



DECEMBER 5, 2014

NVRTC - CUDA Runtime Compilation

NVIDIA CORPORATION

V7.0

TABLE OF CONTENTS

1	Introduction	1
2	Getting Started.....	2
2.1	System Requirements	2
2.2	Installation	2
3	User Interface	3
3.1	API	3
3.1.1	Error Handling	3
3.1.2	General Information Query.....	4
3.1.3	Compilation.....	4
3.2	Supported Compile Options.....	9
4	Language	12
4.1	Execution Space	12
4.2	Separate Compilation	12
4.3	Integer Size.....	12
4.4	Predefined Macros.....	12
4.5	Predefined types	13
4.6	Builtin Functions	13
5	Basic Usage	14
6	Known Issues.....	17
7	Notes.....	17
Appendix A	Example: SAXPY.....	18
A.1	Code (saxpy.cpp).....	18
A.2	Build Instruction.....	21

1 INTRODUCTION

NVRTC is a runtime compilation library for CUDA C++. It accepts CUDA C++ source code in character string form and creates handles that can be used to obtain the PTX. The PTX string generated by NVRTC can be loaded by [cuModuleLoadData](#) and [cuModuleLoadDataEx](#), and linked with other modules by [cuLinkAddData](#) of the CUDA Driver API. This facility can often provide optimizations and performance not possible in a purely offline static compilation.

In the absence of NVRTC (or any runtime compilation support in CUDA), users needed to spawn a separate process to execute `nvcc` at runtime if they wished to implement runtime compilation in their applications or libraries, and, unfortunately, this approach has the following drawbacks:

- The compilation overhead tends to be higher than necessary, and
- End users are required to install `nvcc` and related tools which make it complicated to distribute applications that use runtime compilation.

NVRTC addresses these issues by providing a library interface that eliminates overhead associated with spawning separate processes, disk I/O, etc., while keeping application deployment simple.

NVRTC is a preview feature in CUDA 7.0 and any or all parts of this specification are subject to change in the next CUDA release.

2 GETTING STARTED

2.1 SYSTEM REQUIREMENTS

NVRTC requires the following system configuration:

- Operating System: 64-bit Linux, Windows, or Mac OS X.
- GPU: Any GPU with CUDA Compute Capability 2.0 or higher.
- 7.0 toolkit and driver.

2.2 INSTALLATION

NVRTC is part of the CUDA toolkit release and the components are organized as follows in the CUDA toolkit installation directory:

- On Windows:
 - include\nvrtc.h
 - bin\nvrtc64_70.dll
 - bin\nvrtc-builtins64_70.dll
 - lib\x64\nvrtc.lib
 - doc\pdf\NVRTC_User_Guide.pdf
- On Linux:
 - include/nvrtc.h
 - lib64/libnvrtc.so
 - lib64/libnvrtc.so.7.0
 - lib64/libnvrtc.so.7.0.0
 - lib64/libnvrtc-builtins.so
 - lib64/libnvrtc-builtins.so.7.0
 - lib64/libnvrtc-builtins.so.7.0.0
 - doc/pdf/NVRTC_User_Guide.pdf
- On Mac OS X:
 - include/nvrtc.h
 - lib/libnvrtc.dylib
 - lib/libnvrtc.7.0.dylib
 - lib/libnvrtc-builtins.dylib
 - lib/libnvrtc-builtins.7.0.dylib
 - doc/pdf/NVRTC_User_Guide.pdf

3 USER INTERFACE

3.1 API

This section presents the API of NVRTC for CUDA 7.0. Basic usage of the API is explained in Section 5. Note that the API may change in the production release based on user feedback.

3.1.1 ERROR HANDLING

NVRTC defines the following enumeration type and function for API call error handling.

