

VkFFT -
VULKAN/CUDA/HIP/OPENCL/LEVEL
ZERO/METAL FAST FOURIER TRANSFORM
LIBRARY

API guide with examples

Dmitrii Tolmachev

October 2023, version 1.3.2

Contents

1	Introduction	3
2	Using the VkFFT API	4
2.1	Installing VkFFT	4
2.2	Fourier Transform Setup	8
2.3	Fourier Transform types and their definitions	8
2.3.1	C2C transforms	9
2.3.2	R2C/C2R transforms	9
2.3.3	R2R (DCT) transforms	9
2.4	Memory management, data layouts for different transforms	10
2.4.1	VkFFT buffers	10
2.4.2	VkFFT buffers strides. A special case of R2C/C2R transforms	11
2.5	VkFFT algorithms	12
2.5.1	Bluestein's algorithm	12
2.5.2	The Four-Step FFT algorithm	13
2.5.3	R2C/C2R FFTs	14
2.5.4	R2C/C2R multi-upload FFT algorithm	14
2.5.5	R2R Discrete Cosine/Sine Transforms	14
2.5.6	Register overutilization	14
2.5.7	Zero padding	15
2.5.8	Convolution and cross-correlation support	15
2.6	VkFFT accuracy	15
2.7	VkFFT additional memory allocations	16
2.7.1	LUT allocations	16
2.7.2	The Four-Step FFT algorithm - tempBuffer allocation	16
2.7.3	Bluestein's buffers allocation	16
2.8	VkFFT support for double-double emulation of quad precision	17
3	Runtime code optimization platform	18
3.1	Platform for GPU code generation	18
3.2	Platform structure	19
3.3	Container abstractions	20
4	VkFFT API Reference	21
4.1	Return value VkFFTResult	21
4.2	VkFFT application management functions	24
4.2.1	Function initializeVkFFT()	25
4.2.2	Function VkFFTAppend()	25
4.2.3	Function deleteVkFFT()	26
4.2.4	Function VkFFTGetVersion()	26
4.3	VkFFT configuration	26
4.3.1	Driver API parameters	38
4.3.2	Memory management parameters	39

4.3.3	General FFT parameters	40
4.3.4	Precision parameters (and some things that can affect it):	41
4.3.5	Advanced parameters (code will work fine without using them)	42
4.3.6	Rader control parameters	44
4.3.7	Bluestein control parameters	44
4.3.8	Zero padding parameters	45
4.3.9	Convolution parameters	45
4.3.10	Register overutilization	46
4.3.11	Extra advanced parameters (filled automatically)	46
5	VkFFT Benchmark/Precision Suite and utils_VkFFT helper routines	48
5.1	utils_VkFFT helper routines	49
6	VkFFT Code Examples	52
6.1	Driver initializations	52
6.2	Simple FFT application example: 1D (one dimensional) C2C (complex to complex) FP32 (single precision) FFT	57
6.3	Advanced FFT application example: ND, C2C/R2C/R2R, different precisions, batched FFT	59
6.4	Advanced FFT application example: out-of-place R2C FFT with custom strides	60
6.5	Advanced FFT application example: 3D FFT with innermost batching	62
6.6	Advanced FFT application example: 3D zero-padded FFT	62
6.7	Convolution application example: 3x3 matrix-vector convolution in 1D	63
6.8	Convolution application example: R2C cross-correlation between two sets of N images	65
6.9	Simple FFT application binary reuse application	66

1 Introduction

This document describes VkFFT - Vulkan/CUDA/HIP/OpenCL/Level Zero/Metal Fast Fourier Transform library. It describes the features and current limitations of VkFFT, explains the API and compares it to other FFT libraries (like FFTW and cuFFT) on the set of examples. It is by no means the final version, so if there is something unclear - feel free to contact me (dtolm96@gmail.com), so I can update it.

2 Using the VkFFT API

This chapter will cover the basics of VkFFT. Fourier transform of a sequence is called Discrete Fourier Transform (DFT). It is defined by the following formula:

$$X_k = \sum_{n=0}^{N-1} x_n e^{-\frac{2\pi i}{N}nk} = \text{DFT}_N(x_n, k), \quad (1)$$

where x_n is the input sequence, N is the length of the input sequence and $k \in [0, N-1]$, $k \in \mathbb{Z}$ is the output index, corresponding to frequency in Fourier space. Corresponding to that, inverse DFT is defined as following:

$$x_n = \sum_{k=0}^{N-1} X_k e^{\frac{2\pi i}{N}nk} = \text{IDFT}_N(X_k, n) \quad (2)$$

VkFFT follows the same definitions as FFTW and cuFFT - forward FFT has the exponent sign -1 , while the inverse has the exponent sign 1 . Note, that inverse transform by default is unnormalized, so to get the input sequence after FFT + iFFT, the user has to divide the result by N .

2.1 Installing VkFFT

VkFFT is distributed as a header-only library (starting from the version 1.3, it is no longer a single header). The installation process consists of the following steps:

1. Copy vkFFT.h file and vkFFT directory into one of the directories included in the user's project.
2. Define VKFFT_BACKEND as a number corresponding to the API used in the user's project: 0 - Vulkan, 1 - CUDA, 2 - HIP, 3 - OpenCL, 4 - Level Zero, 5 - Metal. Definition is done like:

```
-DVKFFT_BACKEND=X
```

in GCC or as

```
set(VKFFT_BACKEND 1 CACHE STRING "0 - Vulkan, 1 - CUDA, 2 -  
↪ HIP, 3 - OpenCL, 4 - Level Zero, 5 - Metal")
```

in CMake.

3. Define VKFFT_MAX_FFT_DIMENSIONS as the number of max dimensions to be supported by VkFFT. Default is 4.
4. Depending on the API backend, the project must use additional libraries for run-time compilation:
 - (a) Vulkan API: SPIRV, glslang and Vulkan. Define VK_API_VERSION to the available Vulkan version. Sample CMakeLists can look like this:

```

find_package(Vulkan REQUIRED)
target_compile_definitions(${PROJECT_NAME} PUBLIC
    ↪ -DVK_API_VERSION=11)#10 - Vulkan 1.0, 11 - Vulkan
    ↪ 1.1, 12 - Vulkan 1.2
target_include_directories(${PROJECT_NAME} PUBLIC
    ↪ ${CMAKE_CURRENT_SOURCE_DIR}/glslang-
    ↪ master/glslang/Include/)
add_subdirectory(${CMAKE_CURRENT_SOURCE_DIR}/glslang-
    ↪ master)

target_include_directories(${PROJECT_NAME} PUBLIC
    ↪ ${CMAKE_CURRENT_SOURCE_DIR}/vkFFT/)
add_library(VkFFT INTERFACE)
target_compile_definitions(VkFFT INTERFACE
    ↪ -DVKFFT_BACKEND=0)

target_link_libraries(${PROJECT_NAME} PUBLIC SPIRV
    ↪ glslang Vulkan::Vulkan VkFFT)

```

(b) CUDA API: CUDA and NVRTC. Sample CMakeLists can look like this:

```

find_package(CUDA 9.0 REQUIRED)
enable_language(CUDA)
set_property(TARGET ${PROJECT_NAME} PROPERTY
    ↪ CUDA_ARCHITECTURES 35 60 70 75 80 86)
target_compile_options(${PROJECT_NAME} PUBLIC
    ↪ "$<${COMPILE_LANGUAGE:CUDA}:SHELL:
    ↪ -DVKFFT_BACKEND=${VKFFT_BACKEND}
    ↪ -gencode arch=compute_60,code=compute_60
    ↪ -gencode arch=compute_70,code=compute_70
    ↪ -gencode arch=compute_75,code=compute_75
    ↪ -gencode arch=compute_80,code=compute_80
    ↪ -gencode arch=compute_86,code=compute_86>")
set_target_properties(${PROJECT_NAME} PROPERTIES
    ↪ CUDA_SEPARABLE_COMPILATION ON)
set_target_properties(${PROJECT_NAME} PROPERTIES
    ↪ CUDA_RESOLVE_DEVICE_SYMBOLS ON)

find_library(CUDA_NVRTC_LIB libnVRTC nVRTC HINTS
    ↪ "${CUDA_TOOLKIT_ROOT_DIR}/lib64"
    ↪ "${LIBNVRTC_LIBRARY_DIR}"
    ↪ "${CUDA_TOOLKIT_ROOT_DIR}/lib/x64" /usr/lib64
    ↪ /usr/local/cuda/lib64)

```

```

add_definitions(-
    ↪ DCUDA_TOOLKIT_ROOT_DIR="${CUDA_TOOLKIT_ROOT_DIR}")

target_include_directories(${PROJECT_NAME} PUBLIC
    ↪ ${CMAKE_CURRENT_SOURCE_DIR}/vkFFT/)
add_library(VkFFT INTERFACE)
target_compile_definitions(VkFFT INTERFACE
    ↪ -DVKFFT_BACKEND=1)

target_link_libraries(${PROJECT_NAME} PUBLIC
    ↪ ${CUDA_LIBRARIES} cuda ${CUDA_NVRTC_LIB} VkFFT)

```

(c) HIP API: HIP and HIPRTC. Sample CMakeLists can look like this:

```

list(APPEND CMAKE_PREFIX_PATH /opt/rocm/hip /opt/rocm)
find_package(hip)

target_include_directories(${PROJECT_NAME} PUBLIC
    ↪ ${CMAKE_CURRENT_SOURCE_DIR}/vkFFT/)
add_library(VkFFT INTERFACE)
target_compile_definitions(VkFFT INTERFACE
    ↪ -DVKFFT_BACKEND=2)
#target_compile_definitions(${PROJECT_NAME} PUBLIC
    ↪ -DVKFFT_OLD_ROCM) #ROCm versions before 4.5 needed
    ↪ kernel include of hiprtc

target_link_libraries(${PROJECT_NAME} PUBLIC hip::host
    ↪ VkFFT)

```

(d) OpenCL API: OpenCL. Sample CMakeLists can look like this:

```

find_package(OpenCL REQUIRED)

target_include_directories(${PROJECT_NAME} PUBLIC
    ↪ ${CMAKE_CURRENT_SOURCE_DIR}/vkFFT/)
add_library(VkFFT INTERFACE)
target_compile_definitions(VkFFT INTERFACE
    ↪ -DVKFFT_BACKEND=3)

target_link_libraries(${PROJECT_NAME} PUBLIC
    ↪ OpenCL::OpenCL VkFFT)

```

(e) Level Zero API: Level Zero; Clang and llvm-spirv must be in the system path (for kernel compilation). Sample CMakeLists can look like this:

```

set(LevelZero_LIBRARY "/usr/lib/x86_64-linux-gnu/")

