() const
The A matrix transpose.

hipblasOperation_t **getOpB**() const
The B matrix transpose.

hipDataType **getTypeA**() const
The A matrix datatype.

hipDataType **getTypeB**() const
The B matrix datatype.

hipDataType **getTypeC**() const

The C matrix datatype.

hipDataType **getTypeD**() const

The D matrix datatype.

hipblasComputeType_t **getTypeCompute**() const

The compute datatype.

4.1.4 GemmEpilogue

struct **GemmEpilogue**

hipblasLt extension Epilogue for gemm problems. (Deprecated)

This structure sets the epilogue of a gemm problem.

Public Members

hipblasLtEpilogue_t **mode** = HIPBLASLT_EPILOGUE_DEFAULT

The mode of epilogue. Default is gemm.

hipDataType **bias_data_type** = HIPBLASLT_DATATYPE_INVALID

The bias datatype. Only works if mode is set to bias related epilogues.

int **aux_ld** = 0

The aux leading dimension. Only works if mode is set to aux related epilogues.

int **aux_stride** = 0

The aux batch stride. Only works if mode is set to aux related epilogues.

4.1.5 GemmEpilogueV2

class **GemmEpilogueV2**

hipblasLt extension Epilogue for gemm problems.

This class sets the epilogue of a gemm problem.

Public Functions

void **setMode**(*hipblasLtEpilogue_t* mode)

Set the mode of epilogue. Default is gemm.

void **setBiasDataType**(hipDataType biasDataType)

Set the bias datatype. Only works if mode is set to bias related epilogues.

void **setAuxLeadingDimension**(int auxLeadingDimension)

Set the aux leading dimension. Only works if mode is set to aux related epilogues.

void **setAuxBatchStride**(int auxBatchStride)

Set the aux batch stride. Only works if mode is set to aux related epilogues.

void **setScalingAType**(int scalingAType)

0 is scalar, 1 is vector. Only works if DataTypeA = DataTypeB = FP8.

void **setScalingBType**(int scalingBType)

0 is scalar, 1 is vector. Only works if DataTypeA = DataTypeB = FP8.

hipblasLtEpilogue_t **getMode**() const

The mode of epilogue. Default is gemm.

hipDataType **getBiasDataType**() const

The bias datatype. Only works if mode is set to bias related epilogues.

int **getAuxLeadingDimension**() const

The aux leading dimension. Only works if mode is set to aux related epilogues.

int **getAuxBatchStride**() const

The aux batch stride. Only works if mode is set to aux related epilogues.

int **getScalingAType**() const

0 is scalar, 1 is vector. Only works if DataTypeA = DataTypeB = FP8.

int **getScalingBType**() const

0 is scalar, 1 is vector. Only works if DataTypeA = DataTypeB = FP8.

4.1.6 GemmInputs

struct **GemmInputs**

hipblasLt extension Inputs for gemm problems. (Deprecated)

This structure sets the input pointers of a gemm problem.

Public Members

const void ***a** = nullptr

The a matrix input pointer.

const void ***b** = nullptr

The b matrix input pointer.

const void ***c** = nullptr

The c matrix input pointer.

const void ***d** = nullptr

The d matrix input pointer.

const void ***alpha** = nullptr

The alpha value.

const void ***beta** = nullptr

The beta value.

const void ***bias** = nullptr
The bias input pointer.

const void ***scaleA** = nullptr
The Scale A input pointer.

const void ***scaleB** = nullptr
The Scale B input pointer.

const void ***scaleC** = nullptr
The Scale C input pointer.

const void ***scaleD** = nullptr
The Scale D input pointer.

const void ***scaleAux** = nullptr
The Scale AUX input pointer.

const void ***scaleAlphaVec** = nullptr
The scaleAlpha vector input pointer.

const void ***aux** = nullptr
The aux input pointer.

4.1.7 GemmInputsV2

class **GemmInputsV2**

hipblasLt extension Inputs for gemm problems.

This class sets the input pointers of a gemm problem.

Public Functions

void **setA**(const void *a)
Set the a matrix input pointer.

void **setB**(const void *b)
Set the b matrix input pointer.

void **setC**(const void *c)
Set the c matrix input pointer.

void **setD**(const void *d)
Set the d matrix input pointer.

void **setAlpha**(const void *alpha)
Set the alpha value.

void **setBeta**(const void *beta)
Set the beta value.

void **setBias**(const void *bias)
Set the bias input pointer.

void **setScaleA**(const void *scaleA)
Set the Scale A input pointer.

void **setScaleB**(const void *scaleB)
Set the Scale B input pointer.

void **setScaleC**(const void *scaleC)
Set the Scale C input pointer.

void **setScaleD**(const void *scaleD)
Set the Scale D input pointer.

void **setScaleAux**(const void *scaleAux)
Set the Scale AUX input pointer.

void **setScaleAlphaVec**(const void *scaleAlphaVec)
Set the scaleAlpha vector input pointer.

void **setAux**(const void *aux)
Set the aux input pointer.

void **setAmaxD**(const void *amaxD)
Set the AmaxD input pointer.

const void ***getA**() const
The a matrix input pointer.

const void ***getB**() const
The b matrix input pointer.

const void ***getC**() const
The c matrix input pointer.

const void ***getD**() const
The d matrix input pointer.

const void ***getAlpha**() const
The alpha value.

const void ***getBeta**() const
The beta value.

const void ***getBias**() const
The bias input pointer.

const void ***getScaleA**() const
The Scale A input pointer.

const void ***getScaleB**() const
The Scale B input pointer.

const void ***getScaleC**() const
The Scale C input pointer.

const void ***getScaleD()** const
The Scale D input pointer.

const void ***getScaleAux()** const
The Scale AUX input pointer.

const void ***getScaleAlphaVec()** const
The scaleAlpha vector input pointer.

const void ***getAux()** const
The aux input pointer.

const void ***getAmaxD()** const
The AmaxD input pointer.

4.2 hipBLASLtExt Gemm Class Reference

4.2.1 GemmPreference

class **GemmPreference**
hipblasLt extension preference for gemm problems. (deprecated)
Currently only supports setting max workspace size.

Public Functions

void **setMaxWorkspaceBytes**(size_t workspaceBytes)
This function sets the max workspace size.

Parameters

workspaceBytes – [in] Set the max workspace size in bytes.

const size_t **getMaxWorkspaceBytes()** const
This function returns the set max workspace size.

Return values

size_t – Returns the set max workspace size.

4.2.2 GemmPreferenceV2

class **GemmPreferenceV2**
hipblasLt extension preference for gemm problems.
Currently only supports setting max workspace size.

Public Functions

void **setMaxWorkspaceBytes**(size_t workspaceBytes)
This function sets the max workspace size.

Parameters

workspaceBytes – [in] Set the max workspace size in bytes.

```
const size_t getMaxWorkspaceBytes() const
```

This function returns the set max workspace size.

Return values

size_t – Returns the set max workspace size.

4.2.3 GemmInstance

class **GemmInstance**

hipblasLt extension instance for gemm problems.

Subclassed by *hipblaslt_ext::Gemm*, *hipblaslt_ext::GroupedGemm*

Public Functions

```
hipblasStatus_t algoGetHeuristic(const int requestedAlgoCount, const GemmPreference &pref,  
                                std::vector<hipblasLtMatmulHeuristicResult_t> &heuristicResults)
```

Retrieve the possible algorithms (deprecated)

This function retrieves the possible algorithms for the matrix multiply operation *hipblasLtMatmul()* function with the given data and compute type. The output is placed in heuristicResult in the order of increasing estimated compute time. Note that the wall duration increases if the requestedAlgoCount increases.

Parameters

- **requestedAlgoCount** – [in] number of requested algorithms.
- **pref** – [in] hipblasLt extension preference for gemm problems.
- **heuristicResults** – [out] The algorithm heuristic vector.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If query was successful. Inspect heuristicResults.size > 0, but may heuristicResults.size < requestedAlgoCount state for the status of the results.
- **HIPBLAS_STATUS_NOT_SUPPORTED** – If no heuristic function available for current configuration.
- **HIPBLAS_STATUS_INVALID_VALUE** – If no solution is found.

```
hipblasStatus_t algoGetHeuristic(const int requestedAlgoCount, const GemmPreferenceV2 &pref,  
                                std::vector<hipblasLtMatmulHeuristicResult_t> &heuristicResults)
```

Retrieve the possible algorithms.

This function retrieves the possible algorithms for the matrix multiply operation *hipblasLtMatmul()* function with the given data and compute type. The output is placed in heuristicResult in the order of increasing estimated compute time.

Parameters

- **requestedAlgoCount** – [in] number of requested algorithms.
- **pref** – [in] hipblasLt extension preference for gemm problems.
- **heuristicResults** – [out] The algorithm heuristic vector.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If query was successful. Inspect heuristicResults.size > 0, but may heuristicResults.size < requestedAlgoCount state for the status of the results.

- **HIPBLAS_STATUS_NOT_SUPPORTED** – If no heuristic function available for current configuration.
- **HIPBLAS_STATUS_INVALID_VALUE** – If no solution is found.

`hipblasStatus_t isAlgoSupported(hipblasLtMatmulAlgo_t &algo, size_t &workspaceSizeInBytes)`

Check if the algorithm supports the problem. (For hipblaslt extension API)

This function updates the problem saved inside the algorithm if the problem is supported. The required `workspaceSizeInBytes` is also returned.

Parameters

- **algo** – [in] The algorithm heuristic.
- **workspaceSizeInBytes** – [out] Return the required workspace size.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If query was successful. The problem is supported by the algorithm. results.
- **HIPBLAS_STATUS_INVALID_VALUE** – The problem is not supported.

`hipblasStatus_t isAlgoSupported(hipblasLtMatmulAlgo_t &algo, GemmTuning &tuning, size_t &workspaceSizeInBytes)`

Check if the algorithm supports the problem. (deprecated)

This function updates the problem saved inside the algorithm if the problem is supported. The required `workspaceSizeInBytes` is also returned.

Parameters

- **algo** – [in] The algorithm heuristic.
- **tuning** – [in] The tuning parameters.
- **workspaceSizeInBytes** – [out] Return the required workspace size.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If query was successful. The problem is supported by the algorithm. results.
- **HIPBLAS_STATUS_INVALID_VALUE** – The problem is not supported.

`hipblasStatus_t isAlgoSupported(hipblasLtMatmulAlgo_t &algo, GemmTuningV2 &tuning, size_t &workspaceSizeInBytes)`

Check if the algorithm supports the problem. (For hipblaslt extension API)

This function updates the problem saved inside the algorithm if the problem is supported. The required `workspaceSizeInBytes` is also returned.

Parameters

- **algo** – [in] The algorithm heuristic.
- **tuning** – [in] The tuning parameters.
- **workspaceSizeInBytes** – [out] Return the required workspace size.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If query was successful. The problem is supported by the algorithm. results.
- **HIPBLAS_STATUS_INVALID_VALUE** – The problem is not supported.

hipblasStatus_t **initialize**(const *hipblasLtMatmulAlgo_t* &algo, void *workspace, bool useUserArgs = true, hipStream_t stream = 0)

Create kernel arguments from a given *hipblaslt_ext::GemmInstance*.

This function creates kernel arguments from a given *hipblaslt_ext::GemmInstance* then saves the arguments inside the instance.

Parameters

- **algo** – [in] Handle for matrix multiplication algorithm to be used. See `hipblaslt.h:hipblasLtMatmulAlgo_t`. When NULL, an implicit heuristics query with default search preferences will be performed to determine actual algorithm to use.
- **workspace** – [in] Pointer to the workspace buffer allocated in the GPU memory. Pointer must be 16B aligned (that is, lowest 4 bits of address must be 0).
- **useUserArgs** – [in] Use user args, this does not affect vanilla gemm. (May be deprecated in the future)
- **stream** – [in] The HIP stream where all the GPU work will be submitted. (May be deprecated in the future)

Return values

- **HIPBLAS_STATUS_SUCCESS** – If the operation completed successfully.
- **HIPBLAS_STATUS_INVALID_VALUE** – If the `gemm_count = 0`.

hipblasStatus_t **initialize**(const *hipblasLtMatmulAlgo_t* &algo, GemmTuning &tuning, void *workspace, bool useUserArgs = true, hipStream_t stream = 0)

Create kernel arguments from a given *hipblaslt_ext::GemmInstance*. (deprecated)

This function creates kernel arguments from a given *hipblaslt_ext::GemmInstance* then saves the arguments inside the instance.

Parameters

- **algo** – [in] Handle for matrix multiplication algorithm to be used. See `hipblaslt.h:hipblasLtMatmulAlgo_t`. When NULL, an implicit heuristics query with default search preferences will be performed to determine actual algorithm to use.
- **tuning** – [in] Structure with user tuning parameters. Note that not every algo supports user tuning parameters. Will return `HIPBLAS_STATUS_INVALID_VALUE` if not supported. (be 0).
- **workspace** – [in] Pointer to the workspace buffer allocated in the GPU memory. Pointer must be 16B aligned (that is, lowest 4 bits of address must be 0).
- **useUserArgs** – [in] Use user args, this does not affect vanilla gemm. (May be deprecated in the future)
- **stream** – [in] The HIP stream where all the GPU work will be submitted. (May be deprecated in the future)

Return values

- **HIPBLAS_STATUS_SUCCESS** – If the operation completed successfully.
- **HIPBLAS_STATUS_INVALID_VALUE** – If the `gemm_count = 0`.

hipblasStatus_t **initialize**(const *hipblasLtMatmulAlgo_t* &algo, GemmTuningV2 &tuning, void *workspace, bool useUserArgs = true, hipStream_t stream = 0)

Create kernel arguments from a given *hipblaslt_ext::GemmInstance*.

This function creates kernel arguments from a given *hipblaslt_ext::GemmInstance* then saves the arguments inside the instance.

Parameters

- **algo** – [in] Handle for matrix multiplication algorithm to be used. See `hipblasLtMatmulAlgo_t`. When NULL, an implicit heuristics query with default search preferences will be performed to determine actual algorithm to use.
- **tuning** – [in] Structure with user tuning parameters. Note that not every algo supports user tuning parameters. Will return `HIPBLAS_STATUS_INVALID_VALUE` if not supported. (be 0).
- **workspace** – [in] Pointer to the workspace buffer allocated in the GPU memory. Pointer must be 16B aligned (that is, lowest 4 bits of address must be 0).
- **useUserArgs** – [in] Use user args, this does not affect vanilla gemm. (May be deprecated in the future)
- **stream** – [in] The HIP stream where all the GPU work will be submitted. (May be deprecated in the future)

Return values

- `HIPBLAS_STATUS_SUCCESS` – If the operation completed successfully.
- `HIPBLAS_STATUS_INVALID_VALUE` – If the `gemm_count = 0`.