3.1.1.1 `nvrtcResult`

The enumerated type `nvrtcResult` defines API call result codes with enumeration values shown below.

```
typedef enum {
    NVRTC_SUCCESS = 0,
    NVRTC_ERROR_OUT_OF_MEMORY = 1,
    NVRTC_ERROR_PROGRAM_CREATION_FAILURE = 2,
    NVRTC_ERROR_INVALID_INPUT = 3,
    NVRTC_ERROR_INVALID_PROGRAM = 4,
    NVRTC_ERROR_INVALID_OPTION = 5,
    NVRTC_ERROR_COMPILATION = 6,
    NVRTC_ERROR_BUILTIN_OPERATION_FAILURE = 7
} nvrtcResult;
```

NVRTC API functions return `nvrtcResult` to indicate the call result.

3.1.1.2 `nvrtcGetErrorString`

```
const char *nvrtcGetErrorString(nvrtcResult result)
```

`nvrtcGetErrorString` is a helper function that returns a string describing the given `nvrtcResult` code, e.g., `NVRTC_SUCCESS` to “`NVRTC_SUCCESS`”. For unrecognized enumeration values, it returns “`NVRTC_ERROR unknown`”.

PARAMETER

- `[in]` `result`
CUDA Online Compiler API result code.

RETURNS

- Message string for the given `nvrtcResult` code.

3.1.2 GENERAL INFORMATION QUERY

NVRTC defines the following function for general information query.

3.1.2.1 `nVRTCVersion`

```
nVRTCResult nVRTCVersion(int *major, int *minor)
```

`nVRTCVersion` sets the output parameters `major` and `minor` with the CUDA Online Compiler version number.

PARAMETERS

- `[out] major`
CUDA Online Compiler major version number.
- `[out] minor`
CUDA Online Compiler minor version number.

RETURNS

- `NVRTC_SUCCESS`
- `NVRTC_ERROR_INVALID_INPUT`

3.1.3 COMPILATION

NVRTC defines the following type and functions for actual compilation.

3.1.3.1 `nVRTCProgram`

```
typedef struct _nVRTCProgram *nVRTCProgram
```

`nVRTCProgram` is the unit of compilation, and an opaque handle for a program. To compile a CUDA program string, an instance of `nVRTCProgram` must be created first with `nVRTCCreateProgram` (see Section 3.1.3.2), then compiled with `nVRTCCompileProgram` (see Section 3.1.3.4).

3.1.3.2 `nVRTCCreateProgram`

```
nVRTCResult nVRTCCreateProgram(nVRTCProgram *prog,  
                               const char *src,  
                               const char *name,  
                               int numHeaders,  
                               const char **headers,  
                               const char **includeNames)
```

`nVRTCCreateProgram` creates an instance of `nVRTCProgram` with the given input parameters, and sets the output parameter `prog` with it.

PARAMETERS

- *[out]* prog
CUDA Online Compiler program.
- *[in]* src
CUDA program source.
- *[in]* name
CUDA program name.
name can be NULL; “default_program” is used when it is NULL.
- *[in]* numHeaders
Number of headers used.
numHeaders must be greater than or equal to 0.
- *[in]* headers
Sources of the headers.
headers can be NULL when numHeaders is 0.
- *[in]* includeNames
Name of each header by which they can be included in the CUDA program source.
includeNames can be NULL when numHeaders is 0.

RETURNS

- NVRTC_SUCCESS
- NVRTC_ERROR_OUT_OF_MEMORY
- NVRTC_ERROR_PROGRAM_CREATION_FAILURE
- NVRTC_ERROR_INVALID_INPUT
- NVRTC_ERROR_INVALID_PROGRAM

SEE ALSO

- Section 3.1.3.3 for `nVRTC_DestroyProgram`.

3.1.3.3 `nVRTC_DestroyProgram`

```
nVRTC_Result nVRTC_DestroyProgram(nVRTC_Program *prog)
```

`nVRTC_DestroyProgram` destroys the given program.

PARAMETER

- *[in]* prog
CUDA Online Compiler program.

RETURNS

- NVRTC_SUCCESS
- NVRTC_ERROR_INVALID_PROGRAM

SEE ALSO

- Section 3.1.3.2 for `nVRTCCreateProgram`.