```

```

set(LevelZero_INCLUDE_DIR "/usr/include/")
find_library(
    LevelZero_LIB
    NAMES "ze_loader"
    PATHS ${LevelZero_LIBRARY}
    PATH_SUFFIXES "lib" "lib64"
    NO_DEFAULT_PATH
)
find_path(
    LevelZero_INCLUDES
    NAMES "ze_api.h"
    PATHS ${LevelZero_INCLUDE_DIR}
    PATH_SUFFIXES "include"
    NO_DEFAULT_PATH
)
target_include_directories(${PROJECT_NAME} PUBLIC
    ↪ ${LevelZero_INCLUDES})
add_library(VkFFT INTERFACE)
target_compile_definitions(VkFFT INTERFACE
    ↪ -DVKFFT_BACKEND=4)

target_link_libraries(${PROJECT_NAME} PUBLIC LevelZero
    ↪ VkFFT)

```

(f) Metal API: Metal. Sample CMakeLists can look like this:

```

add_compile_options(-WMTL_IGNORE_WARNINGS)
find_library(FOUNDATION_LIB Foundation REQUIRED)
find_library(QUARTZ_CORE_LIB QuartzCore REQUIRED)
find_library(METAL_LIB Metal REQUIRED)
target_include_directories(${PROJECT_NAME} PUBLIC
    ↪ "metal-cpp/")

target_include_directories(${PROJECT_NAME} PUBLIC
    ↪ ${CMAKE_CURRENT_SOURCE_DIR}/vkFFT/)
add_library(VkFFT INTERFACE)
target_compile_definitions(VkFFT INTERFACE
    ↪ -DVKFFT_BACKEND=5)

target_link_libraries(${PROJECT_NAME} PUBLIC
    ↪ ${FOUNDATION_LIB} ${QUARTZ_CORE_LIB} ${METAL_LIB}
    ↪ VkFFT)

```


2.2 Fourier Transform Setup

VkFFT follows a plan structure like FFTW/cuFFT with a notable difference - there is a unified interface to all transforms. This means that there are no separate functions like `fftPlan1D/fftPlan2D/fftPlanMany/etc`. The initialization is done through a single configuration struct - `VkFFTConfiguration`. Each parameter of it will be covered in detail in this document. Plans in VkFFT are called `VkFFTApplication` and they are created with a unified `initializeVkFFT` call.

As the code is written in C, don't forget to zero-initialize used structs!

During the `initializeVkFFT(VkFFTApplication* app, VkFFTConfiguration inputLaunchConfiguration)` call `VkFFT` performs kernel generation and compilation from scratch (kernel reuse may be added later). The overall process of initialization looks like this:

1. Get device parameters, perform default initialization of internal copy of configuration struct inside the `VkFFTApplication`, then fill in user-defined parameters from `inputLaunchConfiguration`. `VkFFTApplication` is passed as a pointer, so `initializeVkFFT` modifies the user-provided application.
2. By default, there are two internal FFT plans created - inverse and forward. Multidimensional FFT is done as a combination of 1D FFTs in each axis direction. For each axis, the `VkFFTPlanAxis` function is called.
3. `VkFFTPlanAxis` configures parameters for each axis. It may perform additional memory allocations (see: memory allocated by `VkFFT`).
4. `shaderGenVkFFT` generates corresponding to the axis code in a char buffer (each axis may require more than one kernel: see Four-step FFT, Bluestein's algorithm for FFT).
5. Code is then compiled with the run-time compiler of the specified backend.

Once the plan is no longer need, a call to the `deleteVkFFT` function frees all the allocated resources. There are no processes launched that continue to work outside of the `VkFFT` related function calls.

2.3 Fourier Transform types and their definitions

VkFFT supports commonly used Complex to complex (C2C), real to complex (R2C), complex to real (C2R) transformations and real to real (R2R) Discrete Cosine/Sine Transformations of types I, II, III and IV. VkFFT uses the same definitions as FFTW, except for the multidimensional FFT axis ordering: in FFTW dimensions are ordered with the decrease in consecutive elements stride, while VkFFT does the opposite - the first axis is the non-strided axis (the one that has elements located consecutively in memory with no gaps, usually named as the X-axis). So, in FFTW dimensions are specified as ZYX and in VkFFT as XYZ. This felt more logical to me - no matter if there are 1, 2, 3 or more dimensions, the user can always find the axis with the same stride at the same position. This choice doesn't require any modification in the user's data management - just provide the FFT dimensions in the reverse order to `VkFFT`.

The 1.3 version of VkFFT supports arbitrary number of dimensions, defined as `VKFFT_MAX_FFT_DIM` compile constant. VkFFT also supports two forms of batching: the number of coordinates and the number of systems. The choice of two distinct batching ways is made to support matrix-vector convolutions, where the kernel is presented as a matrix. Overall, the layout of VkFFT can be described as WHD+CN - width, height, depth, other dimensions, coordinate and number of systems (in order of increasing strides, starting with 1 for width). Often, the coordinate part of the layout is not used, so the main batching is done by specifying N.

VkFFT assumes that complex numbers are stored consecutively in memory: RIRIRI... where R denotes the real part of the complex number and I denotes the imaginary part. There is no difference between using a float2/double2/half2 container or access memory as float/double/half/double-double as long as the byte order remains the same.

This section and the next one will cover the basics of VkFFT data layouts and memory management.

2.3.1 C2C transforms

The base FFT algorithm - C2C in VkFFT has the same definition as FFTW. Forward FFT has the exponent sign -1 , while the inverse has the exponent sign 1 . By default, the inverse transform is unnormalized. $N_x N_y N_z$ complex numbers map to $N_x N_y N_z$ complex numbers and no additional padding is required. The resulting data order will be the same as in FFTW/cuFFT, unless special parameters are provided in configuration (see: advanced memory management)

2.3.2 R2C/C2R transforms

R2C/C2R transforms can be explained as C2C transforms with imaginary part set to zero. They exploit Hermitian symmetry of the result: $X_k = X_{N-k}^*$ on the non-strided axis (the one that has elements located consecutively in memory with no gaps). This results in a reduction of required memory to store the complex result - we may only store $\text{floor}(\frac{N_x}{2}) + 1$ complex numbers instead of N_x . However, this results in memory requirements mismatch between input and output in R2C: $\text{floor}(\frac{N_x}{2}) + 1$ complex elements will require $N_x + 2$ real numbers worth of memory for even N_x and $N_x + 1$ real numbers worth of memory for odd N_x . For C2R the situation is reversed. There are two approaches to this problem: pad each sequence of the non-strided axis with zeros to the required length or use out-of-place mode. More information on how to do this will be given in the next section.

2.3.3 R2R (DCT) transforms

R2R transforms in VkFFT are implemented in the form of Discrete cosine transforms and Discrete sine transforms of types I, II, III and IV. Their definitions and transforms results match FFTW:

1. DCT-I: $X_k = x_0 + (-1)^k x_{N-1} + 2 \sum_{n=1}^{N-2} x_n \cos(\frac{\pi}{N-1} nk)$, inverse of DCT-I (itself)

2. DCT-II: $X_k = 2 \sum_{n=0}^{N-1} x_n \cos(\frac{\pi}{N}(n + \frac{1}{2})k)$, inverse of DCT-III
3. DCT-III: $X_k = x_0 + 2 \sum_{n=1}^{N-1} x_n \cos(\frac{\pi}{N}n(k + \frac{1}{2}))$, inverse of DCT-II
4. DCT-IV: $X_k = 2 \sum_{n=0}^{N-1} x_n \cos(\frac{\pi}{N}(n + \frac{1}{2})(k + \frac{1}{2}))$, inverse of DCT-IV (itself)
5. DST-I: $X_k = 2 \sum_{n=0}^{N-1} x_n \sin(\frac{\pi}{N+1}(n+1)(k+1))$, inverse of DST-I (itself)
6. DST-II: $X_k = 2 \sum_{n=0}^{N-1} x_n \sin(\frac{\pi}{N}(n + \frac{1}{2})(k+1))$, inverse of DST-III
7. DST-III: $X_k = (-1)^k x_{N-1} + 2 \sum_{n=0}^{N-2} x_n \sin(\frac{\pi}{N}(n+1)(k + \frac{1}{2}))$, inverse of DST-II
8. DST-IV: $X_k = 2 \sum_{n=0}^{N-1} x_n \sin(\frac{\pi}{N}(n + \frac{1}{2})(k + \frac{1}{2}))$, inverse of DST-IV (itself)

R2R transforms are performed by redefinition of them to the C2C transforms (internal C2C sequence length can be different from the input R2R sequence length). R2R transform performs a one-to-one mapping between real numbers, so they don't require stride management, unlike R2C/C2R.

2.4 Memory management, data layouts for different transforms

2.4.1 VkFFT buffers

VkFFT allows for explicit control over the data flow, which makes both in-place and out-of-place transforms possible. Buffers are passed to VkFFT as VkBuffer pointer in Vulkan, as double void pointers in CUDA/HIP/Level Zero, as cl_mem pointer in OpenCL and as MTL::Buffer pointer in Metal. This is done to maintain a uniform data pattern because some of the buffers can be allocated automatically.

The main buffer is called buffer and it always has to be provided, either during the plan creation or when the plan is executed. All calculations are performed in this buffer and it is always overwritten. To do calculations out-of-place, VkFFT provides an option to specify inputBuffer/outputBuffer buffer. The logic behind their usage is fairly simple - the user specifies inputBuffer if the input data has to be read from a buffer, different from the main buffer. As the data is read only once and nothing is written back to the inputBuffer, this allows doing truly out-of-place transformations. The same logic applies to outputBuffer with the difference that it is responsible for the absolute last write of the VkFFT. It is possible to use all three buffers to create complex data management paths.

It must be noted, that sometimes FFT can not be done inside one buffer (see: Four-Step FFT algorithm, Bluestein's algorithm). To compute FFT in these cases, there exists tempBuffer buffer and data is transferred between the main buffer and tempBuffer during the FFT

execution. The ordering of transfers between the main buffer and tempBuffer is done in such a way, so the initial data read and final data write are obeying the configuration from the previous paragraph. Users can allocate tempBuffer themselves of some memory that does not have any useful information at the time of FFT execution (the tempBuffer size can depend on the configuration, so this is a rather advanced operation - read more in the advanced memory management section) or allow VkFFT to manage tempBuffer allocation itself (tempBuffer will be freed at the deleteVkFFT call).

To compute convolutions and cross-correlations, a kernel buffer has to be specified. It must have the same layout as the result of the FFT transform.