`hipblasStatus_t run(hipStream_t stream, hipEvent_t start = nullptr, hipEvent_t stop = nullptr)`

Execute the kernel arguments stored inside the *hipblaslt_ext::GemmInstance*.

Parameters

- **stream** – [in] The HIP stream where all the GPU work will be
- **start** – [in] The HIP event which will record the start of the kernel
- **stop** – [in] The HIP event which will record the end of the kernel submitted.

Return values

- `HIPBLAS_STATUS_SUCCESS` – If the operation completed successfully.

Protected Functions

explicit `GemmInstance(hipblasLtHandle_t handle, GemmType type)`

Constructor of *GemmInstance*.

4.2.4 Gemm

class `Gemm` : public *hipblaslt_ext::GemmInstance*

hipblasLt extension instance for gemm.

The instance can be used to create arguments to compute the matrix multiplication of matrices A and B to produce the output matrix D, according to the following operation: $D = \alpha * (A * B) + \beta * C$, where A, B, and C are input matrices, and alpha and beta are input scalars.

Public Functions

explicit **Gemm**(*hipblasLtHandle_t* handle, *hipblasOperation_t* opA, *hipblasOperation_t* opB, *hipDataType* typeA, *hipDataType* typeB, *hipDataType* typeC, *hipDataType* typeD, *hipblasComputeType_t* typeCompute)

Constructor.

This function set the problem from *hipblasLt* structures. For more information about the structures, see *hipblasLtMatmul* for more information.

Parameters

- **handle** – [in] The handle from *hipBLASLt*.
- **opA, opB** – [in] The transpose type of matrix A, B
- **typeA, typeB, typeC, typeD** – [in] The data type of matrix A, B, C, D
- **typeCompute** – [in] The compute type of the gemm problem

explicit **Gemm**(*hipblasLtHandle_t* handle, *hipblasLtMatmulDesc_t* matmul_descr, const void *alpha, const void *A, *hipblasLtMatrixLayout_t* matA, const void *B, *hipblasLtMatrixLayout_t* matB, const void *beta, const void *C, *hipblasLtMatrixLayout_t* matC, void *D, *hipblasLtMatrixLayout_t* matD)

Constructor that sets the gemm problem from *hipblasLt* structures.

This constructor sets the problem from *hipblasLt* structures. For more information about the structures, see *hipblasLtMatmul* for more information.

Parameters

- **handle** – [in] The handle from *hipBLASLt*.
- **matmul_descr** – [in] Handle to a previously created matrix multiplication descriptor of type *hipblasLtMatmulDesc_t*.
- **alpha, beta** – [in] Pointers to the scalars used in the multiplication.
- **matA, matB, matC, matD** – [in] Handles to the previously created matrix layout descriptors of the type *hipblasLtMatrixLayout_t*.
- **A, B, C** – [in] Pointers to the GPU memory associated with the corresponding descriptors *matA*, *matB* and *matC*.
- **D** – [out] Pointer to the GPU memory associated with the descriptor *matD*.

hipblasStatus_t **setProblem**(*int64_t* m, *int64_t* n, *int64_t* k, *int64_t* batch_count, *GemmEpilogue* &epilogue, *GemmInputs* &inputs)

Sets the problem for a gemm problem. (Deprecated)

This function sets the problem with m, n, k, batch_count. It uses the problem type sets from the constructor.

Parameters

- **m, n, k** – [in] The problem size.
- **batch_count** – [in] The batch count.
- **epilogue** – [in] The structure that controls the epilogue.
- **inputs** – [in] The inputs of the problem.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If the operation completed successfully.

- **HIPBLAS_STATUS_EXECUTION_FAILED** – If HIP reported an execution error from the device.
- **HIPBLAS_STATUS_ARCH_MISMATCH** – If the configured operation cannot be run using the selected device.
- **HIPBLAS_STATUS_NOT_SUPPORTED** – If the current implementation on the selected device doesn't support the configured operation.
- **HIPBLAS_STATUS_INVALID_VALUE** – If the parameters are unexpectedly NULL, in conflict or in an impossible configuration.
- **HIBLAS_STATUS_NOT_INITIALIZED** – If hipBLASLt handle has not been initialized.

hipblasStatus_t **setProblem**(int64_t m, int64_t n, int64_t k, int64_t batch_count, *GemmEpilogueV2* &epilogue, *GemmInputsV2* &inputs)

Sets the problem for a gemm problem.

This function sets the problem with m, n, k, batch_count. It uses the problem type sets from the constructor.

Parameters

- **m, n, k** – [in] The problem size.
- **batch_count** – [in] The batch count.
- **epilogue** – [in] The class that controls the epilogue.
- **inputs** – [in] The inputs of the problem.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If the operation completed successfully.
- **HIPBLAS_STATUS_EXECUTION_FAILED** – If HIP reported an execution error from the device.
- **HIPBLAS_STATUS_ARCH_MISMATCH** – If the configured operation cannot be run using the selected device.
- **HIPBLAS_STATUS_NOT_SUPPORTED** – If the current implementation on the selected device doesn't support the configured operation.
- **HIPBLAS_STATUS_INVALID_VALUE** – If the parameters are unexpectedly NULL, in conflict or in an impossible configuration.
- **HIBLAS_STATUS_NOT_INITIALIZED** – If hipBLASLt handle has not been initialized.

hipblasStatus_t **setProblem**(int64_t m, int64_t n, int64_t k, int64_t batch_count, int64_t lda, int64_t ldb, int64_t ldc, int64_t ldd, int64_t strideA, int64_t strideB, int64_t strideC, int64_t strideD, *GemmEpilogue* &epilogue, *GemmInputs* &inputs, *GemmProblemType* &problemtype)

Sets the problem for a gemm problem. (Deprecated)

This function sets the problem with m, n, k, batch_count. It uses the problem type sets from the constructor.

Parameters

- **m, n, k** – [in] The problem size.
- **batch_count** – [in] The batch count.
- **lda, ldb, ldc, ldd** – [in] The leading dimensions of the matrix.
- **strideA, strideB, strideC, strideD** – [in] The batch stride of the matrix.

- **epilogue** – [in] The structure that controls the epilogue.
- **inputs** – [in] The inputs of the problem.
- **problemtype** – [in] The structure that sets the problem type of a gemm problem.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If the operation completed successfully.
- **HIPBLAS_STATUS_EXECUTION_FAILED** – If HIP reported an execution error from the device.
- **HIPBLAS_STATUS_ARCH_MISMATCH** – If the configured operation cannot be run using the selected device.
- **HIPBLAS_STATUS_NOT_SUPPORTED** – If the current implementation on the selected device doesn't support the configured operation.
- **HIPBLAS_STATUS_INVALID_VALUE** – If the parameters are unexpectedly NULL, in conflict or in an impossible configuration.
- **HIBLAS_STATUS_NOT_INITIALIZED** – If hipBLASLt handle has not been initialized.

`hipblasStatus_t setProblem(int64_t m, int64_t n, int64_t k, int64_t batch_count, int64_t lda, int64_t ldb, int64_t ldc, int64_t ldd, int64_t strideA, int64_t strideB, int64_t strideC, int64_t strideD, GemmEpilogueV2 &epilogue, GemmInputsV2 &inputs, GemmProblemTypeV2 &problemtype)`

Sets the problem for a gemm problem.

This function sets the problem with m, n, k, batch_count. It uses the problem type sets from the constructor.

Parameters

- **m, n, k** – [in] The problem size.
- **batch_count** – [in] The batch count.
- **lda, ldb, ldc, ldd** – [in] The leading dimensions of the matrix.
- **strideA, strideB, strideC, strideD** – [in] The batch stride of the matrix.
- **epilogue** – [in] The structure that controls the epilogue.
- **inputs** – [in] The inputs of the problem.
- **problemtype** – [in] The structure that sets the problem type of a gemm problem.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If the operation completed successfully.
- **HIPBLAS_STATUS_EXECUTION_FAILED** – If HIP reported an execution error from the device.
- **HIPBLAS_STATUS_ARCH_MISMATCH** – If the configured operation cannot be run using the selected device.
- **HIPBLAS_STATUS_NOT_SUPPORTED** – If the current implementation on the selected device doesn't support the configured operation.
- **HIPBLAS_STATUS_INVALID_VALUE** – If the parameters are unexpectedly NULL, in conflict or in an impossible configuration.
- **HIBLAS_STATUS_NOT_INITIALIZED** – If hipBLASLt handle has not been initialized.

hipblasStatus_t **setProblem**(*hipblasLtMatmulDesc_t* matmul_desc, const void *alpha, const void *A, *hipblasLtMatrixLayout_t* matA, const void *B, *hipblasLtMatrixLayout_t* matB, const void *beta, const void *C, *hipblasLtMatrixLayout_t* matC, void *D, *hipblasLtMatrixLayout_t* matD)

Sets the gemm problem from hipblasLt structures.

This function sets the problem from hipblasLt structures. For more information about the structures, see hipblasLtMatmul for more information.

Parameters

- **matmul_desc** – [in] Handle to a previously created matrix multiplication descriptor of type *hipblasLtMatmulDesc_t* .
- **alpha, beta** – [in] Pointers to the scalars used in the multiplication.
- **matA, matB, matC, matD** – [in] Handles to the previously created matrix layout descriptors of the type *hipblasLtMatrixLayout_t* .
- **A, B, C** – [in] Pointers to the GPU memory associated with the corresponding descriptors matA, matB and matC .
- **D** – [out] Pointer to the GPU memory associated with the descriptor matD .

Return values

- **HIPBLAS_STATUS_SUCCESS** – If the operation completed successfully.
- **HIPBLAS_STATUS_EXECUTION_FAILED** – If HIP reported an execution error from the device.
- **HIPBLAS_STATUS_ARCH_MISMATCH** – If the configured operation cannot be run using the selected device.
- **HIPBLAS_STATUS_NOT_SUPPORTED** – If the current implementation on the selected device doesn't support the configured operation.
- **HIPBLAS_STATUS_INVALID_VALUE** – If the parameters are unexpectedly NULL, in conflict or in an impossible configuration.
- **HIPBLAS_STATUS_NOT_INITIALIZED** – If hipBLASLt handle has not been initialized.

4.2.5 GroupedGemm

class **GroupedGemm** : public hipblasLt_ext::GemmInstance

hipblasLt extension instance for grouped gemm.

The instance can be used to create arguments to compute the matrix multiplication of matrices A and B to produce the output matrix D, according to the following operation: $D = \alpha * (A * B) + \beta * (C)$, where A, B, and C are input matrices, and alpha and beta are input scalars.

Public Functions

explicit **GroupedGemm**(*hipblasLtHandle_t* handle, *hipblasOperation_t* opA, *hipblasOperation_t* opB, *hipDataType* typeA, *hipDataType* typeB, *hipDataType* typeC, *hipDataType* typeD, *hipblasComputeType_t* typeCompute)

Constructor.

This function set the problem from hipblasLt structures. For more information about the structures, see hipblasLtMatmul for more information.

Parameters

- **handle** – [in] The handle from hipBLASLt.
- **opA, opB** – [in] The transpose type of matrix A, B
- **typeA, typeB, typeC, typeD** – [in] The data type of matrix A, B, C, D
- **typeCompute** – [in] The compute type of the gemm problem

explicit **GroupedGemm**(*hipblasLtHandle_t* handle, std::vector<*hipblasLtMatmulDesc_t*> &matmul_descr, std::vector<void*> &alpha, std::vector<void*> &A, std::vector<*hipblasLtMatrixLayout_t*> &matA, std::vector<void*> &B, std::vector<*hipblasLtMatrixLayout_t*> &matB, std::vector<void*> &beta, std::vector<void*> &C, std::vector<*hipblasLtMatrixLayout_t*> &matC, std::vector<void*> &D, std::vector<*hipblasLtMatrixLayout_t*> &matD)

Constructor that sets the grouped gemm problem from hipblasLt structures.

This constructor sets the problem from hipblasLt structures. For more information about the structures, see `hipblasLtMatmul` for more information.

Parameters

- **handle** – [in] The handle from hipBLASLt.
- **matmul_descr** – [in] Vectors of handle to a previously created matrix multiplication descriptor of type *hipblasLtMatmulDesc_t*.
- **alpha, beta** – [in] Vectors of float used in the multiplication.
- **matA, matB, matC, matD** – [in] Vectors of handle to the previously created matrix layout descriptors of the type *hipblasLtMatrixLayout_t*.
- **A, B, C** – [in] Vectors of pointer to the GPU memory associated with the corresponding descriptors `matA`, `matB` and `matC`.
- **D** – [out] Vector of pointer to the GPU memory associated with the descriptor `matD`.

`hipblasStatus_t` **setProblem**(std::vector<int64_t> &m, std::vector<int64_t> &n, std::vector<int64_t> &k, std::vector<int64_t> &batch_count, std::vector<*GemmEpilogue*> &epilogue, std::vector<*GemmInputs*> &inputs)

Sets the problem for a gemm problem. (Deprecated)

This function sets the problem with `m`, `n`, `k`, `batch_count`. It uses the problem type sets from the constructor.

Parameters

- **m, n, k** – [in] The problem size in vector.
- **batch_count** – [in] The batch count in vector.
- **epilogue** – [in] The structure in vector that controls the epilogue.
- **inputs** – [in] The inputs in vector of the problem.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If the operation completed successfully.
- **HIPBLAS_STATUS_EXECUTION_FAILED** – If HIP reported an execution error from the device.
- **HIPBLAS_STATUS_ARCH_MISMATCH** – If the configured operation cannot be run using the selected device.

- **HIPBLAS_STATUS_NOT_SUPPORTED** – If the current implementation on the selected device doesn't support the configured operation.
- **HIPBLAS_STATUS_INVALID_VALUE** – If the parameters are unexpectedly NULL, in conflict or in an impossible configuration.
- **HIBLAS_STATUS_NOT_INITIALIZED** – If hipBLASLt handle has not been initialized.

hipblasStatus_t **setProblem**(std::vector<int64_t> &m, std::vector<int64_t> &n, std::vector<int64_t> &k, std::vector<int64_t> &batch_count, std::vector<GemmEpilogueV2> &epilogue, std::vector<GemmInputsV2> &inputs)

Sets the problem for a gemm problem.

This function sets the problem with m, n, k, batch_count. It uses the problem type sets from the constructor.