3.1.3.4 `nVRTCCompileProgram`

```
nVRTCResult nVRTCCompileProgram(nVRTCProgram prog,
                                int numOptions,
                                const char **options)
```

`nVRTCCompileProgram` compiles the given program. It supports compile options listed in Section 3.2.

PARAMETERS

- *[in]* `prog`
CUDA Online Compiler program.
- *[in]* `numOptions`
Number of compiler options passed.
- *[in]* `options`
Compiler options in the form of C string array.
`options` can be `NULL` when `numOptions` is 0.

RETURNS

- `NVRTC_SUCCESS`
- `NVRTC_ERROR_OUT_OF_MEMORY`
- `NVRTC_ERROR_INVALID_INPUT`
- `NVRTC_ERROR_INVALID_PROGRAM`
- `NVRTC_ERROR_INVALID_OPTION`
- `NVRTC_ERROR_COMPILATION`
- `NVRTC_ERROR_BUILTIN_OPERATION_FAILURE`

3.1.3.5 `nVRTCGetPTXSize`

```
nVRTCResult nVRTCGetPTXSize(nVRTCProgram prog,
                             size_t *ptxSizeRet)
```

`nVRTCGetPTXSize` sets `ptxSizeRet` with the size of the PTX generated by the previous compilation of `prog` (including the trailing `NULL`).

PARAMETERS

- *[in]* `prog`
CUDA Online Compiler program.

- *[out]* `ptxSizeRet`
Size of the generated PTX (including the trailing `NULL`).

RETURNS

- `NVRTC_SUCCESS`
- `NVRTC_ERROR_INVALID_INPUT`
- `NVRTC_ERROR_INVALID_PROGRAM`

SEE ALSO

- Section 3.1.3.6 for `nVRTCGetPTX`.

3.1.3.6 `nVRTCGetPTX`

```
nVRTCResult nVRTCGetPTX(nVRTCProgram prog, char *ptx)
```

`nVRTCGetPTX` stores the PTX generated by the previous compilation of `prog` in the memory pointed by `ptx`.

PARAMETERS

- *[in]* `prog`
CUDA Online Compiler program.
- *[out]* `ptx`
Compiled result.

RETURNS

- `NVRTC_SUCCESS`
- `NVRTC_ERROR_INVALID_INPUT`
- `NVRTC_ERROR_INVALID_PROGRAM`

SEE ALSO

- Section 3.1.3.5 for `nVRTCGetPTXSize`.

3.1.3.7 `nVRTCGetProgramLogSize`

```
nVRTCResult nVRTCGetProgramLogSize(nVRTCProgram prog,  
                                   size_t *logSizeRet)
```

`nVRTCGetProgramLogSize` sets `logSizeRet` with the size of the log generated by the previous compilation of `prog` (including the trailing `NULL`). Note that compilation log may be generated with warnings and informative messages, even when the compilation of `prog` succeeds.

PARAMETERS

- *[in]* prog
CUDA Online Compiler program.
- *[out]* logSizeRet
Size of the compilation log (including the trailing NULL).

RETURNS

- NVRTC_SUCCESS
- NVRTC_ERROR_INVALID_INPUT
- NVRTC_ERROR_INVALID_PROGRAM

SEE ALSO

- Section 3.1.3.8 for `nVRTCGetProgramLog`.

3.1.3.8 `nVRTCGetProgramLog`

```
nVRTCResult nVRTCGetProgramLog(nVRTCProgram prog, char *log)
```

`nVRTCGetProgramLog` stores the log generated by the previous compilation of `prog` in the memory pointed by `log`.

PARAMETERS

- *[in]* prog
CUDA Online Compiler program.
- *[out]* log
Compilation log.

RETURNS

- NVRTC_SUCCESS
- NVRTC_ERROR_INVALID_INPUT
- NVRTC_ERROR_INVALID_PROGRAM

SEE ALSO

- Section 3.1.3.7 for `nVRTCGetProgramLogSize`.