2.4.2 VkFFT buffers strides. A special case of R2C/C2R transforms

To have better control of memory, the user can specify the strides between consecutive elements of different axes. The W (width) stride is fixed to be 1, so to achieve innermost batching, user needs to initialize an N+1 dimension FFT and omit the innermost dimension using omitDimension[0] = 1 parameter. N (number of systems) stride will be consecutive of C in memory if C is used, otherwise N will propagate the previous non-uniform stride multiplied by the corresponding axis length). Strides are specified not in bytes, but in the element type used - similar to the way how the user would access the corresponding element in the array. If all elements are consecutive in C2C, stride for H will be equal to the FFT length of W axis, stride for D will be multiplication of first two FFT axis lengths and so on, stride for C will be multiplication of all FFT axes lengths, etc. These are the default values of C2C and R2R strides if they are not explicitly specified.

One of the main use-cases of strides comes to solve the R2C/C2R Hermitian symmetry H stride mismatch - for real space, it is equal to N_x real elements and for the frequency space it is equal to $\text{floor}(\frac{N_x}{2}) + 1$ complex numbers. So, with strides it is possible to use a buffer, padded to $2 \cdot (\text{floor}(\frac{N_x}{2}) + 1)$ real elements in H stride (all elements between N_x and $2 \cdot (\text{floor}(\frac{N_x}{2}) + 1)$ will not be read so it does not matter what data is there before the write stage). All other strides are done as a multiplication between the previous stride and the number of elements in the previous axis. These are the default values of R2C/C2R strides if they are not explicitly specified.

It is possible to specify separate sets of strides for all user-defined buffers: bufferStride for the main buffer, inputStride for input buffer, outputStride for output buffer (kernel stride is assumed to be the same as bufferStride, tempBuffer strides are configured automatically).

For an out-of-place R2C FFT, there is no need to pad buffer with real numbers, but user must specify H stride there (as it differs to default one) - N_x real elements for real space and $\text{floor}(\frac{N_x}{2}) + 1$ complex numbers for the frequency space.

An out-of-place C2R FFT is a more tricky transform. In the multidimensional case, the main buffer will be written to and read from multiple times. The intermediate stores have a complex layout, which requires more space than the output real layout, so in order not to modify the input data, there exist two options. First, pad the real data layout has to $2 \cdot (\text{floor}(\frac{N_x}{2}) + 1)$ real elements in H stride (complex buffer will be used as inputBuffer, real buffer as buffer). Second, use the third buffer, so both input and output buffers have their

original layouts (complex buffer will be used as inputBuffer, the main buffer for calculations is buffer and output real buffer as outputBuffer). If you use inverseReturnToInputBuffer option, where R2C is configured to read from input buffer and C2R is configured to write to the input buffer; C2R will modify the buffer it reads from in some cases (see issue #58)

By default, R2C/C2R will assume the complex-padded real strides in the in-place mode and non-padded real strides in out-of-place mode.

2.5 VkFFT algorithms

VkFFT implements a wide range of algorithms to compute different types of FFTs but all of them can be reduced to a mixed-radix Cooley-Tukey FFT algorithm in the Stockham autosort form. The main idea behind it is to decompose the sequence as a set of primes, for each of which FFT can be written down exactly. As of now, VkFFT has radix implementations for primes up to 13, so all C2C sequences decomposable into a multiplication of such primes will be done purely with the Stockham algorithm. Below additional algorithms and their use-cases are described.

2.5.1 Bluestein's algorithm

A complex algorithm that is used in cases where the sequence is not decomposable with implemented radix butterflies (currently - primes up to 13). It is derived by replacing $nk = (n^2 + k^2 - (n - k)^2) / 2$ in 1:

$$X_k = \left(e^{-\pi i \frac{k^2}{N}} \right) \sum_{n=0}^{N-1} \left(x_n e^{-\pi i \frac{n^2}{N}} \right) \left(e^{\pi i \frac{(k-n)^2}{N}} \right) = b_k^* \sum_{n=0}^{N-1} a_n b_{k-n} \quad (3)$$

$$a_n = x_n b_n^* \quad (4)$$

$$b_n = e^{\pi i \frac{n^2}{N}} \quad (5)$$

Here FFT is represented as a convolution between two sequences: a_n and b_n , which can be performed by the means of convolution theorem:

$$F\{a * b\} = F\{a\} \cdot F\{b\} \quad (6)$$

By padding a_n and b_n to a sequence length decomposable with implemented radix butterflies with a size of at least $2N - 1$ (because the length of b_n is $2N - 1$), we can perform FFT of any length. FFT of b_n can be precomputed, so overall this algorithm requires at least 4x the computations and more memory transfers. This algorithm can be combined with all other algorithms implemented in VkFFT. If an FFT can not be done in a single upload, a tempBuffer has to be allocated (because the logical FFT buffer size is bigger than the original system).

2.5.2 The Four-Step FFT algorithm

GPUs and CPUs have a hierarchical memory model - the closer memory to the unit that performs the computations, the faster its speed and the lower the size. So it is advantageous to split FFTs, not to the lowest primes, but to some bigger multiplication of those primes, then upload this subsequence to the closest cache level to the cores and do the final prime split there. The absolute lowest level is the register file, however, it does not allow for thread communications outside the warp. For this purpose, modern GPUs employ shared memory - a fast memory with a bank structure that is visible to all threads in a thread block. The usual sizes of it change on a scale from 16KB to 192KB and it is often beneficial to use it fully. However, if the full sequence can not fit inside the shared memory, FFT has to be done in multiple uploads - with the Four Step FFT algorithm. The main idea behind it is to represent a big 1D sequence as a 2D (or 3D for the three-upload scheme) FFT - we first do FFT along the columns, then the rows, then transpose the result and multiply by a special set of phase vectors. Similar decomposition idea as the main Cooley-Tukey algorithm. However, performing transpositions in-place is a complicated task - especially for a non-trivial ratio between dimensions. It will also require an additional read/write stage, as it can not be merged with the last write of the FFT algorithm. The easiest and the most performant solution is to use a tempBuffer (it is the main reason for having this functionality, actually) and store intermediate FFT results out-of-place. This way the last transposition step can be merged with the write step, as we can overwrite the output buffer without losing data.

To estimate if your sequence size is single upload or not, divide the amount of available shared memory (48KB - Nvidia GPUs with Vulkan/OpenCL API, 64KB - AMD GPUs, 100KB - Nvidia GPUs in CUDA API) by the complex size used for calculations (4 byte - half precision, 8 byte - single precision, 16 byte - double precision, 32 byte - double-double). For 64KB of shared memory, we get 8192 as max single upload single-precision non-strided FFT, 4096 for double precision. For strided axes (H and D parts of the layout) these numbers have to be divided by 4 and 2 respectively to achieve coalescing, resulting in 2048 length for single upload in both precisions. For more information on coalescing see: coalescing API reference.

In the case of the Four-Step FFT algorithm, tempBuffer size has to be at least the same as the default main buffer size. It does not matter how many uploads are in the Four Step FFT algorithm - only a single tempBuffer is required. In this document, all systems that can fit in the shared memory entirely and be done without the Four Step FFT algorithm (and multiple uploads) are called single upload systems.

If the last transposition is not required (the output data is allowed to be in not unshuffled form) it can be disabled during the configuration phase. This way tempBuffer will not be needed and all computations will be done in-place (unless Bluestein's algorithm is used). An example use-case of this is convolutions - if the kernel is computed with the same operation ordering, point-wise multiplication in the frequency domain is not dependent on the correct data ordering and the inverse FFT will restore the original layout.

2.5.3 R2C/C2R FFTs

A typical approach to a single upload R2C/C2R system is to just set the imaginary part to zero inside the shared memory and do a simple C2C transform. This doesn't affect the amount of memory transferred from VRAM and is not a bad approach as FFT is a memory-bound algorithm, however, this can be improved in multidimensional (in HDCN part of the layout) case by the composition of a single C2C sequence from two real sequences and some write for R2C/read for C2R post-processing. Both of these algorithms are implemented in VkFFT. Note, that R2C/C2R only affects the non-strided axis (W). All strided axes are still done as C2C.

2.5.4 R2C/C2R multi-upload FFT algorithm

For even sequences there exists an easy mapping between R2C/C2R FFTs and the C2C of half the size. In this case, all even indices (starting from 0) are read as the real values of a complex number and all odd indices are read as the imaginary values. This C2C sequence can be done with the help of the Four-Step FFT algorithm. When FFT is done, separate post-processing for R2C/pre-processing for C2R is applied.

2.5.5 R2R Discrete Cosine/Sine Transforms

There exist many different mappings between DCT and FFT. DSTs are reformulated as DCTs inside the VkFFT, so they use the same algorithms. As of now, VkFFT has the following algorithms implemented (all single-upload for now):

- DCT-I - mapping between R2R and C2C of the $2N - 2$ length. For non-strided axis can use an optimization similar to the R2C/C2R multidimensional case (setting the imaginary part to the next FFT sequence).
- DCT-II/DCT-III - mapping between R2R and C2C of the same length. For non-strided axis can use an optimization similar to the R2C/C2R multidimensional case (setting the imaginary part to the next FFT sequence).
- DCT-IV - for even sizes, mapping between R2R and C2C sequence of half-length. For odd sizes mapping to the FFT of the same length (for non-strided axis can use an optimization similar to the R2C/C2R multidimensional case (setting the imaginary part to the next FFT sequence)).

2.5.6 Register overutilization

Not an FFT algorithm by itself, but an optimization to do bigger sequences in a single upload instead of switching to the Four Step FFT algorithm. The main idea behind it is to use a register file (which is often bigger than the amount of shared memory) to store the sequence and use shared memory only as a communication buffer. This is useful in Vulkan and OpenCL APIs on Nvidia GPU, as they are only allowed to allocate 48KB of shared memory with a register file having the size of 256KB.

2.5.7 Zero padding

Not an FFT algorithm by itself, but a memory management optimization. If the user's system has parts that are known to be zero - for example, when an open system is modeled, to avoid a circular part of the FFT system has to be padded with zeros up to 2x in each direction. VkFFT can omit sequences full of zeros and don't perform the corresponding memory transfers and computations, as the output result will be zero. This way it is possible to get up to two times speed increase in the 2D case and up to 3x increase in the 3D case.

2.5.8 Convolution and cross-correlation support

With the help of the Convolution theorem, which states that the Fourier transform of a convolution is the pointwise product of signals Fourier transforms, it is possible to perform convolution with $N \log N$ complexity, compared to N^2 complexity of the simple multiplication approach. This is extremely useful for kernels spanning more than 50 elements in size. VkFFT can merge the last step FFT, kernel multiplication in the Fourier domain and the first step of inverse FFT to provide substantial memory transfer savings. Moreover, FFTs of big sequences can be performed without data reordering, which results in a better locality.