Parameters

- **m, n, k** – [in] The problem size in vector.
- **batch_count** – [in] The batch count in vector.
- **epilogue** – [in] The structure in vector that controls the epilogue.
- **inputs** – [in] The inputs in vector of the problem.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If the operation completed successfully.
- **HIPBLAS_STATUS_EXECUTION_FAILED** – If HIP reported an execution error from the device.
- **HIPBLAS_STATUS_ARCH_MISMATCH** – If the configured operation cannot be run using the selected device.
- **HIPBLAS_STATUS_NOT_SUPPORTED** – If the current implementation on the selected device doesn't support the configured operation.
- **HIPBLAS_STATUS_INVALID_VALUE** – If the parameters are unexpectedly NULL, in conflict or in an impossible configuration.
- **HIBLAS_STATUS_NOT_INITIALIZED** – If hipBLASLt handle has not been initialized.

hipblasStatus_t **setProblem**(std::vector<int64_t> &m, std::vector<int64_t> &n, std::vector<int64_t> &k, std::vector<int64_t> &batch_count, std::vector<int64_t> &l1da, std::vector<int64_t> &l1db, std::vector<int64_t> &l1dc, std::vector<int64_t> &l1dd, std::vector<int64_t> &strideA, std::vector<int64_t> &strideB, std::vector<int64_t> &strideC, std::vector<int64_t> &strideD, std::vector<GemmEpilogue> &epilogue, std::vector<GemmInputs> &inputs, GemmProblemType &problemtyp)

Sets the problem for a gemm problem. (Deprecated)

This function sets the problem with m, n, k, batch_count. It uses the problem type sets from the constructor.

Parameters

- **m, n, k** – [in] The problem size in vector.
- **batch_count** – [in] The batch count in vector.
- **lda, ldb, ldc, ldd** – [in] The leading dimensions in vector of the matrix.
- **strideA, strideB, strideC, strideD** – [in] The batch stride in vector of the matrix.
- **epilogue** – [in] The structure in vector that controls the epilogue.

- **inputs** – [in] The inputs in vector of the problem.
- **problemtype** – [in] The structure that sets the problem type of a gemm problem.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If the operation completed successfully.
- **HIPBLAS_STATUS_EXECUTION_FAILED** – If HIP reported an execution error from the device.
- **HIPBLAS_STATUS_ARCH_MISMATCH** – If the configured operation cannot be run using the selected device.
- **HIPBLAS_STATUS_NOT_SUPPORTED** – If the current implementation on the selected device doesn't support the configured operation.
- **HIPBLAS_STATUS_INVALID_VALUE** – If the parameters are unexpectedly NULL, in conflict or in an impossible configuration.
- **HIBLAS_STATUS_NOT_INITIALIZED** – If hipBLASLt handle has not been initialized.

```
hipblasStatus_t setProblem(std::vector<int64_t> &m, std::vector<int64_t> &n, std::vector<int64_t> &k,  
    std::vector<int64_t> &batch_count, std::vector<int64_t> &lda,  
    std::vector<int64_t> &ldb, std::vector<int64_t> &ldc, std::vector<int64_t>  
    &ldd, std::vector<int64_t> &strideA, std::vector<int64_t> &strideB,  
    std::vector<int64_t> &strideC, std::vector<int64_t> &strideD,  
    std::vector<GemmEpilogueV2> &epilogue, std::vector<GemmInputsV2>  
    &inputs, GemmProblemTypeV2 &problemtype)
```

Sets the problem for a gemm problem.

This function sets the problem with m, n, k, batch_count. It uses the problem type sets from the constructor.

Parameters

- **m, n, k** – [in] The problem size in vector.
- **batch_count** – [in] The batch count in vector.
- **lda, ldb, ldc, ldd** – [in] The leading dimensions in vector of the matrix.
- **strideA, strideB, strideC, strideD** – [in] The batch stride in vector of the matrix.
- **epilogue** – [in] The structure in vector that controls the epilogue.
- **inputs** – [in] The inputs in vector of the problem.
- **problemtype** – [in] The structure that sets the problem type of a gemm problem.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If the operation completed successfully.
- **HIPBLAS_STATUS_EXECUTION_FAILED** – If HIP reported an execution error from the device.
- **HIPBLAS_STATUS_ARCH_MISMATCH** – If the configured operation cannot be run using the selected device.
- **HIPBLAS_STATUS_NOT_SUPPORTED** – If the current implementation on the selected device doesn't support the configured operation.
- **HIPBLAS_STATUS_INVALID_VALUE** – If the parameters are unexpectedly NULL, in conflict or in an impossible configuration.
- **HIBLAS_STATUS_NOT_INITIALIZED** – If hipBLASLt handle has not been initialized.

```
hipblasStatus_t setProblem(std::vector<hipblasLtMatmulDesc_t> &matmul_descr, std::vector<void*>
    &alpha, std::vector<void*> &A, std::vector<hipblasLtMatrixLayout_t> &matA,
    std::vector<void*> &B, std::vector<hipblasLtMatrixLayout_t> &matB,
    std::vector<void*> &beta, std::vector<void*> &C,
    std::vector<hipblasLtMatrixLayout_t> &matC, std::vector<void*> &D,
    std::vector<hipblasLtMatrixLayout_t> &matD)
```

Sets the grouped gemm problem from hipblasLt structures.

This function sets the problem from hipblasLt structures. For more information about the structures, see hipblasLtMatmul for more information.

Parameters

- **matmul_descr** – [in] Vectors of handle to a previously created matrix multiplication descriptor of type *hipblasLtMatmulDesc_t*.
- **alpha, beta** – [in] Vectors of float used in the multiplication.
- **matA, matB, matC, matD** – [in] Vectors of handle to the previously created matrix layout descriptors of the type *hipblasLtMatrixLayout_t*.
- **A, B, C** – [in] Vectors of pointer to the GPU memory associated with the corresponding descriptors matA, matB and matC.
- **D** – [out] Vector of pointer to the GPU memory associated with the descriptor matD.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If the operation completed successfully.
- **HIPBLAS_STATUS_EXECUTION_FAILED** – If HIP reported an execution error from the device.
- **HIPBLAS_STATUS_ARCH_MISMATCH** – If the configured operation cannot be run using the selected device.
- **HIPBLAS_STATUS_NOT_SUPPORTED** – If the current implementation on the selected device doesn't support the configured operation.
- **HIPBLAS_STATUS_INVALID_VALUE** – If the parameters are unexpectedly NULL, in conflict or in an impossible configuration.
- **HIPBLAS_STATUS_NOT_INITIALIZED** – If hipBLASLt handle has not been initialized.

```
hipblasStatus_t getDefaultDeviceUserArguments(void *hostDeviceUserArgs)
```

A helper function to initialize DeviceUserArguments using the set problem(s) saved in the gemm object.

Parameters

- **hostDeviceUserArgs** – [in] The DeviceUserArguments structure allocated in host. Note that the user must put the correct type of the DeviceUserArguments.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If the operation completed successfully.

```
hipblasStatus_t run(void *deviceUserArgs, hipStream_t stream)
```

Run the kernel using DeviceUserArguments.

Parameters

- **deviceUserArgs** – [in] Pointer to the DeviceUserArguments buffer allocated in the GPU memory. Pointer must be 16B aligned (that is, lowest 4 bits of
- **stream** – [in] The HIP stream where all the GPU work will be submitted.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If the operation completed successfully.
- **HIPBLAS_STATUS_INVALID_VALUE** – If the `gemm_count = 0`.

`hipblasStatus_t run(hipStream_t stream, hipEvent_t start = nullptr, hipEvent_t stop = nullptr)`

Execute the kernel arguments stored inside the `hipblasLt_ext::GemmInstance`.

Parameters

- **stream** – [in] The HIP stream where all the GPU work will be
- **start** – [in] The HIP event which will record the start of the kernel
- **stop** – [in] The HIP event which will record the end of the kernel submitted.

Return values

HIPBLAS_STATUS_SUCCESS – If the operation completed successfully.

4.3 hipBLASLtExt API Reference

4.3.1 getAllAlgos()

`hipblasStatus_t hipblasLt_ext::getAllAlgos(hipblasLtHandle_t handle, GemmType typeGemm, hipblasOperation_t opA, hipblasOperation_t opB, hipDataType typeA, hipDataType typeB, hipDataType typeC, hipDataType typeD, hipblasComputeType_t typeCompute, std::vector<hipblasLtMatmulHeuristicResult_t> &heuristicResults)`

Retrieve the possible algorithms.

This function retrieves the possible algorithms for the matrix multiply operation `hipblasLtMatmul()` function with the given data and compute type. The output is placed in `heuristicResults` in the order of increasing estimated compute time. It should use `matmulIsAlgoSupported()` to check if the algorithm support the problem before execute `hipblasLtMatmul()`.

Parameters

- **handle** – [in] Pointer to the allocated hipBLASLt handle for the hipBLASLt context. See `hipblasLtHandle_t`.
- **typeGemm** – [in] *Gemm* type. ex. GEMM, GROUPED_GEMM.
- **opA, opB** – [in] Transpose settings of A, B.
- **typeA, typeB, typeC, typeD** – [in] The data type of matrix A, B, C, D.
- **typeCompute** – [in] The compute type.
- **heuristicResults** – [out] The algorithm heuristic vector.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If query was successful. Inspect returned `AlgoCount > 0`.state for the status of the results.
- **HIPBLAS_STATUS_NOT_SUPPORTED** – If no heuristic function available for current configuration.
- **HIPBLAS_STATUS_INVALID_VALUE** – If no solution is found.

4.3.2 getIndexFromAlgo()

int hipblasLt_ext::getIndexFromAlgo(*hipblasLtMatmulAlgo_t* &algo)

Retrieve the algorithm index.

Parameters

algo – [in] The algorithm.

Return values

int – The index of the algorithm, can be used to get heuristic results from *getAlgosFromIndex*. Returns -1 if the index stored in algo < 0. Note that the index may not be valid if the algo struct is not initialized properly.

4.3.3 getAlgosFromIndex()

hipblasStatus_t hipblasLt_ext::getAlgosFromIndex(*hipblasLtHandle_t* handle, std::vector<int> &algoIndex, std::vector<*hipblasLtMatmulHeuristicResult_t*> &heuristicResults)

Retrieve the possible algorithms.

This function retrieves the possible algorithms for the matrix multiply operation *hipblasLtMatmul()* function with the given index. The output is placed in heuristicResult in the order of increasing estimated compute time. A recorded solution index cannot be used across different version of library. It should use *matmulIsAlgoSupported()* to check if the algorithm support the problem before execute *hipblasLtMatmul()*.

Parameters

- **handle** – [in] Pointer to the allocated hipBLASLt handle for the hipBLASLt context. See *hipblasLtHandle_t*.
- **algoIndex** – [in] The algorithm index vector.
- **heuristicResults** – [out] The algorithm heuristic vector.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If query was successful. Inspect heuristicResults.size() > 0.state for the status of the results.
- **HIPBLAS_STATUS_NOT_SUPPORTED** – If no heuristic function available for current configuration.
- **HIPBLAS_STATUS_INVALID_VALUE** – If query indexes are all out of bound of solution map.

4.3.4 matmulIsAlgoSupported()

hipblasStatus_t hipblasLt_ext::matmulIsAlgoSupported(*hipblasLtHandle_t* handle, *hipblasLtMatmulDesc_t* matmulDesc, const void *alpha, *hipblasLtMatrixLayout_t* Adesc, *hipblasLtMatrixLayout_t* Bdesc, const void *beta, *hipblasLtMatrixLayout_t* Cdesc, *hipblasLtMatrixLayout_t* Ddesc, *hipblasLtMatmulAlgo_t* &algo, size_t &workspaceSizeInBytes)

Check if the algorithm supports the problem. (For hipblasLt API)

This function updates the problem saved inside the algorithm if the problem is supported. The required workspaceSizeInBytes is also returned.

Parameters

- **handle** – [in] Pointer to the allocated hipBLASLt handle for the hipBLASLt context. See *hipblasLtHandle_t*.
- **matmulDesc** – [in] Handle to a previously created matrix multiplication descriptor of type *hipblasLtMatmulDesc_t*.
- **alpha, beta** – [in] Pointers to the scalars used in the multiplication.
- **Adesc, Bdesc, Cdesc, Ddesc** – [in] Handles to the previously created matrix layout descriptors of the type *hipblasLtMatrixLayout_t*.
- **algo** – [in] The algorithm heuristic.
- **workspaceSizeInBytes** – [out] Return the required workspace size.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If query was successful. The problem is supported by the algorithm. results.
- **HIPBLAS_STATUS_INVALID_VALUE** – The problem is not supported.

4.4 hipblasLtExt usage

Here are the three use-cases supported by the hipBLASLtExt APIs.

4.4.1 Gemm

hipblasLt has its own instance.

You must assign the problem type when constructing or importing the problem from hipBLAS API.

```
HIPBLASLT_EXPORT explicit Gemm(hipblasLtHandle_t    handle,
                               hipblasOperation_t   opA,
                               hipblasOperation_t   opB,
                               hipDataType          typeA,
                               hipDataType          typeB,
                               hipDataType          typeC,
                               hipDataType          typeD,
                               hipblasComputeType_t typeCompute);

HIPBLASLT_EXPORT explicit Gemm(hipblasLtHandle_t    handle,
                               hipblasLtMatmulDesc_t matmul_descr,
                               const void*         alpha,
                               const void*         A,
                               hipblasLtMatrixLayout_t matA,
                               const void*         B,
                               hipblasLtMatrixLayout_t matB,
                               const void*         beta,
                               const void*         C,
                               hipblasLtMatrixLayout_t matC,
                               void*              D,
                               hipblasLtMatrixLayout_t matD);
```

After the instance is created, you can set the problem using the API. The API may require the following structures:

`GemmProblemType` allows you to change the problem type after the instance is initialized. (Deprecated, use `GemmProblemTypeV2` instead)

```

struct GemmProblemType
{
    hipblasOperation_t    op_a;
    hipblasOperation_t    op_b;
    hipDataType           type_a;
    hipDataType           type_b;
    hipDataType           type_c;
    hipDataType           type_d;
    hipblasComputeType_t  type_compute;
};

```

GemmEpilogue allows the user to control the epilogue of the problem. (Deprecated, use GemmEpilogueV2 instead)

```

struct GemmEpilogue
{
    hipblasLtEpilogue_t mode = HIPBLASLT_EPILOGUE_DEFAULT;
    hipDataType          bias_data_type;
    int                  aux_ld;
    int                  aux_stride;
};

```

GemmInputs specifies the problem inputs. (Deprecated, use GemmInputsV2 instead)

```

struct GemmInputs
{
    void* a = nullptr;
    void* b = nullptr;
    void* c = nullptr;
    void* d = nullptr;
    void* alpha = nullptr;
    void* beta = nullptr;
    // Epilogue inputs
    void* bias = nullptr;
    void* aux = nullptr;
};

```

setProblem APIs:

```

HIPBLASLT_EXPORT hipblasStatus_t setProblem(
    int64_t m, int64_t n, int64_t k, int64_t batch_count, GemmEpilogueV2& epilogue,
    ↪GemmInputsV2& inputs);

```

You can also set the leading dimensions and strides, and reassign the data type with the following API:

```

HIPBLASLT_EXPORT hipblasStatus_t setProblem(int64_t          m,
                                             int64_t          n,
                                             int64_t          k,
                                             int64_t          batch_count,
                                             int64_t          lda,
                                             int64_t          ldb,
                                             int64_t          ldc,
                                             int64_t          ldd,
                                             int64_t          strideA,