3.2 SUPPORTED COMPILE OPTIONS

NVRTC supports the compile options below. Option names with two preceding dashes ('--') are *long option names* and option names with one preceding dash ('-') are *short option names*. Short option names can be used instead of long option names. When a compile option takes an argument, an assignment operator ('=') is used to separate the compile option argument from the compile option name, e.g., "--gpu-architecture=compute_20". Alternatively, the compile option name and the argument can be specified in separate strings without an assignment operator, .e.g., "--gpu-architecture" and "compute_20". Single-character short option names, such as -D, -U, and -I, do not require an assignment operator, and the compile option name and the argument can be present in the same string with or without spaces between them. For instance, "-D=<macrodef>", "-D<macrodef>", and "-D <macrodef>" are all supported.

- **Compilation targets**
 - --gpu-architecture=<GPU architecture> (-arch=<GPU architecture>)
Specify the name of the class of GPU architectures for which the input must be compiled.
 - Valid GPU architectures:
 - compute_20
 - compute_30
 - compute_35
 - compute_50
 - Default: compute_20
- **Separate compilation and whole-program compilation**
 - --device-c (-dc)
Generate relocatable code that can be linked with other relocatable device code. It is equivalent to --relocatable-device-code=true.
 - --device-w (-dw)
Generate non-relocatable code. It is equivalent to --relocatable-device-code=false.
 - --relocatable-device-code=[true, false] (-rdc)
Enable (disable) the generation of relocatable device code.
 - Default: false
- **Debugging support**
 - --device-debug (-G)
Generate debug information.
 - --generate-line-info (-lineinfo)
Generate line-number information.
- **Code generation**
 - --maxrregcount=<N> (-maxrregcount=<N>)
Specify the maximum amount of registers that GPU functions can use. Until a function-

specific limit, a higher value will generally increase the performance of individual GPU threads that execute this function. However, because thread registers are allocated from a global register pool on each GPU, a higher value of this option will also reduce the maximum thread block size, thereby reducing the amount of thread parallelism. Hence, a good `maxrregcount` value is the result of a trade-off. If this option is not specified, then no maximum is assumed. Value less than the minimum registers required by ABI will be bumped up by the compiler to ABI minimum limit.

- `--ftz=[true, false] (-ftz=[true, false])`
When performing single-precision floating-point operations, flush denormal values to zero or preserve denormal values. `--use_fast_math` implies `--ftz=true`.
 - Default: false
- `--prec-sqrt=[true, false] (-prec-sqrt=[true, false])`
For single-precision floating-point square root, use IEEE round-to-nearest mode or use a faster approximation. `--use_fast_math` implies `--prec-sqrt=false`.
 - Default: true
- `--prec-div=[true, false] (-prec-div=[true, false])`
For single-precision floating-point division and reciprocals, use IEEE round-to-nearest mode or use a faster approximation. `--use_fast_math` implies `--prec-div=false`.
 - Default: true
- `--fmad=[true, false] (-fmad=[true, false])`
Enables (disables) the contraction of floating-point multiplies and adds/subtracts into floating-point multiply-add operations (FMAD, FFMA, or DFMA). `--use_fast_math` implies `--fmad=true`.
 - Default: true
- `--use_fast_math (-use_fast_math)`
Make use of fast math operations. `--use_fast_math` implies `--ftz=true --prec-div=false --prec-sqrt=false --fmad=true`.
- Preprocessing
 - `--define-macro=<macrodef> (-D<macrodef>)`
macrodef can be either *name* or *name=definitions*.
 - *name*
Predefine *name* as a macro with definition 1.
 - *name=definition*
The contents of *definition* are tokenized and preprocessed as if they appeared during translation phase three in a `#define` directive. In particular, the *definition* will be truncated by embedded new line characters.
 - `--undefine-macro=<name> (-U<name>)`
Cancel any previous definition of *name*.