2.6 VkFFT accuracy

To measure how VkFFT (single/double/half/double-double precision) results compare to cuFFT/rocFFT (single/double/half precision) and FFTW (double/quad precision), multiple sets of systems covering full supported C2C/R2C+C2R/R2R FFT range are filled with random complex data on the scale of $[-1,1]$ and one transform was performed on each system. Samples 11(single), 12(double), 13(half), 14(non-power of 2 C2C, single), 15(R2C+C2R, single), 16(DCT-I/II/III/IV, single), 17(DCT-I/II/III/IV, double), 18(non-power of 2 C2C, double), 19(double-double emulation of quad precision, C2C) are available in VkFFT Benchmark Suite to perform VkFFT verification on any of the target platforms. Overall, the Cooley-Tukey algorithm (Stockham autosort) exhibits logarithmic relative error scaling, similar to those of other GPU FFT libraries. Typically, the more computationally expensive algorithm is - the worse its precision is. So, Bluestein's algorithm has lower accuracy than Stockham autosort algorithm.

Single precision in VkFFT supports two modes of calculation - by using the on-chip Special Function Units that can compute sines and cosines on the go or by using the precomputed on CPU look-up tables. For Nvidia and AMD GPUs, SFU provide great precision, while Intel iGPUs and mobile GPUs must use LUT to perform FFTs correctly.

Double precision in VkFFT also supports two modes of calculation - by using polynomial sincos approximation and computing them on-chip or by using precomputed LUT as well. The second option is the better one, as polynomial sincos approximation is too compute-heavy for modern GPUs. It is selected by default on all devices.

Half precision is supported in all backends. Half precision is implemented only as a memory trick - all on-chip computations are done in single precision. Half precision can use SFU or LUT as well.

Double-double emulation of quad precision uses precomputed on CPU in quad precision twiddle factors and implements only additions and multiplications on GPU (for now).

VkFFT also supports mixed-precision operations, where memory storing is done at lower precision, compared to the on-chip calculations. For example, it is possible to read data in single precision, do calculations in double and store data back in single precision.

2.7 VkFFT additional memory allocations

In this section, all GPU memory allocations that are done by VkFFT are described. There are up to three situations when VkFFT allocates memory. All of the VkFFT allocated memory is freed at the deleteVkFFT call.

2.7.1 LUT allocations

This memory is used to store precomputed twiddle factors and phase vectors used during the computation. This buffer can have:

- twiddle factors for each radix stage of Stockham FFT calculation
- phase vectors used in the Four Step FFT algorithm between stages
- phase vectors used in DCT/DST-II/III/IV to perform a mapping between R2R and C2C
- phase vectors used in post-processing for R2C/pre-processing for C2R for even length sequences as C2C of half size

VkFFT manages LUT allocations by itself and they are performed during the initializeVkFFT call. LUT are allocated per axis, though some of them can be reused if the axes have the same LUT. Inverse and forward FFT plans share the same LUT (conjugation is performed on-chip).

2.7.2 The Four-Step FFT algorithm - tempBuffer allocation

To perform the merging of the transposition with the last upload of an axis, VkFFT requires additional memory to mimic an out-of-place execution. This memory is located in tempBuffer and has to be of at least the same size as the main buffer. It is possible for the users to allocate it themselves, though if this is not done, VkFFT can do the allocation automatically (the size of the tempBuffer will be the same as the main buffer, unless the logical dimensions of FFT are bigger than user-defined - then, it will allocate the system with the minimal size, that can cover maximal logical system size used in any of the axes - see next subsection).

2.7.3 Bluestein's buffers allocation

To do Bluestein's FFT algorithm, precomputed sequences $b_n = e^{\pi i \frac{n^2}{N}}$, $FFT(b_n)$ and $iFFT(b_n)$ are required. For each axis, they can be different and are computed separately (unless VkFFT can determine that they match, then the buffers are allocated only once). Notably, as Bluestein's algorithm pads the sequence length to at least $2N - 1$, if it can not be done in a single upload and the Four Step algorithm has to be used, the intermediate storage

required will be bigger than the main buffer size. In this case, tempBuffer must always be allocated. As the padded sequence can be different for each of the dimensions, the required size of the tempBuffer will also vary. VkFFT determines the biggest size needed among axes and allocated tempBuffer of this size.

2.8 VkFFT support for double-double emulation of quad precision

Since VkFFT 1.3.2 experimental support for double-double emulation of quad precision has been added. Double-double number is defined as an unevaluated sum of two double numbers, second one being on the order of ~ 1 ULP of the first one (like 1 and $1e-16$). This boosts the significand from 53 to 106 bits and leaves the exponent the same (11bits). For example, true quad precision would have 15 bits of exponent and 113 bits significand. The range of double-double stays the same as in double precision ($1e308$ vs $1e4932$), while precision boosts from $\sim 1e-16$ to $\sim 1e-32$ (vs $\sim 1e-34$ of true quad precision). The double-double format of VkFFT is equivalent to `__ibm128` type in gcc. VkFFT uses quadmath library to precompute twiddle factors and initialize kernels for double-double in the true quad precision. Currently, conversion between quad and double-double is only possible on CPU in true quad (GPU can only do operations on data in double-double format). This will be changed later.

To enable double-double support, link against quadmath and set `VkFFT_use_FP128_double_double` compile flag to true.

The performance impact is rather harsh - the code becomes compute-bound even on modern HPC clusters. The simple radix algorithms work well and are usable ($\sim 2x$ impact compared to double precision), while Bluestein and big primes Rader algorithms incur a high toll on operations count. This can be improved in the future, but probably not to the extent of being VRAM bandwidth-bound on modern accelerators.

The main usage for the double-double precision should be mixed-precision calculations, where most of the time is spent in regular FP64 and only some parts that hit precision boundaries are calculated in double-double. For example, in VkFFT it is possible to perform a double precision FFT with all on-chip calculations done in double-double. This way the final errors should be on the scale of 1 ULP.

3 Runtime code optimization platform

3.1 Platform for GPU code generation

While GPUs show extraordinary performance, they are harder to program for compared to regular CPUs. Firstly, the algorithms must be optimized for SIMT (single instruction, multiple threads) execution, which is the soul of the GPU model. Secondly, rapid GPU development makes all vendor architectures (Nvidia, AMD, Intel) largely different. Below we will list some examples:

- Shared memory size. If your code uses the full 64KB of shared memory on AMD MI250X, it won't be optimal on Nvidia H100 which has 256KB of it. The same goes for L2/L3 cache size.
- Warp size. If your code uses 32 core warps (usual on Nvidia), it will not work or work inefficiently for 64 core warps of AMD CDNA.
- Coalesced memory. For Nvidia and AMD we have to coalesce 32-byte memory requests, for Intel, it is equal to 64 bytes. Multiple of these requests are then combined into 128-byte transactions to the chip. However, this is a big simplification. There exist undocumented issues for Nvidia and AMD, when bandwidth drops when code tries to do distant, but coalesced memory accesses, i.e. data are still grouped in 32-byte transactions but target addresses with $> 2^{18}$ bytes between them (and the bigger the distance - the bigger the drop). The load-store unit (LSU) can not group 32-byte transactions. It is possible to regain full bandwidth by switching to 128-byte memory coalescing on Nvidia, while on AMD memory management follows different rules and these requests just happen to hit the same physical memory pin and thus be serialized. All this information is usually never shared by vendors, so programmers have to find solutions by trial and error.
- Compute unit composition. Some compute units have special function units that can evaluate single precision sines and cosines in one cycle (AMD and Nvidia). On Intel GPUs, they are often imprecise. The ratio between FP32 and FP64 units also imposes limitations on the algorithm - compute units with low FP64 core count will benefit from data precomputation in the lookup tables more. Some compute units have tensor cores for accelerated matrix products.
- There are many other nuanced differences (like the number of supported by hardware queues/streams, support for CUDA MPS execution model, register spilling control) which result in some complicated kernels being unpredictable in performance before you try. So having an option to tweak some of its parameters loosely related to hardware is a huge benefit.

This brings us to a big obstacle: all vendors have their own computing platforms. CUDA for Nvidia, HIP for AMD, Level Zero for Intel, Metal for Apple. To facilitate the development of code for all these different in details but similar in SIMD nature GPUs, I used all the knowledge obtained via numerous experiments with VkFFT to generalize GPU programming as much as possible - and create a unified platform for GPU code generation.

3.2 Platform structure

The project structure is split into four main parts:

- Application manager. We use a container-based model of management - the user fills out the configuration struct and then initializes the application with only one initializer function. So this part of the platform performs initial configuration management, pass of the parameters to the corresponding to the particular subprogram plan initializer, binary load/store for future reuse and resources deallocation. Here we have application dispatch handlers, which abstract the actual application-dependent code dispatch handlers in the same way as we have one function for initialization.
- Plan manager. After configuration propagation, the plan manager performs extensive configuration of parameters for a particular task code generation. It also performs general API-dependent parameter initialization for supported backends. The plan manager also performs code compilation using supported by the target API runtime compiler (NVRTC by Nvidia CUDA, HIPRTC by AMD HIP, glslang by Vulkan, native OpenCL, clang and llvm-spirv for Intel Level Zero). It manages additional memory allocations and backend-specific memory management. Finally, it has plan-specific dispatch routines and a plan resource deallocation manager.
- Code manager. This part of the platform is solely responsible for GPU code generation. It takes a special input configuration structure created from user configuration and API-derived parameters and produces a char array containing code to be executed. The code manager has a block structure that developers can use. This allows reusing parts of code between kernels. The block structure has the following designations:
 - Level 2 kernels: problem-specific kernels and layout configuration. These functions are called by the plan manager and they are intended to have a clear description of what is going on in the function. There is no pure code here - only calls to append lower-level functions code to the kernel code and additional inlining constants configuration that can be performed by the propagated from the plan manager information. This kernel also performs calls to additional kernel configuration routines like shared memory and register inlining inside a kernel and overall kernel decoration (like calls to inlining structure definitions, calls to kernel name inlining) - these routines are defined in lower kernels levels.
 - Level 1 kernels: less abstract simple routine kernels, like matrix-vector products, copy kernels, radix components of FFT, read/write managers. All of them must have a specified structure so they can be reused, for example, considering all input data to be in shared memory with a specific layout - then level 2 kernels will have a simple job of providing data to level 1 kernels in a specified manner and get the output in an also specified manner.
 - Level 0 kernels: memory management kernels. Here we have the absolute simplest kernel blocks, like memory data transfer request calls, synchronization, inlining of register/shared memory initialization code, etc. Many APIs have their own definitions for particular math operations, so we can define functions for inlining additions multipli-

cations and other basic arithmetics here as well. The target of this level is to provide level 1 and 2 kernels with an easy way to reuse things that can be generalized between the widest range of math algorithms.