```

(continues on next page)

(continued from previous page)

```

int64_t      strideB,
int64_t      strideC,
int64_t      strideD,
GemmEpilogueV2& epilogue,
GemmInputsV2& inputs,
GemmProblemTypeV2& problemtype);

```

You can also import problems from hipblasLt APIs after the instance is created. Note that this may overwrite the problem type of the instance.

```

HIPBLASLT_EXPORT hipblasStatus_t setProblem(hipblasLtMatmulDesc_t matmul_descr,
const void* alpha,
const void* A,
hipblasLtMatrixLayout_t matA,
const void* B,
hipblasLtMatrixLayout_t matB,
const void* beta,
const void* C,
hipblasLtMatrixLayout_t matC,
void* D,
hipblasLtMatrixLayout_t matD);

```

You can get heuristics and make kernel arguments with the instance. If the properties of the gemm and the inputs don't change, you can call the run API to launch the kernel directly.

```

// Pseudo code
hipblaslt_ext::GemmPreferenceV2 pref;
pref.setMaxWorkspaceBytes(10000000);
// Default epilogue mode is HIPBLASLT_EPILOGUE_DEFAULT
hipblaslt_ext::GemmEpilogueV2 epilogue;
hipblaslt_ext::GemmInputsV2 inputs;
inputs.setA(d_a);
inputs.setB(d_b);
inputs.setC(d_c);
inputs.setD(d_d);
inputs.setAlpha(&alpha);
inputs.setBeta(&beta);

hipblaslt_ext::Gemm gemm(handle,
HIPBLAS_OP_N,
HIPBLAS_OP_N,
HIP_R_16F,
HIP_R_16F,
HIP_R_16F,
HIP_R_16F,
HIPBLAS_COMPUTE_32F);
std::vector<hipblasLtMatmulHeuristicResult_t> heuristic;
gemm.setProblem(1, 1, 1, 1, epilogue, inputs); // m, n, k, batch
gemm.algoGetHeuristic(gemm, pref, heuristic);
gemm.initialize(heuristic[0].algo, d_workspace, stream);
for(int i = 0; i < 10; i++)
{

```

(continues on next page)

(continued from previous page)

```
gemm.run(stream);
}
```

4.4.2 Grouped gemm

hipblasLtExt supports grouped gemm. It shares the same class with normal gemm.

After the problem is set, you can check the problem type with function `getGemmType()`.

```
enum class GemmType
{
    HIPBLASLT_GEMM          = 1,
    HIPBLASLT_GROUPED_GEMM = 2
};
```

The grouped gemm class also has the `setProblem` APIs.

```
HIPBLASLT_EXPORT hipblasStatus_t setProblem(
    int64_t m, int64_t n, int64_t k, int64_t batch_count, GemmEpilogueV2& epilogue,
    GemmInputsV2& inputs);

HIPBLASLT_EXPORT hipblasStatus_t setProblem(std::vector<int64_t>& m,
                                             std::vector<int64_t>& n,
                                             std::vector<int64_t>& k,
                                             std::vector<int64_t>& batch_count,
                                             std::vector<GemmEpilogueV2>& epilogue,
                                             std::vector<GemmInputsV2>& inputs);

HIPBLASLT_EXPORT hipblasStatus_t setProblem(std::vector<int64_t>& m,
                                             std::vector<int64_t>& n,
                                             std::vector<int64_t>& k,
                                             std::vector<int64_t>& batch_count,
                                             std::vector<int64_t>& lda,
                                             std::vector<int64_t>& ldb,
                                             std::vector<int64_t>& ldc,
                                             std::vector<int64_t>& ldd,
                                             std::vector<int64_t>& strideA,
                                             std::vector<int64_t>& strideB,
                                             std::vector<int64_t>& strideC,
                                             std::vector<int64_t>& strideD,
                                             std::vector<GemmEpilogueV2>& epilogue,
                                             std::vector<GemmInputsV2>& inputs,
                                             GemmProblemTypeV2& problemtype);

HIPBLASLT_EXPORT hipblasStatus_t setProblem(std::vector<hipblasLtMatmulDesc_t>& matmul_
    descr,
                                             std::vector<void*>& alpha,
                                             std::vector<void*>& A,
                                             std::vector<hipblasLtMatrixLayout_t>& matA,
                                             std::vector<void*>& B,
                                             std::vector<hipblasLtMatrixLayout_t>& matB,
                                             std::vector<void*>& beta,
                                             std::vector<void*>& C,
```

(continues on next page)

(continued from previous page)

```

std::vector<hipblasLtMatrixLayout_t>& matC,
std::vector<void*>& D,
std::vector<hipblasLtMatrixLayout_t>& matD);

```

For the following API, the argument “epilogue” supports broadcasting. They are broadcasted to the length of the problem size by duplicating the last element.

```

HIPBLASLT_EXPORT hipblasStatus_t setProblem(std::vector<int64_t>& m,
std::vector<int64_t>& n,
std::vector<int64_t>& k,
std::vector<int64_t>& batch_count,
std::vector<int64_t>& lda,
std::vector<int64_t>& ldb,
std::vector<int64_t>& ldc,
std::vector<int64_t>& ldd,
std::vector<int64_t>& strideA,
std::vector<int64_t>& strideB,
std::vector<int64_t>& strideC,
std::vector<int64_t>& strideD,
std::vector<GemmEpilogueV2>& epilogue,
std::vector<GemmInputsV2>& inputs,
GemmProblemTypeV2& problemtype);

```

Note that currently only problemtype size equal to 1 (Only one GemmProblemTypeV2 for all problems) is supported.

```

// Pseudo code
std::vector<int64_t> m, n, k;
// ...
for(size_t i = 0; i < problem_size, i++)
{
    // ...
}
std::vector<GemmProblemTypeV2> problemtypes;
problemtypes.push_back(problemtype);
groupedgemm.setProblem(m, n, k, batch_count, lda, ldb, ldc, ldd, strideA, strideB,
↳strideC, strideD, epilogue, inputs, problemtypes);

```

UserArguments

Grouped gemm supports the use of external device memory to run the kernel. This is helpful if some of the arguments are from the output of the pervious kernel. To change the size (m, n, k, batch) related arguments, refer to *Fixed MK*.

```

struct UserArguments
{
    uint32_t m; //!< size m
    uint32_t n; //!< size n
    uint32_t batch; //!< size batch
    uint32_t k; //!< size k
    void* d; //!< The d matrix input pointer.
    void* c; //!< The c matrix input pointer.
    void* a; //!< The a matrix input pointer.
    void* b; //!< The b matrix input pointer.
}

```

(continues on next page)

(continued from previous page)

```

uint32_t strideD1; //!< The d leading dimension.
uint32_t strideD2; //!< The d batch stride
uint32_t strideC1; //!< The c leading dimension.
uint32_t strideC2; //!< The c batch stride
uint32_t strideA1; //!< The a leading dimension.
uint32_t strideA2; //!< The a batch stride
uint32_t strideB1; //!< The b leading dimension.
uint32_t strideB2; //!< The b batch stride
int8_t  alpha[16]; //!< The alpha value.
int8_t  beta[16];  //!< The beta value.
// Epilogue inputs
void*   bias;      //!< The bias input pointer.
int     biasType;  //!< The bias datatype. Only works if mode is set to bias related.
↪epilogues.
uint32_t reserved;
void*   e;         //!< The aux input pointer. Only works if mode is set to aux related.
↪epilogues.
uint32_t strideE1; //!< The aux leading dimension. Only works if mode is set to aux.
↪related epilogues.
uint32_t strideE2; //!< The aux batch stride. Only works if mode is set to aux.
↪related epilogues.
float   act0;      //!< The activation value 1. Some activations might use it.
float   act1;      //!< The activation value 2.
int     activationType; //!< The activation type. Only works if mode is set to.
↪activation related epilogues.
} __attribute__((packed));

```

We add the two functions for UserArguments related API. The first API is a helper function that helps the user to initialize the structure “UserArguments” from the saved problems inside the grouped gemm object. The second API is an overload function with an additional UserArguments device pointer input.

```

HIPBLASLT_EXPORT hipblasStatus_t getDefaultValueForDeviceUserArguments(void* ↪
↪hostDeviceUserArgs);

HIPBLASLT_EXPORT hipblasStatus_t run(void* deviceUserArgs, hipStream_t stream);

```

The following is a simple example of how this API works.

```

// Pseudo code
// Step 1: Get all algorithms
std::vector<hipblasLtMatmulHeuristicResult_t> heuristicResult;
CHECK_HIPBLASLT_ERROR(hipblaslt_ext::getAllAlgos(handle,
                                                    HIPBLASLT_GEMM,
                                                    HIPBLAS_OP_N,
                                                    HIPBLAS_OP_N,
                                                    in_out_datatype,
                                                    in_out_datatype,
                                                    in_out_datatype,
                                                    in_out_datatype,
                                                    HIPBLAS_COMPUTE_32F,
                                                    heuristicResult));

```

(continues on next page)

(continued from previous page)

```

hipblaslt_ext::GemmPreferenceV2 pref;
pref.setMaxWorkspaceBytes(10000000);
// Step 2: Setup problem
std::vector<int64_t> m(gemm_count);
std::vector<int64_t> n(gemm_count);
std::vector<int64_t> k(gemm_count);
std::vector<int64_t> batch_count(gemm_count);
std::vector<hipblaslt_ext::GemmEpilogueV2> epilogue(gemm_count);
std::vector<hipblaslt_ext::GemmInputsV2> inputs(gemm_count);
for(int i = 0; i < gemm_count; i++)
{
    m[i] = 1;
    n[i] = 1;
    k[i] = 1;
    batch_count[i] = 1;
    epilogue[i].setMode(HIPBLASLT_EPILOGUE_GELU);
    inputs[i].setA(d_a[i]);
    inputs[i].setB(d_b[i]);
    inputs[i].setC(d_c[i]);
    inputs[i].setD(d_d[i]);
    inputs[i].setAlpha(&alpha[i]);
    inputs[i].setBeta(&beta[i]);
}

// Step 3: Create grouped gemm instance
hipblaslt_ext::GroupedGemm groupedGemm(handle,
                                       HIPBLAS_OP_N,
                                       HIPBLAS_OP_N,
                                       HIP_R_16F,
                                       HIP_R_16F,
                                       HIP_R_16F,
                                       HIP_R_16F,
                                       HIPBLAS_COMPUTE_32F);

// Step 4: Set problem
groupedGemm.setProblem(m, n, k, batch_count, epilogue, inputs); // m, n, k, batch

// Step 5: Get default value from the instance
hipblaslt_ext::UserArguments* dUAFloat = new hipblaslt_ext::UserArguments[gemm_count];
groupedGemm.getDefaultValueForDeviceUserArguments((void*)dUAFloat);
// Once you get the default value here, you can make several copies and change the values
// from the host

// Next Copy them to the device memory
hipblaslt_ext::UserArguments* d_dUAFloat = nullptr;
hipMalloc(&d_dUAFloat, sizeof(hipblaslt_ext::UserArguments) * gemm_count);
hipMemcpy(d_dUAFloat, dUAFloat, sizeof(hipblaslt_ext::UserArguments) * gemm_count,
↳hipMemcpyHostToDevice);

validIdx.clear();
for(int j = 0; j < heuristicResult.size(); j++)
{

```

(continues on next page)

(continued from previous page)

```

size_t workspace_size = 0;
if(groupedGemm.isAlgoSupported(heuristicResult[j].algo, workspace_size)
    == HIPBLAS_STATUS_SUCCESS)
{
    validIdx.push_back(j);
}
}

// Step 6: Initialize and run
if(validIdx.size() > 1)
{
    groupedGemm.initialize(heuristicResult[validIdx[0]].algo, d_workspace, stream);
    for(int i = 0; i < 10; i++)
    {
        groupedGemm.run(userArgs, stream);
    }
}

```

4.4.3 The base class (GemmInstance)

This is the base class for Gemm and GroupedGemm.

```

// Gets heuristic from the instance.
HIPBLASLT_EXPORT hipblasStatus_t algoGetHeuristic(const int
↳ requestedAlgoCount,
↳
↳ pref,
↳
↳ std::vector
↳ <hipblasLtMatmulHeuristicResult_t>& heuristicResults);

// Returns SUCCESS if the algo is supported, also returns the required workspace size in
↳ bytes.
HIPBLASLT_EXPORT hipblasStatus_t isAlgoSupported(hipblasLtMatmulAlgo_t& algo, size_t&
↳ workspaceSizeInBytes);

// Initializes the instance before calling run. Requires every time the problem is set.
HIPBLASLT_EXPORT hipblasStatus_t initialize(const hipblasLtMatmulAlgo_t& algo, void*
↳ workspace, bool useUserArgs = true, hipStream_t stream = 0);

// Run the problem.
HIPBLASLT_EXPORT hipblasStatus_t run(hipStream_t stream);

```

4.4.4 Get all algorithms

Get all algorithms allows you to get all the algorithms for a specific problem type. It requires the transpose of A, B, data type of the inputs, and the compute type.

```

HIPBLASLT_EXPORT
hipblasStatus_t hipblaslt_ext::getAllAlgos(hipblasLtHandle_t
↳ handle,
↳
↳ hipblasLtExtGemmTypeEnum_t
↳ typeGemm,

```

(continues on next page)

(continued from previous page)

```

↪ opA, hipblasOperation_t
↪ opB, hipblasOperation_t
↪ typeA, hipDataType
↪ typeB, hipDataType
↪ typeC, hipDataType
↪ typeD, hipDataType
↪ typeCompute, hipblasComputeType_t
↪ heuristicResults);
std::vector<hipblasLtMatmulHeuristicResult_t>&

```

This API doesn't require any problem size or epilogue as input, but uses another API `isAlgoSupported` to check if the algorithm supports a problem.

```

hipblaslt_ext::matmulIsAlgoSupported()
gemm.isAlgoSupported()

```

The API returns the required workspace size in bytes on successful completion.

```

// Get all algorithms
CHECK_HIPBLASLT_ERROR(hipblaslt_ext::getAllAlgos(handle,
HIPBLASLT_GEMM,
HIPBLAS_OP_N,
HIPBLAS_OP_N,
in_out_datatype,
in_out_datatype,
in_out_datatype,
in_out_datatype,
HIPBLAS_COMPUTE_32F,
heuristicResult));

validIdx.clear();
for(int j = 0; j < heuristicResult.size(); j++)
{
    size_t workspace_size = 0;
    if(hipblaslt_ext::matmulIsAlgoSupported(handle,
matmul,
&(alpha),
matA,
matB,
&(beta),
matC,
matD,
heuristicResult[j].algo,
workspace_size)
== HIPBLAS_STATUS_SUCCESS)

```

(continues on next page)

(continued from previous page)

```

{
    validIdx.push_back(j);
    heuristicResult[j].workspaceSize = workspace_size;
}
else
{
    heuristicResult[j].workspaceSize = 0;
}
}

```

4.4.5 Algorithm index

This extension API allows you to get the algorithm index using `hipblasLtMatmulAlgo_t`.

```
HIPBLASLT_EXPORT int getIndexFromAlgo(hipblasLtMatmulAlgo_t& algo);
```

It also supports you to get the heuristic results by giving an index vector.

```

HIPBLASLT_EXPORT
hipblasStatus_t
    getAlgosFromIndex(hipblasLtHandle_t                handle,
                     std::vector<int>&                 algoIndex,
                     std::vector<hipblasLtMatmulHeuristicResult_t>& heuristicResults);

```

4.5 Sample codes

Here are the sample codes demonstrating use cases of the extension APIs.