- `--include-path=<dir> (-I<dir>)`
Add the directory *dir* to the list of directories to be searched for headers. These paths are searched after the list of headers given to `nVRTCCreateProgram`.
- `--pre-include=<header> (-include=<header>)`
Preinclude *header* during preprocessing.
- Language Dialect
 - `--std=c++11 (-std=c++11)`
Set the language dialect to C++11.
 - `--builtin-move-forward=[true, false] (-builtin-move-forward=[true, false])`
Provide builtin definitions of `std::move` and `std::forward`, when C++11 language dialect is selected.
 - Default: true
 - `--builtin-initializer-list=[true, false] (-builtin-initializer-list=[true, false])`
Provide builtin definitions of `std::initializer_list` class and member functions, when C++11 language dialect is selected.
 - Default: true
- Misc.
 - `--disable-warnings (-w)`
Inhibit all warning messages.
 - `--restrict (-restrict)`
Programmer assertion that all kernel pointer parameters are restrict pointers.
 - `--device-as-default-execution-space (-default-device)`
Treat entities with no execution space annotation as `__device__` entities. Function declarations without explicit execution space annotations, such as `__global__`, `__host__`, and `__device__`, will be treated as having an implicit `__device__` annotation. Namespace scope variable declarations without an explicit memory space annotation, such as `__device__`, `__constant__`, `__shared__`, and `__managed__`, will be treated as having an implicit `__device__` annotation.

4 LANGUAGE

Unlike the offline `nvcc` compiler, NVRTC is meant for compiling only device CUDA C++ code. It does not accept host code or host compiler extensions in the input code, unless otherwise noted.

4.1 EXECUTION SPACE

NVRTC uses `__host__` as the default execution space, and it generates an error if it encounters any host code in the input. That is, if the input contains entities with explicit `__host__` annotations or no execution space annotation, NVRTC will emit an error. `__host__ __device__` functions are treated as device functions.

NVRTC provides a compile option, `--device-as-default-execution-space (-default-device)`, that enables an alternative compilation mode, in which entities with no execution space annotations are treated as `__device__` entities. See the description of the compile option for details.

4.2 SEPARATE COMPILATION

NVRTC itself does not provide any linker. Users can, however, use [cuLinkAddData](#) in the CUDA Driver API to link the generated relocatable PTX code with other relocatable code. To generate relocatable PTX code, the compile option `--relocatable-device-code=true` or `--device-c` is required.

4.3 INTEGER SIZE

Different operating systems define integer type sizes differently. Linux implement LP64, and Windows implements LLP64.

	short	int	long	long long	pointers / size_t
LLP64	16	32	32	64	64
LP64	16	32	64	64	64

Table 1 Integer sizes in bits for LLP64 and LP64

NVRTC implements LP64 on Linux and LLP64 on Windows.

4.4 PREDEFINED MACROS

- `__CUDACC_RTC__`
Useful for distinguishing between runtime and offline (`nvcc`) compilation in user code.
- `__CUDACC__`
Defined with same semantics as with offline `nvcc`.
- `__CUDA_ARCH__`
Defined with same semantics as with offline `nvcc`.
- `NULL`
null pointer constant.

- `__cplusplus`

4.5 PREDEFINED TYPES

- `clock_t`
- `size_t`
- `ptrdiff_t`
- Predefined types such as `dim3`, `char4`, etc. that are available in the CUDA Runtime headers when compiling offline with `nvcc` are also available, unless otherwise noted.

4.6 BUILTIN FUNCTIONS

Builtin functions provided by the CUDART headers when compiling offline with `nvcc` are available, unless otherwise noted.