All of the different level kernels have full access to the target GPU hardware parameters - hence (depending on the quality of implemented kernels) it is possible to autotune the algorithm at run-time and simplify the compiler job as much as possible. This involves loop unrolling, constants inlining, etc. Having an additional level (level 1) dedicated to the simple algorithms allows for kernel merging.

- Structs definitions. Here we define all structs used by the platform. It also contains some helper structs that can be used to manage GPU with different APIs and collect generated applications in a library struct for easier access through a map key (if the user's code uses many applications, this turns out to be extremely helpful).

3.3 Container abstractions

The main design improvement that enables the platform functionality and greatly expands its usability are container abstractions. The code generator now operates on special data containers (implemented as C-unions), that can hold integer/float/complex numbers either as known during the plan creation values or as strings of variable names - no more `sprintf` calls with big chunks of unreadable code.

An example operation that uses containers is a MUL operation that performs a multiplication $A = B \times C$. If all containers have known values, A can be precomputed during plan creation. If A , B and C are, for example, register names, we print to the kernel an operation of multiplication. Each operation like MUL has its own representation for all supported APIs in the `vkFFT_MathUtils.h` file. An important distinction is that it is not a full compiler, but a simple tool that unifies syntax of different APIs and supports JIT optimizations. However, as the code manager outputs code that resembles binary, one of the future exploration possibilities it may be interesting to try to make level 0 and 1 kernels output assembly without additional compilation by the plan manager.

Usage of container abstractions brings exciting opportunities to code creation. Behaving as multi-API templates it will be possible to implement double-double FP128 mathematics support on GPU without modifying any of the already written algorithm code. Another big benefit is automatic type casting between precisions and option to perform precomputation in higher precision than the GPU code executes in.

I have also explored other ways to implement these container abstractions - including switching to C++ from C99 and using templates and operator overloading, however this turned out to be just a syntactic sugar, as if it did not reduce number of lines of code and did not increase readability. So the code is still C with almost instant compilations as a bonus. All this is subject to change in the future if people present convincing arguments.

VkFFT supports two base number modes (used for all FP calculations on CPU) - long double or quad (`__float128`). The latter is needed to calculate the twiddle factors in double-double emulation of quad precision.

4 VkFFT API Reference

This section covers error codes, API functions that can be used by the user and configuration parameters.

4.1 Return value VkFFTResult

All VkFFT Library return values except for VKFFT_SUCCESS are used in case of a failure and provide information on what has gone wrong. VkFFTResult is unified among different backends, though some of its values may not be used in specific backends. Possible return values of VkFFTResult are defined as following:

```
typedef enum VkFFTResult {
VKFFT_SUCCESS = 0, // The VkFFT operation was successful
VKFFT_ERROR_MALLOC_FAILED = 1, // Some malloc call inside
↳ VkFFT has failed. Report this to the GitHub repo
VKFFT_ERROR_INSUFFICIENT_CODE_BUFFER = 2, // Generated
↳ kernel is bigger than default kernel array. Increase it
↳ with maxCodeLength parameter of configuration.
VKFFT_ERROR_INSUFFICIENT_TEMP_BUFFER = 3, // Temporary
↳ string used in kernel generation is bigger than default
↳ temporary string array. Increase it with maxTempLength
↳ parameter of configuration.
VKFFT_ERROR_PLAN_NOT_INITIALIZED = 4, // Code attempts to use
↳ uninitialized plan (it is zero inside
↳ VkFFTApplication)
VKFFT_ERROR_NULL_TEMP_PASSED = 5, // Internal kernel
↳ generation error
VKFFT_ERROR_MATH_FAILED = 6, // Math instruction in code
↳ generator failed
VKFFT_ERROR_FFTdim_GT_MAX_FFT_DIMENSIONS = 7, // User
↳ specified a number of dimensions higher than the code is
↳ compiled to handle with VKFFT_MAX_FFT_DIMENSIONS
VKFFT_ERROR_NONZERO_APP_INITIALIZATION = 8, // The app
↳ pointer is not pointing to the zero-filled memory block
VKFFT_ERROR_INVALID_PHYSICAL_DEVICE = 1001, // No physical
↳ device is provided (Vulkan API)
VKFFT_ERROR_INVALID_DEVICE = 1002, // No device is provided
↳ (All APIs)
VKFFT_ERROR_INVALID_QUEUE = 1003, // No queue is provided
↳ (Vulkan API)
VKFFT_ERROR_INVALID_COMMAND_POOL = 1004, // No command pool
↳ is provided (Vulkan API)
VKFFT_ERROR_INVALID_FENCE = 1005, // No fence is provided
↳ (Vulkan API)
```

```

VKFFT_ERROR_ONLY_FORWARD_FFT_INITIALIZED = 1006,    // VkFFT
↳ tries to access inverse FFT plan, when appliction is
↳ created with makeForwardPlanOnly flag
VKFFT_ERROR_ONLY_INVERSE_FFT_INITIALIZED = 1007,    // VkFFT
↳ tries to access forward FFT plan, when appliction is
↳ created with makeInversePlanOnly flag
VKFFT_ERROR_INVALID_CONTEXT = 1008,                // No context is
↳ provided (OpenCL API)
VKFFT_ERROR_INVALID_PLATFORM = 1009,                // No platform is
↳ provided (OpenCL API)
VKFFT_ERROR_EMPTY_FILE = 1011,
VKFFT_ERROR_EMPTY_FFTdim = 2001,    // Number of dimensions is
↳ not provided in the configuration
VKFFT_ERROR_EMPTY_size = 2002,    // Array of dimensions is not
↳ provided in the configuration
VKFFT_ERROR_EMPTY_bufferSize = 2003,    // Buffer size has to
↳ be provided during the application creation
VKFFT_ERROR_EMPTY_buffer = 2004,    // Buffer has te be
↳ specified either at the application creation stage or
↳ during launch through VkFFTLaunchParams struct
VKFFT_ERROR_EMPTY_tempBufferSize = 2005,    // Same error as
↳ VKFFT_ERROR_EMPTY_bufferSize if userTempBuffer is enabled
VKFFT_ERROR_EMPTY_tempBuffer = 2006,    // Same error as
↳ VKFFT_ERROR_EMPTY_buffer if userTempBuffer is enabled
VKFFT_ERROR_EMPTY_inputBufferSize = 2007,    // Same error as
↳ VKFFT_ERROR_EMPTY_bufferSize if isInputFormatted is enabled
VKFFT_ERROR_EMPTY_inputBuffer = 2008,    // Same error as
↳ VKFFT_ERROR_EMPTY_buffer if isInputFormatted is enabled
VKFFT_ERROR_EMPTY_outputBufferSize = 2009,    // Same error as
↳ VKFFT_ERROR_EMPTY_bufferSize if isOutputFormatted is
↳ enabled
VKFFT_ERROR_EMPTY_outputBuffer = 2010,    // Same error as
↳ VKFFT_ERROR_EMPTY_buffer if isOutputFormatted is enabled
VKFFT_ERROR_EMPTY_kernelSize = 2011,    // Same error as
↳ VKFFT_ERROR_EMPTY_bufferSize if performConvolution is
↳ enabled
VKFFT_ERROR_EMPTY_kernel = 2012,    // Same error as
↳ VKFFT_ERROR_EMPTY_buffer if performConvolution is enabled
VKFFT_ERROR_EMPTY_applicationString = 2013,    //
↳ loadApplicationString is zero when
↳ loadApplicationFromString is enabled
VKFFT_ERROR_EMPTYRY_useCustomBluesteinPaddingPattern_arrays =
↳ 2014,    // pointers to primeSizes or paddedSizes arrays are
↳ zero when useCustomBluesteinPaddingPattern is enabled

```

```

VKFFT_ERROR_EMPTY_app = 2015,    // app pointer is zero
VKFFT_ERROR_UNSUPPORTED_RADIX = 3001,    // VkFFT has
↳ encountered unsupported radix (more than 13) during
↳ decomposition and Bluestein's FFT fallback did not work
VKFFT_ERROR_UNSUPPORTED_FFT_LENGTH = 3002,    // VkFFT can not do
↳ this sequence length currently - it requires mor than
↳ three-upload Four step FFT
VKFFT_ERROR_UNSUPPORTED_FFT_LENGTH_R2C = 3003,    // VkFFT can
↳ not do this sequence length currently - odd multi-upload
↳ R2C/C2R FFTs
VKFFT_ERROR_UNSUPPORTED_FFT_LENGTH_R2R = 3004,    // VkFFT can
↳ not do this sequence length currently - multi-upload R2R
↳ transforms
VKFFT_ERROR_UNSUPPORTED_FFT_OMIT = 3005,    // VkFFT can not
↳ omit sequences in convolution calculations and R2C/C2R case
VKFFT_ERROR_FAILED_TO_ALLOCATE = 4001,    // VkFFT failed to
↳ allocate GPU memory
VKFFT_ERROR_FAILED_TO_MAP_MEMORY = 4002,    // 4002-4052 are
↳ handlers for errors of used backend APIs. They may indicate
↳ a driver failure. If they are thrown - report to the GitHub
↳ repo
VKFFT_ERROR_FAILED_TO_ALLOCATE_COMMAND_BUFFERS = 4003,
VKFFT_ERROR_FAILED_TO_BEGIN_COMMAND_BUFFER = 4004,
VKFFT_ERROR_FAILED_TO_END_COMMAND_BUFFER = 4005,
VKFFT_ERROR_FAILED_TO_SUBMIT_QUEUE = 4006,
VKFFT_ERROR_FAILED_TO_WAIT_FOR_FENCES = 4007,
VKFFT_ERROR_FAILED_TO_RESET_FENCES = 4008,
VKFFT_ERROR_FAILED_TO_CREATE_DESCRIPTOR_POOL = 4009,
VKFFT_ERROR_FAILED_TO_CREATE_DESCRIPTOR_SET_LAYOUT = 4010,
VKFFT_ERROR_FAILED_TO_ALLOCATE_DESCRIPTOR_SETS = 4011,
VKFFT_ERROR_FAILED_TO_CREATE_PIPELINE_LAYOUT = 4012,
VKFFT_ERROR_FAILED_SHADER_PREPROCESS = 4013,
VKFFT_ERROR_FAILED_SHADER_PARSE = 4014,
VKFFT_ERROR_FAILED_SHADER_LINK = 4015,
VKFFT_ERROR_FAILED_SPIRV_GENERATE = 4016,
VKFFT_ERROR_FAILED_TO_CREATE_SHADER_MODULE = 4017,
VKFFT_ERROR_FAILED_TO_CREATE_INSTANCE = 4018,
VKFFT_ERROR_FAILED_TO_SETUP_DEBUG_MESSENGER = 4019,
VKFFT_ERROR_FAILED_TO_FIND_PHYSICAL_DEVICE = 4020,
VKFFT_ERROR_FAILED_TO_CREATE_DEVICE = 4021,
VKFFT_ERROR_FAILED_TO_CREATE_FENCE = 4022,
VKFFT_ERROR_FAILED_TO_CREATE_COMMAND_POOL = 4023,
VKFFT_ERROR_FAILED_TO_CREATE_BUFFER = 4024,
VKFFT_ERROR_FAILED_TO_ALLOCATE_MEMORY = 4025,