4.5.1 Gemm

```

// Pseudo code for gemm problem
// Get all algorithms
std::vector<hipblasLtMatmulHeuristicResult_t> heuristicResult;
CHECK_HIPBLASLT_ERROR(hipblaslt_ext::getAllAlgos(handle,
                                                  HIPBLASLT_GEMM,
                                                  HIPBLAS_OP_N,
                                                  HIPBLAS_OP_N,
                                                  in_out_datatype,
                                                  in_out_datatype,
                                                  in_out_datatype,
                                                  in_out_datatype,
                                                  HIPBLAS_COMPUTE_32F,
                                                  heuristicResult));

hipblaslt_ext::GemmPreferenceV2 pref;
pref.setMaxWorkspaceBytes(10000000);
hipblaslt_ext::GemmEpilogueV2 epilogue;
epilogue.setMode(HIPBLASLT_EPILOGUE_GELU);
hipblaslt_ext::GemmInputsV2 inputs;
inputs.setA(d_a);
inputs.setB(d_b);

```

(continues on next page)

(continued from previous page)

```

inputs.setC(d_c);
inputs.setD(d_d);
inputs.setAlpha(&alpha);
inputs.setBeta(&beta);

hipblaslt_ext::Gemm gemm(handle,
                        HIPBLAS_OP_N,
                        HIPBLAS_OP_N,
                        HIP_R_16F,
                        HIP_R_16F,
                        HIP_R_16F,
                        HIP_R_16F,
                        HIPBLAS_COMPUTE_32F);

gemm.setProblem(1, 1, 1, 1, epilogue, inputs); // m, n, k, batch

validIdx.clear();
for(int j = 0; j < heuristicResult.size(); j++)
{
    size_t workspace_size = 0;
    if(gemm.isAlgoSupported(heuristicResult[j].algo, workspace_size)
        == HIPBLAS_STATUS_SUCCESS)
    {
        validIdx.push_back(j);
        heuristicResult[j].workspaceSize = workspace_size;
    }
    else
    {
        heuristicResult[j].workspaceSize = 0;
    }
}

if(validIdx.size() > 1)
{
    gemm.initialize(heuristicResult[validIdx[0]].algo, d_workspace, stream);
    for(int i = 0; i < 10; i++)
    {
        gemm.run(stream);
    }
}

```

4.5.2 Grouped gemm

```

// Pseudo code for grouped gemm problem
// Get all algorithms
std::vector<hipblasLtMatmulHeuristicResult_t> heuristicResult;
CHECK_HIPBLASLT_ERROR(hipblaslt_ext::getAllAlgos(handle,
                                                HIPBLASLT_GEMM,
                                                HIPBLAS_OP_N,
                                                HIPBLAS_OP_N,
                                                in_out_datatype,

```

(continues on next page)

(continued from previous page)

```

        in_out_datatype,
        in_out_datatype,
        in_out_datatype,
        HIPBLAS_COMPUTE_32F,
        heuristicResult));

hipblaslt_ext::GemmPreferenceV2 pref;
pref.setMaxWorkspaceBytes(10000000);

std::vector<int64_t> m(gemm_count);
std::vector<int64_t> n(gemm_count);
std::vector<int64_t> k(gemm_count);
std::vector<int64_t> batch_count(gemm_count);
std::vector<hipblaslt_ext::GemmEpilogueV2> epilogue(gemm_count);
std::vector<hipblaslt_ext::GemmInputsV2> inputs(gemm_count);
for(int i = 0; i < gemm_count; i++)
{
    m[i] = 1;
    n[i] = 1;
    k[i] = 1;
    batch_count[i] = 1;
    epilogue[i].setMode(HIPBLASLT_EPILOGUE_GELU);
    inputs[i].setA(d_a[i]);
    inputs[i].setB(d_b[i]);
    inputs[i].setC(d_c[i]);
    inputs[i].setD(d_d[i]);
    inputs[i].setAlpha(&alpha[i]);
    inputs[i].setBeta(&beta[i]);
}

hipblaslt_ext::GroupedGemm groupedGemm(handle,
        HIPBLAS_OP_N,
        HIPBLAS_OP_N,
        HIP_R_16F,
        HIP_R_16F,
        HIP_R_16F,
        HIP_R_16F,
        HIPBLAS_COMPUTE_32F);

groupedGemm.setProblem(m, n, k, batch_count, epilogue, inputs); // m, n, k, batch

validIdx.clear();
for(int j = 0; j < heuristicResult.size(); j++)
{
    size_t workspace_size = 0;
    if(groupedGemm.isAlgoSupported(heuristicResult[j].algo, workspace_size)
        == HIPBLAS_STATUS_SUCCESS)
    {
        validIdx.push_back(j);
    }
}

```

(continues on next page)

(continued from previous page)

```

if(validIdx.size() > 1)
{
    groupedGemm.initialize(heuristicResult[validIdx[0]].algo, d_workspace, stream);
    for(int i = 0; i < 10; i++)
    {
        groupedGemm.run(stream);
    }
}

```

4.5.3 Algorithm index

```

int index = hipblaslt_ext::getIndexFromAlgo(testResults[i].algo);
// Save the index to disk or somewhere else for later use.

// Get the index from previous state.
std::vector<int> algoIndex{index};
std::vector<hipblasLtMatmulHeuristicResult_t> heuristicResults;
// If the index is out of the bound of solutions, getAlgosFromIndex will return HIPBLAS_
↪STATUS_INVALID_VALUE
if(HIPBLAS_STATUS_INVALID_VALUE
    == hipblaslt_ext::getAlgosFromIndex(handle, algoIndex, heuristicResults))
{
    std::cout << "Indexes are all out of bound." << std::endl;
    break;
}

```

4.5.4 [Grouped Gemm] Fixed MK

hipBLASLt extension supports changing the sizes (m, n, k, batch) from the device memory UserArguments, but the setup is a bit different from the normal routing.

Sum of n:

A sum of N needs to be used as an input for the grouped gemm instance.

```

{1000, 1, 1, 1}; // The array of N, the first element is the sum of N

// Below is the values stored in "UserArguments"
{256, 256, 1, 1}; // This is a valid configuration cause 256 + 256 + 1 + 1 < 1000
{512, 512, 1, 1}; // This is NOT a valid configuration cause 512 + 512 + 1 + 1 > 1000

```

For example, we have a grouped gemm with gemm_count = 4. The sum of N must not exceed the “sum of N” set in setProblem API. In this mode, the first element is the “sum of n” in the array of Ns.

```

// Pseudo code
// Step 1: Get all algorithms
std::vector<hipblasLtMatmulHeuristicResult_t> heuristicResult;
CHECK_HIPBLASLT_ERROR(hipblaslt_ext::getAllAlgos(handle,
                                                    HIPBLASLT_GEMM,
                                                    HIPBLAS_OP_N,
                                                    HIPBLAS_OP_N,

```

(continues on next page)

(continued from previous page)

```

        in_out_datatype,
        in_out_datatype,
        in_out_datatype,
        in_out_datatype,
        HIPBLAS_COMPUTE_32F,
        heuristicResult));

hipblaslt_ext::GemmPreferenceV2 pref;
pref.setMaxWorkspaceBytes(10000000);
// Step 2: Setup problem
std::vector<int64_t> m(gemm_count);
std::vector<int64_t> n(gemm_count);
std::vector<int64_t> k(gemm_count);
std::vector<int64_t> batch_count(gemm_count);
std::vector<hipblaslt_ext::GemmEpilogueV2> epilogue(gemm_count);
std::vector<hipblaslt_ext::GemmInputsV2> inputs(gemm_count);

// Step 2.1: Calculate sum of n
int64_t sum_of_n = 0;
for(int i = 0; i < gemm_count; i++)
{
    sum_of_n += n_arr[i];
}

// {sum_of_n, 1, 1, 1, ...}; // The array of N, the first element is the sum of N
for(int i = 0; i < gemm_count; i++)
{
    m[i] = m_arr[i];
    if(i == 0)
        n[i] = sum_of_n;
    else
        n[i] = 1;
    k[i] = k_arr[i];
    batch_count[i] = 1;
    inputs[i].setA(d_a[i]);
    inputs[i].setB(d_b[i]);
    inputs[i].setC(d_c[i]);
    inputs[i].setD(d_d[i]);
    inputs[i].setAlpha(&alpha[i]);
    inputs[i].setBeta(&beta[i]);
}

// Step 3: Create grouped gemm instance
hipblaslt_ext::GroupedGemm groupedGemm(handle,
        HIPBLAS_OP_N,
        HIPBLAS_OP_N,
        HIP_R_16F,
        HIP_R_16F,
        HIP_R_16F,
        HIP_R_16F,
        HIPBLAS_COMPUTE_32F);

```

(continues on next page)

(continued from previous page)

```

// Step 4: Set problem
groupedGemm.setProblem(m, n, k, batch_count, epilogue, inputs); // m, n, k, batch

// Step 5: Get default value from the instance
hipblaslt_ext::UserArguments* dUAFloat = new hipblaslt_ext::UserArguments[gemm_count];
groupedGemm.getDefaultValueForDeviceUserArguments((void*)dUAFloat);
// Once you get the default value here, you can make several copies and change the values
// from the host

// Next Copy them to the device memory
hipblaslt_ext::UserArguments* d_dUAFloat = nullptr;
hipMalloc(&d_dUAFloat, sizeof(hipblaslt_ext::UserArguments) * gemm_count);
hipMemcpy(d_dUAFloat, dUAFloat, sizeof(hipblaslt_ext::UserArguments) * gemm_count,
↳hipMemcpyHostToDevice);

validIdx.clear();
for(int j = 0; j < heuristicResult.size(); j++)
{
    size_t workspace_size = 0;
    if(groupedGemm.isAlgoSupported(heuristicResult[j].algo, workspace_size)
        == HIPBLAS_STATUS_SUCCESS)
    {
        validIdx.push_back(j);
    }
}

int threads = 256;
int blocks = ceil((double)gemm_count / threads);

// Step 6: Initialize and run
if(validIdx.size() > 1)
{
    groupedGemm.initialize(heuristicResult[validIdx[0]].algo, d_workspace);
    for(int i = 0; i < 10; i++)
    {
        hipLaunchKernelGGL(kernelUpdateN,
                            dim3(blocks),
                            dim3(threads),
                            0,
                            stream,
                            gemm_count,
                            d_dUAFloat,
                            d_n_vec); // d_n_vec is a device pointer with Ns
        groupedGemm.run(userArgs, stream);
    }
}

// .....

__global__ void kernelUpdateN(uint32_t gemm_count, void* userArgs, int32_t* sizes_n)
{
    uint64_t id = hipBlockIdx_x * 256 + hipThreadIdx_x;

```

(continues on next page)

(continued from previous page)

```
if(id >= gemm_count)
    return;

hipblaslt_ext::UserArguments* dUAFloat = static_cast<hipblaslt_ext::UserArguments*>
↳(userArgs);
dUAFloat[id].n                = sizes_n[id];
}
```


HIPBLASLTEXT OPERATION API REFERENCE

hipBLASLt has the following extension operation APIs that are independent to gemm operations. These extensions support:

1. **hipblasltExtSoftmax**
Softmax for 2D-tensor. Currently, it performs softmax on the second dimension of input tensor and assumes the input to be contiguous on the second dimension. For sample code, refer to `client_extop_softmax`.
2. **hipblasltExtLayerNorm**
Converts a 2D tensor using LayerNorm to generate a new 2D normalized tensor. it is an independent function used to just call and get result. For sample code, refer to `sample_hipblaslt_ext_op_layernorm`.
3. **hipblasltExtAMax**
Abs maximum value of a 2D tensor. it is an independent function used to just call and get result. For sample code, refer to `sample_hipblaslt_ext_op_amax`.
4. **hipblasltExtAMaxWithScale**
Abs maximum value and scaled output of a 2D tensor. it is an independent function used to just call and get result. For sample code, refer to `sample_hipblaslt_ext_op_amax_with_scale`.

These APIs are explained in detail below.

5.1 hipblasltExtSoftmax()

hipblasStatus_t **hipblasltExtSoftmax**(hipDataType datatype, uint32_t m, uint32_t n, uint32_t dim, void *output, void *input, hipStream_t stream)

Perform softmax on given tensor.

This function computes softmax on given 2D-tensor along specified dimension.

Parameters

- **datatype** – [in] Datatype of input/output tensor, currently support HIP_R_32F only.
- **m** – [in] The first dimension of input/output tensor.
- **n** – [in] The second dimension of input/output tensor. Currently only values less than or equal to 256 are supported.
- **dim** – [in] Specified dimension to perform softmax on. Currently 1 is the only valid value.
- **input** – [in] input tensor buffer.
- **stream** – [in] The HIP stream where all the GPU work will be submitted.
- **output** – [out] output tensor buffer.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If it runs successfully.
- **HIPBLAS_STATUS_INVALID_VALUE** – If `n` is greater than 256.
- **HIPBLAS_STATUS_NOT_SUPPORTED** – If `dim` is not 1 or `datatype` is not `HIP_R_32F`.

5.2 hipblasLtExtLayerNorm()

`hipblasStatus_t hipblasLtExtLayerNorm`(`hipDataType` datatype, void *output, void *mean, void *invvar, void *input, `uint32_t` m, `uint32_t` n, float eps, void *gamma, void *beta, `hipStream_t` stream)

Perform 2-D layernorm on with source input tensor and result output tensor.

This function computes layernorm on given 2D-tensor.