5 BASIC USAGE

This section of the document uses a simple example, *Single-Precision A·X Plus Y (SAXPY)*, shown in Listing 1 to explain what is involved in runtime compilation with NVRTC. For brevity and readability, error checks on the API return values are not shown. The complete code listing is available in Appendix A.

```
const char *saxpy = "                                \n\  
extern \"C\" __global__                               \n\  
void saxpy(float a, float *x, float *y, float *out, size_t n) \n\  
{                                                    \n\  
    size_t tid = blockIdx.x * blockDim.x + threadIdx.x; \n\  
    if (tid < n) {                                    \n\  
        out[tid] = a * x[tid] + y[tid];              \n\  
    }                                                \n\  
}                                                    \n\";
```

Listing 1 CUDA source string for SAXPY

First, an instance of `nVRTCProgram` needs to be created. Listing 2 shows creation of `nVRTCProgram` for SAXPY. As SAXPY does not require any header, 0 is passed as `numHeaders`, and `NULL` as `headers` and `includeNames`.

```
nVRTCProgram prog;  
  
nVRTCCreateProgram(&prog, // prog  
                  saxpy, // buffer  
                  "saxpy.cu", // name  
                  0, // numHeaders  
                  NULL, // headers  
                  NULL); // includeNames
```

Listing 2 nVRTCProgram creation for SAXPY

If SAXPY had any `#include` directives, the contents of the files that are `#include'd` can be passed as elements of `headers`, and their names as elements of `includeNames`. For example, `#include <foo.h>` and `#include <bar.h>` would require 2 as `numHeaders`, { "`<contents of foo.h>`", "`<contents of bar.h>`" } as `headers`, and { "`foo.h`", "`bar.h`" } as `includeNames` (`<contents of foo.h>` and `<contents of bar.h>` must be replaced by the actual contents of `foo.h` and `bar.h`). Alternatively, the compile option `-I` can be used if the header is guaranteed to exist in the file system at runtime.

Once the instance of `nVRTCProgram` for compilation is created, it can be compiled by `nVRTCCompileProgram` as shown in Listing 3. Two compile options are used in this example, `--gpu-architecture=compute_20` and `--fmad=false`, to generate code for the `compute_20` architecture and to disable the contraction of floating-point multiplies and adds/subtracts into floating-point multiply-add operations. Other combinations of compile options can be used as needed and Section 3.2 lists valid compile options.

```
const char *opts[] = {"--gpu-architecture=compute_20",
                    "--fmad=false"};

nVRTCCompileProgram(prog,    // prog
                  2,        // numOptions
                  opts);    // options
```

Listing 3 Compilation of SAXPY for `compute_20` with FMAD enabled

After the compilation completes, users can obtain the program compilation log and the generated PTX as Listing 4 shows. NVRTC does not generate valid PTX when the compilation fails, and it may generate program compilation log even when the compilation succeeds if needed.

```
// Obtain compilation log from the program.
size_t logSize;
nVRTCGetProgramLogSize(prog, &logSize);
char *log = new char[logSize];
nVRTCGetProgramLog(prog, log);

// Obtain PTX from the program.
size_t ptxSize;
nVRTCGetPTXSize(prog, &ptxSize);
char *ptx = new char[ptxSize];
nVRTCGetPTX(prog, ptx);
```

Listing 4 Obtaining generated PTX and program compilation log

A `nVRTCProgram` can be compiled by `nVRTCCompileProgram` multiple times with different compile options, and users can only retrieve the PTX and the log generated by the last compilation.

When the instance of `nVRTCProgram` is no longer needed, it can be destroyed by `nVRTCDestroyProgram` as shown in Listing 6.

```
nVRTC_DestroyProgram(&prog);
```