```



```

VKFFT_ERROR_FAILED_TO_BIND_BUFFER_MEMORY = 4026,
VKFFT_ERROR_FAILED_TO_FIND_MEMORY = 4027,
VKFFT_ERROR_FAILED_TO_SYNCHRONIZE = 4028,
VKFFT_ERROR_FAILED_TO_COPY = 4029,
VKFFT_ERROR_FAILED_TO_CREATE_PROGRAM = 4030,
VKFFT_ERROR_FAILED_TO_COMPILE_PROGRAM = 4031,
VKFFT_ERROR_FAILED_TO_GET_CODE_SIZE = 4032,
VKFFT_ERROR_FAILED_TO_GET_CODE = 4033,
VKFFT_ERROR_FAILED_TO_DESTROY_PROGRAM = 4034,
VKFFT_ERROR_FAILED_TO_LOAD_MODULE = 4035,
VKFFT_ERROR_FAILED_TO_GET_FUNCTION = 4036,
VKFFT_ERROR_FAILED_TO_SET_DYNAMIC_SHARED_MEMORY = 4037,
VKFFT_ERROR_FAILED_TO_MODULE_GET_GLOBAL = 4038,
VKFFT_ERROR_FAILED_TO_LAUNCH_KERNEL = 4039,
VKFFT_ERROR_FAILED_TO_EVENT_RECORD = 4040,
VKFFT_ERROR_FAILED_TO_ADD_NAME_EXPRESSION = 4041,
VKFFT_ERROR_FAILED_TO_INITIALIZE = 4042,
VKFFT_ERROR_FAILED_TO_SET_DEVICE_ID = 4043,
VKFFT_ERROR_FAILED_TO_GET_DEVICE = 4044,
VKFFT_ERROR_FAILED_TO_CREATE_CONTEXT = 4045,
VKFFT_ERROR_FAILED_TO_CREATE_PIPELINE = 4046,
VKFFT_ERROR_FAILED_TO_SET_KERNEL_ARG = 4047,
VKFFT_ERROR_FAILED_TO_CREATE_COMMAND_QUEUE = 4048,
VKFFT_ERROR_FAILED_TO_RELEASE_COMMAND_QUEUE = 4049,
VKFFT_ERROR_FAILED_TO_ENUMERATE_DEVICES = 4050,
VKFFT_ERROR_FAILED_TO_GET_ATTRIBUTE = 4051,
VKFFT_ERROR_FAILED_TO_CREATE_EVENT = 4052,
VKFFT_ERROR_FAILED_TO_CREATE_COMMAND_LIST = 4053,
VKFFT_ERROR_FAILED_TO_DESTROY_COMMAND_LIST = 4054,
VKFFT_ERROR_FAILED_TO_SUBMIT_BARRIER = 4055
} VkFFTResult;

```

4.2 VkFFT application management functions

VkFFT has a unified plan management model - all different transform types/ dimensionalities/ precision use the same calls with configuration done through `VkFFTConfiguration` struct. This section shows how to initialize/use/free `VkFFT` with this unified model, while the next one will go into how to configure `VkFFTConfiguration` correctly. All of the functions operate on `VkFFTApplication` and `VkFFTConfiguration` assuming they have been zero-initialized before usage, so do not forget to do this when initializing (otherwise `VkFFT` will throw an error):

```

VkFFTConfiguration configuration = {};
VkFFTApplication app = {};

```

4.2.1 Function initializeVkFFT()

```
VkFFTResult initializeVkFFT(VkFFTApplication* app,  
    ↪  VkFFTConfiguration inputLaunchConfiguration)
```

Creates an FFT application (collection of forward and inverse plans). As forward and inverse FFTs may have different memory layouts, can have different normalizations - they are done as separate internal plans inside `VkFFTApplication`. This call assumes the application to be zero-initialized, so can be only done once on a particular application, until it is deleted.

If the `initializeVkFFT` call fails, it frees all allocated by `VkFFT` CPU/GPU resources and sets the application to zero. `VkFFTResult` is returned with an error code corresponding to what went wrong.

In case of success, `VkFFTApplication` will contain initialized plans with compiled kernels ready for execution with `VKFFT_SUCCESS` returned.

4.2.2 Function VkFFTAppend()

```
VkFFTResult VkFFTAppend(VkFFTApplication* app, int inverse,  
    ↪  VkFFTLaunchParams* launchParams)
```

Performs FFT in the `int` inverse direction (-1 for forward FFT, 1 for inverse FFT). FFT plans are selected from the `VkFFTApplication` collection automatically. `VkFFTApplication` must be initialized with `initializeVkFFT` call before. `VkFFTLaunchParams` struct allows for pre-launch configuration of some parameters, namely:

- `buffer` - similar to how FFTW/cuFFT expects input/output data pointers in `*execC2C` (and other) function calls, `VkFFT` allows specifying memory used for computations at launch. It must have the same size/layout/strides as defined during the application creation.
- `inputBuffer/outputBuffer/tempBuffer/kernel` - other buffers can also be specified at launch. In addition to them having the same size/layout/strides as defined during the application creation, the application must be created with flags enabling the corresponding buffer usage: `isInputFormatted/isOutputFormatted/userTempBuffer/performConvolution` respectively.
- `bufferOffset/tempBufferOffset/inputBufferOffset/outputBufferOffset/kernelOffset` - specify if `VkFFT` has to offset the first element position inside the corresponding buffer. In bytes. Default 0. `specifyOffsetsAtLaunch` parameter must be enabled during the `initializeVkFFT` call before.

Depending on the API, the execution model may vary and require additional information at launch:

- Vulkan API: `VkFFT` appends a sequence of `vkCmdDispatch` calls to the user-defined `VkCommandBuffer` (with respective push constants/descriptor sets/pipelines/memory barriers bindings). `VkCommandBuffer` must be provided as a pointer in `VkFFTLaunchParams`.

VkCommandBuffer must be in the writing stage, started with vkBeginCommandBuffer call. After VkFFTAppend has finished, provided VkCommandBuffer will contain a sequence of operations performing FFT. The first call of the sequence has no input memory barrier, the last call has one, ensuring FFT has finished execution.

- CUDA/HIP API: if the user wants to use streams, they have to be provided during the application configuration stage. VkFFTAppend performs a series of cuLaunchKernel, which are sequential if appended to one stream and synchronized if appended to multiple streams.
- OpenCL API: similar to Vulkan, VkFFT appends a sequence of clEnqueueNDRangeKernel calls to user-defined cl_command_queue. Currently, they are all assumed to be sequential. cl_command_queue must be provided as a pointer in VkFFTLaunchParams.
- Level Zero API: similar to Vulkan, VkFFT appends a sequence of zeCommandListAppendLaunchKernel calls to user-defined command list ze_command_list_handle_t. They have execution barriers between. ze_command_list_handle_t must be provided as a pointer in VkFFTLaunchParams.
- Metal API: similar to Vulkan, VkFFT appends a sequence of dispatchThreads calls to user-defined command encoder MTL::ComputeCommandEncoder. MTL::ComputeCommandEncoder and its MTL::CommandBuffer must be provided as a pointer in VkFFTLaunchParams.

If VkFFT fails during the VkFFTAppend call, it will not free the application and allocated there resources - use a separate call for that.

4.2.3 Function deleteVkFFT()

```
void deleteVkFFT(VkFFTApplication* app)
```

Performs deallocation of resources used in the provided application. Returns application to the zero-initialized state.

4.2.4 Function VkFFTGetVersion()

```
int VkFFTGetVersion()
```

Returns the version of the VkFFT library in the X.XX.XX format (without dots).

4.3 VkFFT configuration

This section will cover all the parameters that can be specified in the VkFFTConfiguration struct. It will start with a short description of the struct (intended to be used as a cheat sheet), then go for each field in detail. Here pfUINT and pfINT are defined as uint64_t and int64_t respectively.

```
typedef struct {  
    // Required parameters:
```

```

pfUINT FFTdim; // FFT dimensionality (1, 2 or 3)
pfUINT size[VKFFT_MAX_FFT_DIMENSIONS]; // WHD+ - system
    ↪ dimensions
#if(VKFFT_BACKEND==0) //Vulkan API
VkPhysicalDevice* physicalDevice; // Pointer to Vulkan
    ↪ physical device, obtained from vkEnumeratePhysicalDevices
VkDevice* device; // Pointer to Vulkan device, created with
    ↪ vkCreateDevice
VkQueue* queue; // Pointer to Vulkan queue, created with
    ↪ vkGetDeviceQueue
VkCommandPool* commandPool; // Pointer to Vulkan command
    ↪ pool, created with vkCreateCommandPool
VkFence* fence; // Pointer to Vulkan fence, created with
    ↪ vkCreateFence
pfUINT isCompilerInitialized; // Specify if glslang compiler
    ↪ has been intialized before (0 - off, 1 - on). Default 0
#elif(VKFFT_BACKEND==1) //CUDA API
CUdevice* device; // Pointer to CUDA device, obtained from
    ↪ cuDeviceGet
cudaStream_t* stream; // Pointer to streams (can be more than
    ↪ 1), where to execute the kernels. Deafult 0
pfUINT num_streams; // Try to submit CUDA kernels in
    ↪ multiple streams for asynchronous execution. Default 0, set
    ↪ to >=1 if you pass values in the stream pointer.
#elif(VKFFT_BACKEND==2) //HIP API
hipDevice_t* device; // Pointer to HIP device, obtained from
    ↪ hipDeviceGet
hipStream_t* stream; // Pointer to streams (can be more than
    ↪ 1), where to execute the kernels. Deafult 0
pfUINT num_streams; // Try to submit HIP kernels in
    ↪ multiple streams for asynchronous execution. Default 0, set
    ↪ to >=1 if you pass values in the stream pointer.
#elif(VKFFT_BACKEND==3) //OpenCL API
cl_platform_id* platform; // NOT REQUIRED
cl_device_id* device; // Pointer to OpenCL device, obtained
    ↪ from clGetDeviceIDs
cl_context* context; // Pointer to OpenCL context, obtained
    ↪ from clCreateContext
#elif(VKFFT_BACKEND==4) //Level Zero API
ze_device_handle_t* device; // Pointer to Level Zero
    ↪ device, obtained from zeDeviceGet
ze_context_handle_t* context; // Pointer to Level Zero
    ↪ context, obtained from zeContextCreate