Parameters

- **datatype** – [in] Datatype of input/output tensor, currently support `HIP_R_32F` only.
- **output** – [out] output tensor buffer. can't be nullptr.
- **mean** – [out] tensor buffer. can't be nullptr.
- **invvar** – [out] tensor buffer. $1 / \text{sqrt}(\text{std})$. can't be nullptr.
- **input** – [in] tensor buffer. can't be nullptr.
- **m** – [in] The first dimension of input/output tensor.
- **n** – [in] The second dimension of input/output tensor.
- **eps** – [in] for sqrt to avoid inf value.
- **gamma** – [in] tensor buffer. nullptr means calculation doesn't involve gamma.
- **beta** – [in] tensor buffer. nullptr means calculation doesn't involve beta.
- **stream** – [in] The HIP stream where all the GPU work will be submitted.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If it runs successfully.
- **HIPBLAS_STATUS_INVALID_VALUE** – If `m` is greater than 4096.
- **HIPBLAS_STATUS_NOT_SUPPORTED** – if `datatype` is not `HIP_R_32F`.

5.3 hipblasLtExtAMax()

`hipblasStatus_t hipblasLtExtAMax`(const `hipDataType` datatype, const `hipDataType` outDatatype, void *output, void *input, `uint32_t` m, `uint32_t` n, `hipStream_t` stream)

Perform absmax on given 2-D tensor and output one value absmax(tensor) value.

This function computes amax on given 2D-tensor.

Parameters

- **datatype** – [in] Datatype of input tensor, currently support `HIP_R_32F` and `HIP_R_16F` only.
- **outDatatype** – [in] Datatype of output tensor, currently support `HIP_R_32F` and `HIP_R_16F` only.
- **output** – [out] Amax tensor buffer. can't be nullptr.

- **input** – [in] 2-D tensor buffer. can't be nullptr.
- **m** – [in] The first dimension of input/output tensor.
- **n** – [in] The second dimension of input/output tensor.
- **stream** – [in] The HIP stream where all the GPU work will be submitted.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If it runs successfully.
- **HIPBLAS_STATUS_INVALID_VALUE** – If m or n is 0, or input or output is nullptr.
- **HIPBLAS_STATUS_NOT_SUPPORTED** – If datatype is not (HIP_R_32F or HIP_R_16F).

5.4 hipblasLtExtAMaxWithScale()

hipblasStatus_t **hipblasLtExtAMaxWithScale**(const hipDataType datatype, const hipDataType outDatatype, const hipDataType scaleDatatype, void *output, void *outputD, void *input, void *inputScale, uint32_t m, uint32_t n, hipStream_t stream)

Perform absmax and scaling on given 2-D tensor. Generate one absmax value and scaled 2-D tensor output.

This function computes amax and scaling on given 2D-tensor.

Parameters

- **datatype** – [in] Datatype of input tensor, currently support HIP_R_32F only.
- **outDatatype** – [in] Datatype of output tensor, currently support HIP_R_32F and HIP_R_16F only.
- **scaleDatatype** – [in] Datatype of outputD tensor, currently support HIP_R_8F_E4M3_FNUZ and HIP_R_8F_E5M2_FNUZ only.
- **output** – [out] Amax tensor buffer. can't be nullptr.
- **outputD** – [out] scaled 2-D tensor buffer. can't be nullptr.
- **input** – [in] 2-D tensor buffer. can't be nullptr.
- **inputScale** – [in] 1-D tensor buffer. can't be nullptr. only support float.
- **m** – [in] The first dimension of input/output tensor.
- **n** – [in] The second dimension of input/output tensor.
- **stream** – [in] The HIP stream where all the GPU work will be submitted.

Return values

- **HIPBLAS_STATUS_SUCCESS** – If it runs successfully.
- **HIPBLAS_STATUS_INVALID_VALUE** – If m or n is 0, or input, inputScale, output, or outputD is nullptr.
- **HIPBLAS_STATUS_NOT_SUPPORTED** – If datatype is not HIP_R_32F, or scaleDatatype is not HIP_R_8F_E4M3_FNUZ or HIP_R_8F_E5M2_FNUZ.

LOGGING AND HEURISTICS

This document contains information for debugging and improving the application performance when using hipBLASLt APIs.

6.1 Logging

You can enable the hipBLASLt logging mechanism by setting the following environment variables before launching the target application:

HIPBLASLT_LOG_LEVEL=<level> - where the level can be:

“0” - Off - Logging is disabled (default)
“1” - Error - Only errors are logged
“2” - Trace - API calls that launch HIP kernels log their parameters and important information
“3” - Hints - Hints that can potentially improve the application’s performance
“4” - Info - Provides general information about the library execution, may contain details about heuristic status
“5” - API Trace - API calls log their parameters and important information

HIPBLASLT_LOG_MASK=<mask> - where mask is a combination of the following:

“0” - Off
“1” - Error
“2” - Trace
“4” - Hints
“8” - Info
“16” - API Trace
“32” - Bench
“64” - Profile

HIPBLASLT_LOG_FILE=<file_name> - where file_name is a path to a logging file. The file name can contain %i, which is replaced with the process ID, for example, <file_name>_%i.log. If HIPBLASLT_LOG_FILE is not defined, the log messages are printed to stdout.

HIPBLASLT_ENABLE_MARKER=1

Setting HIPBLASLT_ENABLE_MARKER to 1 will enable marker trace for rocprof profiling.

6.2 Heuristics cache

hipBLASLt uses heuristics to pick the most suitable matmul kernel for execution based on the problem sizes, GPU configuration, and other parameters. This requires performing some computations on the host CPU, which could take tens of microseconds. To overcome this overhead, it is recommended to query the heuristics once using *hipblasLtMatmulAlgoGetHeuristic()* and then reuse the result for subsequent computations using *hipblasLtMatmul()*.

USING THE HIPBLASLT TUNING UTILITY

hipBLASLt has a simple utility that uses the existing kernel pools and searches for the best solution for a given problem size.

7.1 The template.yaml file

The `template.yaml` file can be found in the `utilities` directory.

```
# Two steps, can comment out Bench or CreateLogic if you want to disable.
Bench:
ProblemType: # Same as the given problem type
  ComputeDataType: s
  ComputeInputDataType: s # Usually the same as DataTypeA and DataTypeB unless you_
↪are using mix precisions.
  DataTypeA: s
  DataTypeB: s
  DataTypeC: s
  DataTypeD: s
  TransposeA: 0
  TransposeB: 0
  UseBias: False
TestConfig:
  ColdIter: 20
  Iter: 100 # You can change this to a larger value for a more stable result, but the_
↪executing time also increases.
  AlgoMethod: "all" # Fixed value
  RotatingBuffer: 512 # It's recommended to set this value larger than the cache size_
↪of the GPU.
TuningParameters:
  # SplitK list control parameter example
  # SplitK: [0, 4, 8] # [0] For disable
ProblemSizes:
- [128, 128, 128] # M, N, K
CreateLogic: {} # Fixed
```

7.2 Running the tuning utility

The `find_exact.py` file is found in the `utilities` directory.

Follow these steps to run the tuning:

1. Run `./install.sh` first. See *Installation* for more details.
2. Ensure the `MatchTable.yaml` file exists in `build/release/library`.
3. Run the command `python3 find_exact.py <your yaml file> <hipblaslt_root_folder>/build/release <output folder>`

You should see the following message appear. The following example is for NN FP32 tuning:

```
Running benchmarks
--Running size: result_NN_SSS_128x128x128.txt
```

After the tuning completes, the following message will appear:

```
Creating exact logic
--Reading matching table: <hipblaslt_root_folder>/build/release/library/MatchTable.yaml
--Reading bench files
--Found file <output folder>/0_Bench/result_NN_SSS_88x12x664.txt
Writing logic yaml files: 100%| 1/1 [00:05<00:00, 5.69s/it]
```

The final structure of the output folder is as follows:

```

  ∨ hiptune
    ∨ 0_Bench
      ≡ result_NN_SSS_88x12x664.txt
    ∨ 1_LogicYaml
      ! aquavanjaram_Cijk_Ailk_Blijk_SB_Bias_HAS_SAV_UserArgs.yaml
```

The `0_Bench` folder stores the raw benchmark results. The `1_LogicYaml` folder stores the output, which is a tuned Equality logic yaml file.

USER OFFLINE TUNING

`hipblaslt-bench` can be used to find the best-performing GEMM kernel for a given set of GEMM problems. Use the command line interface to access this functionality. (See *Clients* for more details.)

`hipblaslt-bench` provides the best solution index for a given problem size. This index can be used directly in future GEMM calls through the User Offline Tuning mechanism. However, these indices cannot be reused across library releases or across different device architectures.

8.1 Using `hipblaslt-bench` to run the tuning with the best GEMM kernel

To find and use the best GEMM kernel for a problem, follow these steps:

1. Generate the tuning command line by setting the environment variable `HIPBLASLT_LOG_MASK=32` before calling any `hipBLASLt` APIs. For more details on how to use `hipblaslt-bench`, see *Logging and heuristics*.

In the Bash shell, set the following environment variable:

```
export HIPBLASLT_LOG_MASK=32
```

2. The following command uses `sample_hipblaslt_gemm.cpp` as an example:

```
./sample_hipblaslt_gemm
```

The tuning command displays the following log entry:

```
hipblaslt-bench --api_method c -m 1024 -n 512 -k 1024 --lda 1024 --ldb 1024 --ldc_
↪1024 --ldd 1024 --stride_a 0 --stride_b 0 --stride_c 0 --stride_d 0 --alpha 1.
↪0000000 --beta 1.0000000 --transA N --transB N --batch_count 1 --a_type f16_r --b_
↪type f16_r --c_type f16_r --d_type f16_r --scale_type f32_r --bias_type f32_r --
↪compute_type f32_r --algo_method index --solution_index 56073
```

2. Set the environment variable `HIPBLASLT_TUNING_FILE=<file_name>` to tune and store the tuning result of the best solution indices for the GEMM problems. The `<file_name>` points to the tuning file.

In the Bash shell, set the following environment variable:

```
export HIPBLASLT_TUNING_FILE=tuning.txt
```

The default settings for the following parameters in `hipblaslt-bench` will be changed in the tuning environment.

```
--iters |-i <value>           (Default value is: 1000)
--cold_iters |-j <value>      (Default value is: 1000)
--algo_method <value>         (Default value is: all)
```

(continues on next page)

(continued from previous page)

```
--requested_solution <value>    (Default value is: -1)
--rotating <value>              (Default value is: 512)
```

After the tuning completes, the expected output is displayed as follows:

```
./hipblaslt-bench --api_method c -m 1024 -n 512 -k 1024 --lda 1024 --ldb 1024 --ldc_
↪1024 --ldd 1024 --stride_a 0 --stride_b 0 --stride_c 0 --stride_d 0 --alpha 1.
↪000000 --beta 1.000000 --transA N --transB N --batch_count 1 --a_type f16_r --b_
↪type f16_r --c_type f16_r --d_type f16_r --scale_type f32_r --bias_type f32_r --
↪compute_type f32_r --algo_method index --solution_index 56073
```

Winner:

```
transA,transB,grouped_gemm,batch_count,m,n,k,alpha,lda,stride_a,beta,ldb,stride_b,
↪ldc,stride_c,ldd,stride_d,a_type,b_type,c_type,d_type,compute_type,scaleA,scaleB,
↪scaleC,scaleD,amaxD,activation_type,bias_vector,bias_type,rotating_buffer,
↪hipblaslt-Gflops,hipblaslt-GB/s,us,solution_index
N,N,0,1,1024,512,1024,1,1024,1048576,1,1024,524288,1024,524288,1024,524288,f16_r,
↪f16_r,f16_r,f16_r,f32_r,0,0,0,0,0,none,0,f32_r,512,66613.8,363.509,16.1189,56537
```

3. Set the environment variable `HIPBLASLT_TUNING_OVERRIDE_FILE=<file_name>` to load the tuning file and override the default kernel selection with the optimal kernel choices, where `<file_name>` points to the tuning file.

In the Bash shell, set the following environment variable:

```
export HIPBLASLT_TUNING_OVERRIDE_FILE=tuning.txt
```

For example, a user can use `hipblaslt-bench` (with `algo_method` set to `heuristic`) to obtain the solutions for a problem, which include the best tuning solution index.

```
./hipblaslt-bench --api_method c -m 1024 -n 512 -k 1024 --lda 1024 --ldb 1024 --ldc_
↪1024 --ldd 1024 --stride_a 0 --stride_b 0 --stride_c 0 --stride_d 0 --alpha 1.
↪000000 --beta 1.000000 --transA N --transB N --batch_count 1 --a_type f16_r --b_
↪type f16_r --c_type f16_r --d_type f16_r --scale_type f32_r --bias_type f32_r --
↪compute_type f32_r --algo_method heuristic --requested_solution 1 --print_kernel_
↪info
```

```
transA,transB,grouped_gemm,batch_count,m,n,k,alpha,lda,stride_a,beta,ldb,stride_b,
↪ldc,stride_c,ldd,stride_d,a_type,b_type,c_type,d_type,compute_type,scaleA,scaleB,
↪scaleC,scaleD,amaxD,activation_type,bias_vector,bias_type,rotating_buffer,
↪hipblaslt-Gflops,hipblaslt-GB/s,us,solution_index
[0]:
N,N,0,1,1024,512,1024,1,1024,1048576,1,1024,524288,1024,524288,1024,524288,f16_r,
↪f16_r,f16_r,f16_r,f32_r,0,0,0,0,0,none,0,f32_r,512,37575.2,205.047,28.5758,56537
```

USING STREAM-K WITH HIPBLASLT

hipBLASLt supports the Stream-K library, which reduces library sizes for a wide range of General Matrix-Matrix Multiplication (GEMM) shapes and sizes. It also provides more consistent performance, which might be better in some cases. Stream-K partitions an equal share of the aggregate inner-loop iterations among physical processing elements, which provides a near-perfect utilization of computing resources. For more information about Stream-K, see [Stream-K: Work-centric Parallel Decomposition for Dense Matrix-Matrix Multiplication on the GPU](#).

9.1 Configuring the kernel selection strategy

The `TENSILE_SOLUTION_SELECTION_METHOD` environment variable controls the hipBLASLt kernel selection strategy for GEMM operations. Set this variable to 2 to enable the Stream-K library or leave it set to 0 to continue to use the default settings.

- `TENSILE_SOLUTION_SELECTION_METHOD=0` (Default)
 - Kernels are selected from the standard tuned libraries.
 - The heuristic best kernel is selected from the standard tuning grid.
 - User-driven tuning (tunable ops) only accesses kernels from the standard grid and free-size libraries.
 - This option does NOT use any Stream-K kernels.
- `TENSILE_SOLUTION_SELECTION_METHOD=2` (Stream-K)
 - This enables the optional Stream-K library to use a GEMM scheduling algorithm that results in consistently good peak GEMM performance with far fewer tuned kernels.
 - The heuristic best kernel is selected from the Stream-K library.
 - User-driven tuning (tunable ops) considers all kernels from the standard grid, the free-size library, and the Stream-K library.