Listing 6 Destruction of nVRTCProgram

The generated PTX can be further manipulated by the CUDA Driver API for execution or linking. Listing 5 shows an example code sequence for execution of the generated PTX.

```
CUdevice cuDevice;
CUcontext context;
CUmodule module;
CUfunction kernel;
cuInit(0);
cuDeviceGet(&cuDevice, 0);
cuCtxCreate(&context, 0, cuDevice);
cuModuleLoadDataEx(&module, ptx, 0, 0, 0);
cuModuleGetFunction(&kernel, module, "saxpy");

size_t n = NUM_THREADS * NUM_BLOCKS;
size_t bufferSize = n * sizeof(float);
float a = ...;
float *hX = ..., *hY = ..., *hOut = ...;
CUdeviceptr dX, dY, dOut;
cuMemAlloc(&dX, bufferSize);
cuMemAlloc(&dY, bufferSize);
cuMemAlloc(&dOut, bufferSize);
cuMemcpyHtoD(dX, hX, bufferSize);
cuMemcpyHtoD(dY, hY, bufferSize);

void *args[] = { &a, &dX, &dY, &dOut, &n };
cuLaunchKernel(kernel,
               NUM_THREADS, 1, 1, // grid dim
               NUM_BLOCKS, 1, 1, // block dim
               0, NULL, // shared mem and stream
               args, // arguments
               0);

cuCtxSynchronize();
cuMemcpyDtoH(hOut, dOut, bufferSize);
```

Listing 5 Execution of SAXPY using the PTX generated by NVRTC

6 KNOWN ISSUES

The following CUDA C++ features are not yet implemented when compiling with NVRTC:

- Dynamic parallelism (kernel launches from within device code).
- Literals of type `wchar_t`, `char16_t` and `char32_t`.

7 NOTES

- Template instantiations: Since NVRTC compiles only device code, all templates must be instantiated within device code (including `__global__` function templates).
- NVRTC follows the IA64 ABI; function names will be mangled unless the function declaration is marked with `extern "C"` linkage. To look up a kernel with the driver API, users must provide a string name, which is hard if the name is mangled. Using `extern "C"` linkage for a `__global__` function will allow use of the unmangled name when using the driver API to find the kernel's address.

APPENDIX A EXAMPLE: SAXPY

A.1 CODE (SAXPY.CPP)

```
#include <nVRTC.h>
#include <cuda.h>
#include <iostream>

#define NUM_THREADS 128
#define NUM_BLOCKS 32
#define NVRTC_SAFE_CALL(x) \
do { \
    nVRTCResult result = x; \
    if (result != NVRTC_SUCCESS) { \
        std::cerr << "\nerror: " #x " failed with error " \
        << nVRTCGetErrorString(result) << '\n'; \
        exit(1); \
    } \
} while(0)
#define CUDA_SAFE_CALL(x) \
do { \
    CUresult result = x; \
    if (result != CUDA_SUCCESS) { \
        const char *msg; \
        cuGetErrorName(result, &msg); \
        std::cerr << "\nerror: " #x " failed with error " \
        << msg << '\n'; \
        exit(1); \
    } \
} while(0)

const char *saxpy = " \n\
extern \"C\" __global__ \n\
void saxpy(float a, float *x, float *y, float *out, size_t n) \n\
{ \n\
    size_t tid = blockIdx.x * blockDim.x + threadIdx.x; \n\
    if (tid < n) { \n\
        out[tid] = a * x[tid] + y[tid]; \n\
    } \n\
} \n\
\n";

int main()
{
    // Create an instance of nVRTCProgram with the SAXPY code string.
    nVRTCProgram prog;
```

```

NVRTC_SAFE_CALL(
    nVRTCCreateProgram(&prog,          // prog
                      saxpy,          // buffer
                      "saxpy.cu",     // name
                      0,               // numHeaders
                      NULL,            // headers
                      NULL));          // includeNames

// Compile the program for compute_20 with fmad disabled.
const char *opts[] = {"--gpu-architecture=compute_20",
                     "--fmad=false"};

nVRTCResult compileResult = nVRTCCompileProgram(prog, // prog
                                                2,     // numOptions
                                                opts); // options

// Obtain compilation log from the program.
size_t logSize;
NVRTC_SAFE_CALL(nVRTCGetProgramLogSize(prog, &logSize));
char *log = new char[logSize];
NVRTC_SAFE_CALL(nVRTCGetProgramLog(prog, log));
std::cout << log << '\n';
delete[] log;
if (compileResult != NVRTC_SUCCESS) {
    exit(1);
}

// Obtain PTX from the program.
size_t ptxSize;
NVRTC_SAFE_CALL(nVRTCGetPTXSize(prog, &ptxSize));
char *ptx = new char[ptxSize];
NVRTC_SAFE_CALL(nVRTCGetPTX(prog, ptx));

// Destroy the program.
NVRTC_SAFE_CALL(nVRTCDestroyProgram(&prog));

// Load the generated PTX and get a handle to the SAXPY kernel.
CUdevice cuDevice;
CUcontext context;
CUmodule module;
CUfunction kernel;
CUDA_SAFE_CALL(cuInit(0));
CUDA_SAFE_CALL(cuDeviceGet(&cuDevice, 0));
CUDA_SAFE_CALL(cuCtxCreate(&context, 0, cuDevice));
CUDA_SAFE_CALL(cuModuleLoadDataEx(&module, ptx, 0, 0, 0));
CUDA_SAFE_CALL(cuModuleGetFunction(&kernel, module, "saxpy"));