```

```

ze_command_queue_handle_t* commandQueue;    // Pointer to
↳ Level Zero command queue with compute and copy
↳ capabilities, obtained from zeCommandQueueCreate
uint32_t commandQueueID;    // ID of the commandQueue with
↳ compute and copy capabilities
#elif(VKFFT_BACKEND==5)
MTL::Device* device;    // Pointer to Metal device, obtained
↳ from MTL::CopyAllDevices
MTL::CommandQueue* queue;    // Pointer to Metal queue, obtained
↳ from device->newCommandQueue()
#endif

// Data parameters (buffers can be specified at launch):
pfUINT userTempBuffer;    // Buffer allocated by app
↳ automatically if needed to reorder Four step algorithm.
↳ Setting to non zero value enables manual user allocation (0
↳ - off, 1 - on)
pfUINT bufferNum;    // Multiple buffer sequence storage is
↳ Vulkan only. Default 1
pfUINT tempBufferNum;    // Multiple buffer sequence storage is
↳ Vulkan only. Default 1, buffer allocated by app
↳ automatically if needed to reorder Four step algorithm.
↳ Setting to non zero value enables manual user allocation
↳
pfUINT inputBufferNum;    // Multiple buffer sequence storage is
↳ Vulkan only. Default 1, if isInputFormatted is enabled
pfUINT outputBufferNum;    // Multiple buffer sequence storage
↳ is Vulkan only. Default 1, if isOutputFormatted is enabled
pfUINT kernelNum;    // Multiple buffer sequence storage is
↳ Vulkan only. Default 1, if performConvolution is enabled
pfUINT* bufferSize;    // Array of buffers sizes in bytes
pfUINT* tempBufferSize;    // Array of temp buffers sizes in
↳ bytes. Default set to bufferSize sum, buffer allocated by
↳ app automatically if needed to reorder Four step algorithm.
↳ Setting to non zero value enables manual user allocation
pfUINT* inputBufferSize;    // Array of input buffers sizes in
↳ bytes, if isInputFormatted is enabled
pfUINT* outputBufferSize;    // Array of output buffers sizes in
↳ bytes, if isOutputFormatted is enabled
pfUINT* kernelSize;    // Array of kernel buffers sizes in bytes,
↳ if performConvolution is enabled
#if(VKFFT_BACKEND==0) //Vulkan API
VkBuffer* buffer;    // Pointer to array of buffers (or one
↳ buffer) used for computations

```

```

VkBuffer* tempBuffer;    // Needed if reorderFourStep is enabled
↳ to transpose the array. Same sum size or bigger as buffer
↳ (can be split in multiple). Default 0. Setting to non zero
↳ value enables manual user allocation
VkBuffer* inputBuffer;  // Pointer to array of input buffers
↳ (or one buffer) used to read data from if isInputFormatted
↳ is enabled
VkBuffer* outputBuffer; // Pointer to array of output
↳ buffers (or one buffer) used to write data to if
↳ isOutputFormatted is enabled
VkBuffer* kernel;       // Pointer to array of kernel buffers (or
↳ one buffer) used to read kernel data from if
↳ performConvolution is enabled
#elif(VKFFT_BACKEND==1) //CUDA API
void** buffer; // Pointer to device buffer used for
↳ computations
void** tempBuffer; // Needed if reorderFourStep is enabled to
↳ transpose the array. Same size as buffer. Default 0.
↳ Setting to non zero value enables manual user allocation
void** inputBuffer; // Pointer to device buffer used to
↳ read data from if isInputFormatted is enabled
void** outputBuffer; // Pointer to device buffer used to
↳ write data to if isOutputFormatted is enabled
void** kernel; // Pointer to device buffer used to read kernel
↳ data from if performConvolution is enabled
#elif(VKFFT_BACKEND==2) //HIP API
void** buffer; // Pointer to device buffer used for
↳ computations
void** tempBuffer; // Needed if reorderFourStep is enabled to
↳ transpose the array. Same size as buffer. Default 0.
↳ Setting to non zero value enables manual user allocation
void** inputBuffer; // Pointer to device buffer used to
↳ read data from if isInputFormatted is enabled
void** outputBuffer; // Pointer to device buffer used to
↳ write data to if isOutputFormatted is enabled
void** kernel; // Pointer to device buffer used to read kernel
↳ data from if performConvolution is enabled
#elif(VKFFT_BACKEND==3) //OpenCL API
cl_mem* buffer; // Pointer to device buffer used for
↳ computations
cl_mem* tempBuffer; // Needed if reorderFourStep is enabled to
↳ transpose the array. Same size as buffer. Default 0.
↳ Setting to non zero value enables manual user allocation

```

```

cl_mem* inputBuffer;    // Pointer to device buffer used to
↳ read data from if isInputFormatted is enabled
cl_mem* outputBuffer;   // Pointer to device buffer used to
↳ write data to if isOutputFormatted is enabled
cl_mem* kernel; // Pointer to device buffer used to read kernel
↳ data from if performConvolution is enabled
#elif(VKFFT_BACKEND==4)
void** buffer; // Pointer to device buffer used for
↳ computations
void** tempBuffer; // Needed if reorderFourStep is enabled to
↳ transpose the array. Same size as buffer. Default 0.
↳ Setting to non zero value enables manual user allocation
void** inputBuffer;    // Pointer to device buffer used to
↳ read data from if isInputFormatted is enabled
void** outputBuffer;   // Pointer to device buffer used to
↳ read data from if isOutputFormatted is enabled
void** kernel; // Pointer to device buffer used to read kernel
↳ data from if performConvolution is enabled
#elif(VKFFT_BACKEND==5)
MTL::Buffer** buffer; // Pointer to device buffer used for
↳ computations
MTL::Buffer** tempBuffer; // Needed if reorderFourStep is
↳ enabled to transpose the array. Same size as buffer.
↳ Default 0. Setting to non zero value enables manual user
↳ allocation
MTL::Buffer** inputBuffer; // Pointer to device buffer used to
↳ read data from if isInputFormatted is enabled
MTL::Buffer** outputBuffer; // Pointer to device buffer
↳ used to read data from if isOutputFormatted is enabled
MTL::Buffer** kernel; // Pointer to device buffer used to
↳ read kernel data from if performConvolution is enabled
#endif
pfUINT bufferOffset; // Specify if VkFFT has to offset the
↳ first element position inside the buffer. In bytes. Default
↳ 0
pfUINT tempBufferOffset; // Specify if VkFFT has to offset
↳ the first element position inside the temp buffer. In
↳ bytes. Default 0
pfUINT inputBufferOffset; // Specify if VkFFT has to offset
↳ the first element position inside the input buffer. In
↳ bytes. Default 0
pfUINT outputBufferOffset; // Specify if VkFFT has to offset
↳ the first element position inside the output buffer. In
↳ bytes. Default 0

```



```

pfUINT kernelOffset;    // Specify if VkFFT has to offset the
↳ first element position inside the kernel. In bytes. Default
↳ 0
pfUINT specifyOffsetsAtLaunch; // Specify if offsets will be
↳ selected with launch parameters VkFFTLaunchParams (0 - off,
↳ 1 - on). Default 0

// Optional: (default 0 if not stated otherwise)
pfUINT coalescedMemory;    // In bytes, for Nvidia and AMD is
↳ equal to 32, Intel is equal 64, scaled for half precision.
↳ Going to work regardless, but if specified by user
↳ correctly, the performance will be higher.
pfUINT aimThreads; // Aim at this many threads per block.
↳ Default 128
pfUINT numSharedBanks; // How many banks shared memory has.
↳ Default 32
pfUINT inverseReturnToInputBuffer; // return data to the input
↳ buffer in inverse transform (0 - off, 1 - on).
↳ isInputFormatted must be enabled
pfUINT numberBatches; // N - used to perform multiple batches
↳ of initial data. Default 1
pfUINT useUint64; // Use 64-bit addressing mode in generated
↳ kernels
pfUINT omitDimension[VKFFT_MAX_FFT_DIMENSIONS]; // Disable
↳ FFT for this dimension (0 - FFT enabled, 1 - FFT disabled).
↳ Default 0. Doesn't work for R2C first axis for now. Doesn't
↳ work with convolutions.
pfUINT performBandwidthBoost; // Try to reduce coalesced
↳ number by a factor of X to get bigger sequence in one
↳ upload for strided axes. Default: -1 for DCT/DST, 2 for
↳ Bluestein's algorithm (or -1 if DCT/DST), 0 otherwise

pfUINT doublePrecision; // Perform calculations in double
↳ precision (0 - off, 1 - on).
pfUINT quadDoubleDoublePrecision; // Perform calculations in
↳ double-double emulation of quad precision (0 - off, 1 -
↳ on).
pfUINT quadDoubleDoublePrecisionDoubleMemory; // Perform
↳ calculations in double-double emulation of quad precision,
↳ while all memory storage is done in FP64.
pfUINT halfPrecision; // Perform calculations in half
↳ precision (0 - off, 1 - on)