9.2 Configuring the kernel selection strategy

You can control The Stream-K kernel launch behavior using the environment variables listed in the following table. By default, Stream-K uses a model to predict the optimal grid size to use when launching a GEMM kernel at runtime. However, you can choose how many workgroups to launch a GEMM kernel with using the Stream-K settings below:

Environment Variable	Description
TENSILE_STREAMK_DYN	Set this variable to 3 to use the default setting. With this setting, the program automatically picks the number of work groups to launch for optimal performance. Set this variable to 0 to disable the dynamic grid and always use all available compute units.
TENSILE_STREAMK_FIXED_GRID	This variable overrides the default grid size and launches Stream-K GEMM kernels using the specified number of workgroups.
TENSILE_STREAMK_MAX_CUS	This variable sets the maximum number of compute units to use for Stream-K kernels. By default, Stream-K kernels are allowed to use all compute units on the device, but this setting lets you limit the number of units that can be used.

9.3 Recommendations for using Stream-K

Stream-K is especially advantageous in certain situations. Follow these guidelines when choosing a kernel selection strategy, based on your application and the desired performance.

- **Wide range of GEMM sizes:** Stream-K is a better choice for applications that handle a variety of GEMM shapes and sizes.
- **Non-uniform dimensions:** Stream-K is particularly beneficial for GEMMs where one dimension is significantly larger than the others.
- **Consistent performance:** Stream-K provides more consistent peak performance than the default selection method by evenly distributing work across the available compute units.

9.3.1 Managing Stream-K resource use

Follow these guidelines to optimize how Stream-K uses resources:

- **Promoting concurrency:** Use TENSILE_STREAMK_FIXED_GRID to limit the number of workgroups. This prevents GEMM from monopolizing the GPU resources and allows other kernels to run concurrently.

The following example limits the GEMM kernels to 64 workgroups:

```
export TENSILE_STREAMK_FIXED_GRID=64
```

- **Limiting compute units:** Use TENSILE_STREAMK_MAX_CUS to restrict the number of compute units the Stream-K kernels can use.

This example limits the GEMM kernels to 32 compute units:

```
export TENSILE_STREAMK_MAX_CUS=32
```

LIBRARY SOURCE CODE ORGANIZATION

The hipBLASLt source code is available in the following two directories:

- The `library` directory contains all source code for the library.
- The `clients` directory contains all test code and code to build clients.

10.1 library directory

Here are the sub-directories in the `library` directory:

- `library/include`

Contains C98 include files for the external API. These files also contain Doxygen comments that document the API.

- `library/src/amd_detail`

Contains implementation of hipBLASLt interface that is compatible with rocBLASLt APIs.

- `library/src/include`

Contains internal include files for converting C++ exceptions to hipBLAS status.

10.2 clients directory

Here are the sub-directories in `clients` directory:

- `clients/samples`

Contains sample code for calling hipBLASLt functions

10.3 Infrastructure

- CMake is used to build and package hipBLASLt. There are `CMakeLists.txt` files throughout the code.
- Doxygen/Breathe/Sphinx/ReadTheDocs are used to produce documentation. The documentation is sourced from:
 - Doxygen comments in `include` files in the `library/include` directory
 - Files in the `docs/source` directory
- Jenkins is used to automate Continuous Integration testing.
- `clang-format` is used to format C++ code.

The following client executables are available for use with hipBLASLt:

- `hipblaslt-test`
- `hipblaslt-bench`

To build these clients, follow the instructions on the [Build and Install hipBLASLt github page](#) . After building the hipBLASLt clients, you can find them in the directory `hipBLASLt/build/release/clients/staging`.

The next section covers a brief explanation and the usage of each hipBLASLt client.

11.1 hipblaslt-test

`hipblaslt-test` is the main regression gtest for hipBLASLt. All test items must pass.

To run full test items:

```
./hipblaslt-test
```

To run partial test items with filter:

```
./hipblaslt-test --gtest_filter=<test pattern>
```

For demo “quick” test:

```
./hipblaslt-test --gtest_filter=*quick*
```

11.2 hipblaslt-bench

`hipblaslt-bench` is used to measure performance and to verify the correctness of hipBLASLt functions.

It has a command line interface.

To run fp32 GEMM with validation, see command:

```
./hipblaslt-bench --precision f32_r -v  
transA,transB,M,N,K,alpha,lda, stride_a,beta,ldb, stride_b,ldc, stride_c,ldd, stride_d,d_  
→type,compute_type,activation_type,bias_vector,hipblaslt-Gflops,us  
N,N,128,128,128,1,128,16384,0,128,16384,128,16384,128,16384,f32_r,f32_r,none,0, 415.278, 1  
→10.
```

For more information, see command:

```

./hipblaslt-bench --help
--sizeM |-m <value>      Specific matrix size: the number of rows or columns in matrix.
↳                          (Default value is: 128)
--sizeN |-n <value>      Specific matrix the number of rows or columns in matrix
↳                          (Default value is: 128)
--sizeK |-k <value>      Specific matrix size: the number of columns in A and rows in
↳ B.                        (Default value is: 128)
--lda <value>            Leading dimension of matrix A.
--ldb <value>            Leading dimension of matrix B.
--ldc <value>            Leading dimension of matrix C.
-- added <value>         Leading dimension of matrix D.
-- added <value>         Leading dimension of matrix E.
--any_stride             Do not modify input strides based on leading dimensions
--stride_a <value>       Specific stride of strided_batched matrix A, second dimension
↳ * leading dimension.
--stride_b <value>       Specific stride of strided_batched matrix B, second dimension
↳ * leading dimension.
--stride_c <value>       Specific stride of strided_batched matrix C, second dimension
↳ * leading dimension.
--stride_d <value>       Specific stride of strided_batched matrix D, second dimension
↳ * leading dimension.
--stride_e <value>       Specific stride of strided_batched matrix E, second dimension
↳ * leading dimension.
--alpha <value>          specifies the scalar alpha
↳                          (Default value is: 1)
--beta <value>           specifies the scalar beta
↳                          (Default value is: 0)
--function |-f <value>   BLASLt function to test. Options: matmul
↳                          (Default value is: matmul)
--precision |-r <value>  Precision of matrix A,B,C,D Options: f32_r,f16_r,bf16_r,f64_
↳ r,i32_r,i8_r           (Default value is: f16_r)
--a_type <value>         Precision of matrix A. Options: f32_r,f16_r,bf16_r
--b_type <value>         Precision of matrix B. Options: f32_r,f16_r,bf16_r
--c_type <value>         Precision of matrix C. Options: f32_r,f16_r,bf16_r
--d_type <value>         Precision of matrix D. Options: f32_r,f16_r,bf16_r
--compute_type <value>  Precision of computation. Options: s,f32_r,x,xf32_r,f64_r,i32_
↳ r                       (Default value is: f32_r)
--scale_type <value>     Precision of scalar. Options: f16_r,bf16_r
--initialization <value> Initialize matrix data.Options: rand_int, trig_float,
↳ hpl(floating), special, zero (Default value is: hpl)
--transA <value>         N = no transpose, T = transpose
↳                          (Default value is: N)
--transB <value>         N = no transpose, T = transpose
↳                          (Default value is: N)
--batch_count <value>    Number of matrices. Only applicable to batched and strided_
↳ batched routines       (Default value is: 1)
--HMM                    Parameter requesting the use of HipManagedMemory
--verify |-v            Validate GPU results with CPU?
--iters |-i <value>     Iterations to run inside timing loop
↳                          (Default value is: 10)
--cold_iters |-j <value> Cold Iterations to run before entering the timing loop
↳                          (Default value is: 2)
--algo_method <value>    Use different algorithm search API. Options: heuristic, all,

```

(continues on next page)

(continued from previous page)

```

↪index. (Default value is: heuristic)
--solution_index <value> Used with --algo_method 2. Specify solution index to use in.
↪benchmark. (Default value is: -1)
--requested_solution <value> Requested solution num. Set to -1 to get all solutions.
↪Only valid when algo_method is set to heuristic. (Default value is: 1)
--activation_type <value> Options: None, gelu, relu
↪ (Default value is: none)
--activation_arg1 <value> Reserved.
↪ (Default value is: 0)
--activation_arg2 <value> Reserved.
↪ (Default value is: inf)
--bias_type <value> Precision of bias vector.Options: f16_r,bf16_r,f32_r,
↪default(same with D type)
--bias_source <value> Choose bias source: a, b, d
↪ (Default value is: d)
--bias_vector Apply bias vector
--scaleA Apply scale for A buffer
--scaleB Apply scale for B buffer
--scaleAlpha_vector Apply scaleAlpha vector
--amaxScaleA Apple scale for A buffer by abs max of A buffer
--amaxScaleB Apple scale for B buffer by abs max of B buffer
--use_e Apply AUX output/ gradient input
--gradient Enable gradient
--grouped_gemm Use grouped_gemm.
--use_user_args Use UserArguments located in device memory for grouped gemm.
--device <value> Set default device to be used for subsequent program runs
↪ (Default value is: 0)
--c_equal_d C and D are stored in same memory
--workspace <value> Set fixed workspace memory size instead of using hipblaslt.
↪managed memory (Default value is: 0)
--log_function_name Function name precedes other items.
--function_filter <value> Simple strstr filter on function name only without wildcards
--api_method <value> Use extension API. c: C style API. mix: declaration with C.
↪hipblasLtMatmul Layout/Desc but set, initialize, and run the problem with C++.
↪extension API. cpp: Using C++ extension API only. Options: c, mix, cpp. (Default
↪value is: c)
--print_kernel_info Print solution, kernel name and solution index.
--rotating <value> Use rotating memory blocks for each iteration, size in MB.
↪ (Default value is: 0)
--use_gpu_timer Use hipEventElapsedTime to profile elapsed time.
↪ (Default value is: false)
--splitk <value> [Tuning parameter] Set split K for a solution, 0 is use
↪solution's default value. (Only support GEMM + api_method mix or cpp)
--wgm <value> [Tuning parameter] Set workgroup mapping for a solution, 0 is
↪use solution's default value. (Only support GEMM + api_method mix or cpp)
--help |-h produces this help message
--version <value> Prints the version number

```


LICENSE

Copyright (C) 2024 - 2025 Advanced Micro Devices, Inc.

Permission is hereby granted, free of charge, to any person obtaining a copy of this software and associated documentation files (the “Software”), to deal in the Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software, and to permit persons to whom the Software is furnished to do so, subject to the following conditions:

The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Software.

THE SOFTWARE IS PROVIDED “AS IS”, WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.

H

- hipblaslt_ext::Gemm (C++ class), 37
- hipblaslt_ext::Gemm::Gemm (C++ function), 38
- hipblaslt_ext::Gemm::setProblem (C++ function), 38–40
- hipblaslt_ext::GemmEpilogue (C++ struct), 29
- hipblaslt_ext::GemmEpilogue::aux_ld (C++ member), 29
- hipblaslt_ext::GemmEpilogue::aux_stride (C++ member), 29
- hipblaslt_ext::GemmEpilogue::bias_data_type (C++ member), 29
- hipblaslt_ext::GemmEpilogue::mode (C++ member), 29
- hipblaslt_ext::GemmEpilogueV2 (C++ class), 29
- hipblaslt_ext::GemmEpilogueV2::getAuxBatchStride (C++ function), 30
- hipblaslt_ext::GemmEpilogueV2::getAuxLeadingDimension (C++ function), 30
- hipblaslt_ext::GemmEpilogueV2::getBiasDataType (C++ function), 30
- hipblaslt_ext::GemmEpilogueV2::getMode (C++ function), 30
- hipblaslt_ext::GemmEpilogueV2::getScalingAType (C++ function), 30
- hipblaslt_ext::GemmEpilogueV2::getScalingBType (C++ function), 30
- hipblaslt_ext::GemmEpilogueV2::setAuxBatchStride (C++ function), 29
- hipblaslt_ext::GemmEpilogueV2::setAuxLeadingDimension (C++ function), 29
- hipblaslt_ext::GemmEpilogueV2::setBiasDataType (C++ function), 29
- hipblaslt_ext::GemmEpilogueV2::setMode (C++ function), 29
- hipblaslt_ext::GemmEpilogueV2::setScalingAType (C++ function), 29
- hipblaslt_ext::GemmEpilogueV2::setScalingBType (C++ function), 30
- hipblaslt_ext::GemmInputs (C++ struct), 30
- hipblaslt_ext::GemmInputs::a (C++ member), 30
- hipblaslt_ext::GemmInputs::alpha (C++ member), 30
- hipblaslt_ext::GemmInputs::aux (C++ member), 31
- hipblaslt_ext::GemmInputs::b (C++ member), 30
- hipblaslt_ext::GemmInputs::beta (C++ member), 30
- hipblaslt_ext::GemmInputs::bias (C++ member), 30
- hipblaslt_ext::GemmInputs::c (C++ member), 30
- hipblaslt_ext::GemmInputs::d (C++ member), 30
- hipblaslt_ext::GemmInputs::scaleA (C++ member), 31
- hipblaslt_ext::GemmInputs::scaleAlphaVec (C++ member), 31
- hipblaslt_ext::GemmInputs::scaleAux (C++ member), 31
- hipblaslt_ext::GemmInputs::scaleB (C++ member), 31
- hipblaslt_ext::GemmInputs::scaleC (C++ member), 31
- hipblaslt_ext::GemmInputs::scaledD (C++ member), 31
- hipblaslt_ext::GemmInputsV2 (C++ class), 31
- hipblaslt_ext::GemmInputsV2::getA (C++ function), 32
- hipblaslt_ext::GemmInputsV2::getAlpha (C++ function), 32
- hipblaslt_ext::GemmInputsV2::getAmaxD (C++ function), 33
- hipblaslt_ext::GemmInputsV2::getAux (C++ function), 33
- hipblaslt_ext::GemmInputsV2::getB (C++ function), 32
- hipblaslt_ext::GemmInputsV2::getBeta (C++ function), 32
- hipblaslt_ext::GemmInputsV2::getBias (C++ function), 32
- hipblaslt_ext::GemmInputsV2::getC (C++ function), 32
- hipblaslt_ext::GemmInputsV2::getD (C++ function), 32
- hipblaslt_ext::GemmInputsV2::getScaleA (C++