```

```

// Generate input for execution, and create output buffers.
size_t n = NUM_THREADS * NUM_BLOCKS;
size_t bufferSize = n * sizeof(float);
float a = 5.1f;
float *hX = new float[n], *hY = new float[n], *hOut = new float[n];
for (size_t i = 0; i < n; ++i) {
    hX[i] = static_cast<float>(i);
    hY[i] = static_cast<float>(i * 2);
}
CUdeviceptr dX, dY, dOut;
CUDA_SAFE_CALL(cuMemAlloc(&dX, bufferSize));
CUDA_SAFE_CALL(cuMemAlloc(&dY, bufferSize));
CUDA_SAFE_CALL(cuMemAlloc(&dOut, bufferSize));
CUDA_SAFE_CALL(cuMemcpyHtoD(dX, hX, bufferSize));
CUDA_SAFE_CALL(cuMemcpyHtoD(dY, hY, bufferSize));

// Execute SAXPY.
void *args[] = { &a, &dX, &dY, &dOut, &n };
CUDA_SAFE_CALL(
    cuLaunchKernel(kernel,
                   NUM_THREADS, 1, 1,    // grid dim
                   NUM_BLOCKS, 1, 1,    // block dim
                   0, NULL,              // shared mem and stream
                   args, 0));            // arguments
CUDA_SAFE_CALL(cuCtxSynchronize());

// Retrieve and print output.
CUDA_SAFE_CALL(cuMemcpyDtoH(hOut, dOut, bufferSize));
for (size_t i = 0; i < n; ++i) {
    std::cout << a << " * " << hX[i] << " + " << hY[i]
               << " = " << hOut[i] << '\n';
}

// Release resources.
CUDA_SAFE_CALL(cuMemFree(dX));
CUDA_SAFE_CALL(cuMemFree(dY));
CUDA_SAFE_CALL(cuMemFree(dOut));
CUDA_SAFE_CALL(cuModuleUnload(module));
CUDA_SAFE_CALL(cuCtxDestroy(context));
delete[] hX;
delete[] hY;
delete[] hOut;

return 0;
}

```

A.2 BUILD INSTRUCTION

Assuming the environment variable `CUDA_PATH` points to CUDA toolkit installation directory, build this example as:

- **Windows:**

```
cl.exe saxpy.cpp /Fesaxpy ^
    /I "%CUDA_PATH%\include ^
    "%CUDA_PATH%\lib\x64\nvrtc.lib ^
    "%CUDA_PATH%\lib\x64\cuda.lib
```

- **Linux:**

```
g++ saxpy.cpp -o saxpy \
    -I $CUDA_PATH/include \
    -L $CUDA_PATH/lib64 \
    -lnvrtc -lcuda \
    -Wl,-rpath,$CUDA_PATH/lib64
```

- **Mac OS X:**

```
clang++ saxpy.cpp -o saxpy \
    -I $CUDA_PATH/include \
    -L $CUDA_PATH/lib \
    -lnvrtc -framework CUDA \
    -Wl,-rpath,$CUDA_PATH/lib
```