function), 32
 hipblaslt_ext::GemmInputsV2::getScaleAlphaVec (C++ *function*), 33
 hipblaslt_ext::GemmInputsV2::getScaleAux (C++ *function*), 33
 hipblaslt_ext::GemmInputsV2::getScaleB (C++ *function*), 32
 hipblaslt_ext::GemmInputsV2::getScaleC (C++ *function*), 32
 hipblaslt_ext::GemmInputsV2::getScaleD (C++ *function*), 32
 hipblaslt_ext::GemmInputsV2::setA (C++ *function*), 31
 hipblaslt_ext::GemmInputsV2::setAlpha (C++ *function*), 31
 hipblaslt_ext::GemmInputsV2::setAmaxD (C++ *function*), 32
 hipblaslt_ext::GemmInputsV2::setAux (C++ *function*), 32
 hipblaslt_ext::GemmInputsV2::setB (C++ *function*), 31
 hipblaslt_ext::GemmInputsV2::setBeta (C++ *function*), 31
 hipblaslt_ext::GemmInputsV2::setBias (C++ *function*), 31
 hipblaslt_ext::GemmInputsV2::setC (C++ *function*), 31
 hipblaslt_ext::GemmInputsV2::setD (C++ *function*), 31
 hipblaslt_ext::GemmInputsV2::setScaleA (C++ *function*), 32
 hipblaslt_ext::GemmInputsV2::setScaleAlphaVec (C++ *function*), 32
 hipblaslt_ext::GemmInputsV2::setScaleAux (C++ *function*), 32
 hipblaslt_ext::GemmInputsV2::setScaleB (C++ *function*), 32
 hipblaslt_ext::GemmInputsV2::setScaleC (C++ *function*), 32
 hipblaslt_ext::GemmInputsV2::setScaleD (C++ *function*), 32
 hipblaslt_ext::GemmInstance (C++ *class*), 34
 hipblaslt_ext::GemmInstance::algoGetHeuristic (C++ *function*), 34
 hipblaslt_ext::GemmInstance::GemmInstance (C++ *function*), 37
 hipblaslt_ext::GemmInstance::initialize (C++ *function*), 36
 hipblaslt_ext::GemmInstance::isAlgoSupported (C++ *function*), 35
 hipblaslt_ext::GemmInstance::run (C++ *function*), 37
 hipblaslt_ext::GemmPreference (C++ *class*), 33
 hipblaslt_ext::GemmPreference::getMaxWorkspaceBytes (C++ *function*), 33
 hipblaslt_ext::GemmPreference::setMaxWorkspaceBytes (C++ *function*), 33
 hipblaslt_ext::GemmPreferenceV2 (C++ *class*), 33
 hipblaslt_ext::GemmPreferenceV2::getMaxWorkspaceBytes (C++ *function*), 33
 hipblaslt_ext::GemmPreferenceV2::setMaxWorkspaceBytes (C++ *function*), 33
 hipblaslt_ext::GemmProblemType (C++ *struct*), 27
 hipblaslt_ext::GemmProblemType::op_a (C++ *member*), 27
 hipblaslt_ext::GemmProblemType::op_b (C++ *member*), 27
 hipblaslt_ext::GemmProblemType::type_a (C++ *member*), 27
 hipblaslt_ext::GemmProblemType::type_b (C++ *member*), 27
 hipblaslt_ext::GemmProblemType::type_c (C++ *member*), 28
 hipblaslt_ext::GemmProblemType::type_compute (C++ *member*), 28
 hipblaslt_ext::GemmProblemType::type_d (C++ *member*), 28
 hipblaslt_ext::GemmProblemTypeV2 (C++ *class*), 28
 hipblaslt_ext::GemmProblemTypeV2::getOpA (C++ *function*), 28
 hipblaslt_ext::GemmProblemTypeV2::getOpB (C++ *function*), 28
 hipblaslt_ext::GemmProblemTypeV2::getTypeA (C++ *function*), 28
 hipblaslt_ext::GemmProblemTypeV2::getTypeB (C++ *function*), 28
 hipblaslt_ext::GemmProblemTypeV2::getTypeC (C++ *function*), 28
 hipblaslt_ext::GemmProblemTypeV2::getTypeCompute (C++ *function*), 29
 hipblaslt_ext::GemmProblemTypeV2::getTypeD (C++ *function*), 29
 hipblaslt_ext::GemmProblemTypeV2::setOpA (C++ *function*), 28
 hipblaslt_ext::GemmProblemTypeV2::setOpB (C++ *function*), 28
 hipblaslt_ext::GemmProblemTypeV2::setTypeA (C++ *function*), 28
 hipblaslt_ext::GemmProblemTypeV2::setTypeB (C++ *function*), 28
 hipblaslt_ext::GemmProblemTypeV2::setTypeC (C++ *function*), 28
 hipblaslt_ext::GemmProblemTypeV2::setTypeCompute (C++ *function*), 28
 hipblaslt_ext::GemmProblemTypeV2::setTypeD (C++ *function*), 28
 hipblaslt_ext::GemmType (C++ *enum*), 27

<code>hipblasLtExt::GemmType::HIPBLASLT_GEMM</code> (C++ enumerator), 27	<code>hipblasLtMatmulAlgoGetHeuristic</code> (C++ function), 21
<code>hipblasLtExt::GemmType::HIPBLASLT_GROUPED_GEMM</code> (C++ enumerator), 27	<code>hipblasLtMatmulDesc_t</code> (C++ type), 9
<code>hipblasLtExt::getAlgosFromIndex</code> (C++ function), 47	<code>hipblasLtMatmulDescAttributes_t</code> (C++ enum), 9
<code>hipblasLtExt::getAllAlgos</code> (C++ function), 46	<code>hipblasLtMatmulDescAttributes_t::HIPBLASLT_MATMUL_DESC_A_S</code> (C++ enumerator), 9
<code>hipblasLtExt::getIndexFromAlgo</code> (C++ function), 47	<code>hipblasLtMatmulDescAttributes_t::HIPBLASLT_MATMUL_DESC_A_S</code> (C++ enumerator), 10
<code>hipblasLtExt::GroupedGemm</code> (C++ class), 41	<code>hipblasLtMatmulDescAttributes_t::HIPBLASLT_MATMUL_DESC_AMA</code> (C++ enumerator), 10
<code>hipblasLtExt::GroupedGemm::getDefaultValuesForDevice</code> (C++ function), 45	<code>hipblasLtMatmulDescAttributes_t::HIPBLASLT_MATMUL_DESC_B_S</code> (C++ enumerator), 9
<code>hipblasLtExt::GroupedGemm::GroupedGemm</code> (C++ function), 41, 42	<code>hipblasLtMatmulDescAttributes_t::HIPBLASLT_MATMUL_DESC_B_S</code> (C++ enumerator), 10
<code>hipblasLtExt::GroupedGemm::run</code> (C++ function), 45, 46	<code>hipblasLtMatmulDescAttributes_t::HIPBLASLT_MATMUL_DESC_BIA</code> (C++ enumerator), 9
<code>hipblasLtExt::GroupedGemm::setProblem</code> (C++ function), 42–44	<code>hipblasLtMatmulDescAttributes_t::HIPBLASLT_MATMUL_DESC_BIA</code> (C++ enumerator), 9
<code>hipblasLtExt::matmulIsAlgoSupported</code> (C++ function), 47	<code>hipblasLtMatmulDescAttributes_t::HIPBLASLT_MATMUL_DESC_C_S</code> (C++ enumerator), 9
<code>hipblasLtCreate</code> (C++ function), 15	<code>hipblasLtMatmulDescAttributes_t::HIPBLASLT_MATMUL_DESC_COM</code> (C++ enumerator), 10
<code>hipblasLtDestroy</code> (C++ function), 15	<code>hipblasLtMatmulDescAttributes_t::HIPBLASLT_MATMUL_DESC_COM</code> (C++ enumerator), 10
<code>hipblasLtEpilogue_t</code> (C++ enum), 7	<code>hipblasLtMatmulDescAttributes_t::HIPBLASLT_MATMUL_DESC_COM</code> (C++ enumerator), 10
<code>hipblasLtEpilogue_t::HIPBLASLT_EPILOGUE_BGRADA</code> (C++ enumerator), 8	<code>hipblasLtMatmulDescAttributes_t::HIPBLASLT_MATMUL_DESC_D_S</code> (C++ enumerator), 10
<code>hipblasLtEpilogue_t::HIPBLASLT_EPILOGUE_BGRADB</code> (C++ enumerator), 8	<code>hipblasLtMatmulDescAttributes_t::HIPBLASLT_MATMUL_DESC_EPI</code> (C++ enumerator), 9
<code>hipblasLtEpilogue_t::HIPBLASLT_EPILOGUE_BIAS</code> (C++ enumerator), 7	<code>hipblasLtMatmulDescAttributes_t::HIPBLASLT_MATMUL_DESC_EPI</code> (C++ enumerator), 10
<code>hipblasLtEpilogue_t::HIPBLASLT_EPILOGUE_DEFAULT</code> (C++ enumerator), 7	<code>hipblasLtMatmulDescAttributes_t::HIPBLASLT_MATMUL_DESC_EPI</code> (C++ enumerator), 10
<code>hipblasLtEpilogue_t::HIPBLASLT_EPILOGUE_DGELU</code> (C++ enumerator), 7	<code>hipblasLtMatmulDescAttributes_t::HIPBLASLT_MATMUL_DESC_EPI</code> (C++ enumerator), 10
<code>hipblasLtEpilogue_t::HIPBLASLT_EPILOGUE_DGELU_BGRAD</code> (C++ enumerator), 7	<code>hipblasLtMatmulDescAttributes_t::HIPBLASLT_MATMUL_DESC_EPI</code> (C++ enumerator), 10
<code>hipblasLtEpilogue_t::HIPBLASLT_EPILOGUE_GELU</code> (C++ enumerator), 7	<code>hipblasLtMatmulDescAttributes_t::HIPBLASLT_MATMUL_DESC_EPI</code> (C++ enumerator), 10
<code>hipblasLtEpilogue_t::HIPBLASLT_EPILOGUE_GELU_AUX</code> (C++ enumerator), 7	<code>hipblasLtMatmulDescAttributes_t::HIPBLASLT_MATMUL_DESC_EPI</code> (C++ enumerator), 10
<code>hipblasLtEpilogue_t::HIPBLASLT_EPILOGUE_GELU_AUX_BIAS</code> (C++ enumerator), 7	<code>hipblasLtMatmulDescAttributes_t::HIPBLASLT_MATMUL_DESC_MAX</code> (C++ enumerator), 10
<code>hipblasLtEpilogue_t::HIPBLASLT_EPILOGUE_GELU_BIAS</code> (C++ enumerator), 7	<code>hipblasLtMatmulDescAttributes_t::HIPBLASLT_MATMUL_DESC_POI</code> (C++ enumerator), 9
<code>hipblasLtEpilogue_t::HIPBLASLT_EPILOGUE_RELU</code> (C++ enumerator), 7	<code>hipblasLtMatmulDescAttributes_t::HIPBLASLT_MATMUL_DESC_TRA</code> (C++ enumerator), 9
<code>hipblasLtEpilogue_t::HIPBLASLT_EPILOGUE_RELU_BIAS</code> (C++ enumerator), 7	<code>hipblasLtMatmulDescAttributes_t::HIPBLASLT_MATMUL_DESC_TRA</code> (C++ enumerator), 9
<code>hipblasLtExtAMax</code> (C++ function), 66	<code>hipblasLtMatmulDescCreate</code> (C++ function), 17
<code>hipblasLtExtAMaxWithScale</code> (C++ function), 67	<code>hipblasLtMatmulDescDestroy</code> (C++ function), 18
<code>hipblasLtExtLayerNorm</code> (C++ function), 66	<code>hipblasLtMatmulDescGetAttribute</code> (C++ function), 18
<code>hipblasLtExtSoftmax</code> (C++ function), 65	<code>hipblasLtMatmulDescSetAttribute</code> (C++ function), 18
<code>hipblasLtHandle_t</code> (C++ type), 8	<code>hipblasLtMatmulHeuristicResult_t</code> (C++ struct), 11
<code>hipblasLtMatmul</code> (C++ function), 21	<code>hipblasLtMatmulPreference_t</code> (C++ type), 11
<code>hipblasLtMatmulAlgo_t</code> (C++ struct), 8	<code>hipblasLtMatmulPreferenceAttributes_t</code> (C++

enum), 11 (C++ *enumerator*), 8

hipblasLtMatmulPreferenceAttributes_t::HIPBLASLT_MATMUL_PREF_MAX
(C++ *enumerator*), 11

hipblasLtMatmulPreferenceAttributes_t::HIPBLASLT_MATMUL_PREF_MAX_WORKSPACE_BYTES
(C++ *enumerator*), 11

hipblasLtMatmulPreferenceAttributes_t::HIPBLASLT_MATMUL_PREF_SEARCH_MODE
(C++ *enumerator*), 11

hipblasLtMatmulPreferenceCreate (C++ *function*),
19

hipblasLtMatmulPreferenceDestroy (C++ *func-*
tion), 19

hipblasLtMatmulPreferenceGetAttribute (C++
function), 20

hipblasLtMatmulPreferenceSetAttribute (C++
function), 20

hipblasLtMatrixLayout_t (C++ *type*), 11

hipblasLtMatrixLayoutAttribute_t (C++ *enum*),
12

hipblasLtMatrixLayoutAttribute_t::HIPBLASLT_MATRIX_LAYOUT_BATCH_COUNT
(C++ *enumerator*), 12

hipblasLtMatrixLayoutAttribute_t::HIPBLASLT_MATRIX_LAYOUT_COLS
(C++ *enumerator*), 12

hipblasLtMatrixLayoutAttribute_t::HIPBLASLT_MATRIX_LAYOUT_LD
(C++ *enumerator*), 12

hipblasLtMatrixLayoutAttribute_t::HIPBLASLT_MATRIX_LAYOUT_ORDER
(C++ *enumerator*), 12

hipblasLtMatrixLayoutAttribute_t::HIPBLASLT_MATRIX_LAYOUT_ROWS
(C++ *enumerator*), 12

hipblasLtMatrixLayoutAttribute_t::HIPBLASLT_MATRIX_LAYOUT_STRIDED_BATCH_OFFSET
(C++ *enumerator*), 12

hipblasLtMatrixLayoutAttribute_t::HIPBLASLT_MATRIX_LAYOUT_TYPE
(C++ *enumerator*), 12

hipblasLtMatrixLayoutCreate (C++ *function*), 16

hipblasLtMatrixLayoutDestroy (C++ *function*), 16

hipblasLtMatrixLayoutGetAttribute (C++ *func-*
tion), 17

hipblasLtMatrixLayoutSetAttribute (C++ *func-*
tion), 16

hipblasLtMatrixTransform (C++ *function*), 25

hipblasLtMatrixTransformDesc_t (C++ *type*), 13

hipblasLtMatrixTransformDescCreate (C++ *func-*
tion), 23

hipblasLtMatrixTransformDescDestroy (C++
function), 24

hipblasLtMatrixTransformDescGetAttribute
(C++ *function*), 24

hipblasLtMatrixTransformDescSetAttribute
(C++ *function*), 24

hipblasLtPointerMode_t (C++ *enum*), 8

hipblasLtPointerMode_t::HIPBLASLT_POINTER_MODE_ALPHA_DEVICE_VECTOR_BETA_HOST
(C++ *enumerator*), 8

hipblasLtPointerMode_t::HIPBLASLT_POINTER_MODE_DEVICE
(C++ *enumerator*), 8

hipblasLtPointerMode_t::HIPBLASLT_POINTER_MODE_HOST