```



```

pfUINT halfPrecisionMemoryOnly;    // Use half precision only
↳ as input/output buffer. Input/Output have to be allocated
↳ as half, buffer/tempBuffer have to be allocated as float
↳ (out-of-place mode only). Specify isInputFormatted and
↳ isOutputFormatted to use (0 - off, 1 - on)
pfUINT doublePrecisionFloatMemory; // Use FP64 precision for
↳ all calculations, while all memory storage is done in FP32.

pfUINT performR2C; // Perform R2C/C2R decomposition (0 - off,
↳ 1 - on)
pfUINT performDCT; // Perform DCT transformation (X - DCT
↳ type, 1-4)
pfUINT performDST; // Perform DST transformation (X - DST
↳ type, 1-4)
pfUINT disableMergeSequencesR2C; // Disable merging of two
↳ real sequences to reduce calculations (0 - off, 1 - on)
pfUINT normalize; // Normalize inverse transform (0 - off, 1
↳ - on)
pfUINT disableReorderFourStep; // Disables unshuffling of
↳ Four step algorithm. Requires tempbuffer allocation (0 -
↳ off, 1 - on)
pfINT useLUT; // Switches from calculating sincos to using
↳ precomputed LUT tables (0 - off, 1 - on). Configured by
↳ initialization routine
pfINT useLUT_4step; // Switches from calculating sincos to
↳ using precomputed LUT tables for intermediate roots of 1 in
↳ the Four-step FFT algorithm. (-1 - off, 0 - auto, 1 - on).
↳ Configured by initialization routine
pfUINT makeForwardPlanOnly; // Generate code only for
↳ forward FFT (0 - off, 1 - on)
pfUINT makeInversePlanOnly; // Generate code only for
↳ inverse FFT (0 - off, 1 - on)
pfUINT bufferStride[VKFFT_MAX_FFT_DIMENSIONS]; // Buffer
↳ strides - default set to x - x*y - x*y*z values
pfUINT isInputFormatted; // Specify if input buffer is
↳ padded - 0 - padded, 1 - not padded. For example if it is
↳ not padded for R2C if out-of-place mode is selected (only
↳ if numberBatches==1 and numberKernels==1)
pfUINT isOutputFormatted; // Specify if output buffer is
↳ padded - 0 - padded, 1 - not padded. For example if it is
↳ not padded for R2C if out-of-place mode is selected (only
↳ if numberBatches==1 and numberKernels==1)

```

```

pfUINT inputBufferStride[VKFFT_MAX_FFT_DIMENSIONS];    //
↳ Input buffer strides. Used if isInputFormatted is enabled.
↳ Default set to bufferStride values
pfUINT outputBufferStride[VKFFT_MAX_FFT_DIMENSIONS];    //
↳ Output buffer strides. Used if isInputFormatted is enabled.
↳ Default set to bufferStride values
pfUINT swapTo2Stage4Step; // Specify at which number to switch
↳ from 1 upload to 2 upload 4-step FFT, in case if making max
↳ sequence size lower than coalesced sequence helps to combat
↳ TLB misses. Default 0 - disabled.
pfUINT swapTo3Stage4Step; // Specify at which power of 2 to
↳ switch from 2 upload to 3 upload 4-step FFT, in case if
↳ making max sequence size lower than coalesced sequence
↳ helps to combat TLB misses. Default 0 - disabled. Must be
↳ at least 17
pfUINT considerAllAxesStrided; // Will create plan for
↳ non-strided axis similar as a strided axis - used with
↳ disableReorderFourStep to get the same layout for Bluestein
↳ kernel (0 - off, 1 - on)

pfUINT keepShaderCode; // Will keep shader code and print all
↳ executed shaders during the plan execution in order (0 -
↳ off, 1 - on)
pfUINT printMemoryLayout; // Will print order of buffers used
↳ in shaders (0 - off, 1 - on)

pfUINT saveApplicationToString; // Will save all compiled
↳ binaries to VkFFTApplication.saveApplicationString (will
↳ be allocated by VkFFT, deallocated with deleteVkFFT call).
↳ VkFFTApplication.applicationStringSize will contain size
↳ of binary in bytes. Currently disabled in Metal backend. (0
↳ - off, 1 - on)
pfUINT loadApplicationFromString; // Will load all binaries
↳ from loadApplicationString instead of recompiling them
↳ (must be allocated by user, must contain what
↳ saveApplicationToString call generated previously in
↳ VkFFTApplication.saveApplicationString). Currently
↳ disabled in Metal backend. (0 - off, 1 - on). Mutually
↳ exclusive with saveApplicationToString
void* loadApplicationString; // Memory binary array through
↳ which user can load VkFFT binaries, must be provided by
↳ user if loadApplicationFromString = 1. Use rb/wb flags to
↳ load/save.

```

```

pfUINT disableSetLocale;    // disables all VkFFT attempts to
    ↪ set locale to C - user must ensure that VkFFT has C locale
    ↪ during the plan initialization. This option is needed for
    ↪ multithreading. Default 0.

// Optional Bluestein optimizations: (default 0 if not stated
    ↪ otherwise)
pfUINT fixMaxRadixBluestein;    // controls the padding of
    ↪ sequences in Bluestein convolution. If specified, padded
    ↪ sequence will be made of up to fixMaxRadixBluestein primes.
    ↪ Default: 2 for CUDA and Vulkan/OpenCL/HIP up to 1048576
    ↪ combined dimension FFT system, 7 for Vulkan/OpenCL/HIP past
    ↪ after. Min = 2, Max = 13.
pfUINT forceBluesteinSequenceSize; // force the sequence size
    ↪ to pad to in Bluestein's algorithm. Must be at least 2*N-1
    ↪ and decomposable with primes 2-13.
pfUINT useCustomBluesteinPaddingPattern;    // force the
    ↪ sequence sizes to pad to in Bluestein's algorithm, but on a
    ↪ range. This number specifies the number of elements in
    ↪ primeSizes and in paddedSizes arrays. primeSizes - array of
    ↪ non-decomposable as radix scheme sizes - 17, 23, 31 etc.
    ↪ paddedSizes - array of lengths to pad to. paddedSizes[i]
    ↪ will be the padding size for all non-decomposable sequences
    ↪ from primeSizes[i] to primeSizes[i+1] (will use default
    ↪ scheme after last one) - 42, 60, 64 for primeSizes before
    ↪ and 37+ will use default scheme (for example). Default is
    ↪ vendor and API-based specified in
    ↪ autoCustomBluesteinPaddingPattern.
pfUINT* primeSizes; // described in
    ↪ useCustomBluesteinPaddingPattern
pfUINT* paddedSizes;    // described in
    ↪ useCustomBluesteinPaddingPattern

pfUINT fixMinRaderPrimeMult;    // start direct multiplication
    ↪ Rader's algorithm for radix primes from this number. This
    ↪ means that VkFFT will inline custom Rader kernels if
    ↪ sequence is divisible by these primes. Default is 17, as
    ↪ VkFFT has kernels for 2-13. If you make it less than 13,
    ↪ VkFFT will switch from these kernels to Rader.
pfUINT fixMaxRaderPrimeMult;    // switch from Mult Rader's
    ↪ algorithm for radix primes from this number. Current
    ↪ limitation for Rader is maxThreadNum/2+1, realistically you
    ↪ would want to switch somewhere on 30-100 range. Default is
    ↪ vendor-specific (currently ~40)

```

```

pfUINT fixMinRaderPrimeFFT;    // start FFT convolution
↳ version of Rader for radix primes from this number. Better
↳ than direct multiplication version for almost all primes
↳ (except small ones, like 17-23 on some GPUs). Must be
↳ bigger or equal to fixMinRaderPrimeMult. Default 29 on AMD
↳ and 17 on other GPUs.
pfUINT fixMaxRaderPrimeFFT;    // switch to Bluestein's
↳ algorithm for radix primes from this number. Switch may
↳ happen earlier if prime can't fit in shared memory. Default
↳ is 16384, which is bigger than most current GPU's shared
↳ memory.

// Optional zero padding control parameters: (default 0 if not
↳ stated otherwise)
pfUINT performZeroPadding[VKFFT_MAX_FFT_DIMENSIONS]; //
↳ Don't read some data/perform computations if some input
↳ sequences are zero-padded for each axis (0 - off, 1 - on)
pfUINT fft_zeropad_left[VKFFT_MAX_FFT_DIMENSIONS]; // Specify
↳ start boundary of zero block in the system for each axis
pfUINT fft_zeropad_right[VKFFT_MAX_FFT_DIMENSIONS]; //
↳ Specify end boundary of zero block in the system for each
↳ axis
pfUINT frequencyZeroPadding;    // Set to 1 if zero padding of
↳ frequency domain, default 0 - spatial zero padding

// Optional convolution control parameters: (default 0 if not
↳ stated otherwise)
pfUINT performConvolution; // Perform convolution in this
↳ application (0 - off, 1 - on). Disables reorderFourStep
↳ parameter
pfUINT coordinateFeatures; // C - coordinate, or dimension of
↳ features vector. In matrix convolution - size of a vector
pfUINT matrixConvolution; // If equal to 2 perform 2x2, if
↳ equal to 3 perform 3x3 matrix-vector convolution. Overrides
↳ coordinateFeatures
pfUINT symmetricKernel; // Specify if kernel in 2x2 or 3x3
↳ matrix convolution is symmetric
pfUINT numberKernels; // N - only used in convolution step -
↳ specify how many kernels were initialized before. Expands
↳ one input to multiple (batched) output

```

```

pfUINT kernelConvolution;    // Specify if this application is
    ↪ used to create kernel for convolution, so it has the same
    ↪ properties. performConvolution has to be set to 0 for
    ↪ kernel creation

// Register overutilization (experimental): (default 0 if not
    ↪ stated otherwise)
pfUINT registerBoost;    // Specify if register file size is
    ↪ bigger than shared memory and can be used to extend it X
    ↪ times (on Nvidia 256KB register file can be used instead of
    ↪ 32KB of shared memory, set this constant to 4 to emulate
    ↪ 128KB of shared memory). Defaults: Nvidia - 4 in
    ↪ Vulkan/OpenCL, 1 in CUDA backend; AMD - 2 if shared memory
    ↪ >= 64KB, else 4 in Vulkan/OpenCL backend, 1 in HIP backend;
    ↪ Intel - 1 if shared memory >= 64KB, else 2 in
    ↪ Vulkan/OpenCL/Level Zero backends, 1 in Metal; Default 1
pfUINT registerBoostNonPow2;    // Specify if register
    ↪ overutilization should be used on non power of 2 sequences
    ↪ (0 - off, 1 - on)
pfUINT registerBoost4Step;    // Specify if register file
    ↪ overutilization should be used in big sequences (>2^14),
    ↪ same definition as registerBoost. Default 1

//not used techniques:
pfUINT devicePageSize;    // In KB, the size of a page on the
    ↪ GPU. Setting to 0 disables local buffer split in pages
pfUINT localPageSize;    // In KB, the size to split page into
    ↪ if sequence spans multiple devicePageSize pages

// Automatically filled based on device info (still can be
    ↪ reconfigured by user):
pfUINT computeCapabilityMajor;    // CUDA/HIP compute capability
    ↪ of the device
pfUINT computeCapabilityMinor;    // CUDA/HIP compute capability
    ↪ of the device
pfUINT maxComputeWorkGroupCount[VKFFT_MAX_FFT_DIMENSIONS];    //
    ↪ maxComputeWorkGroupCount from VkPhysicalDeviceLimits
pfUINT maxComputeWorkGroupSize[VKFFT_MAX_FFT_DIMENSIONS];    //
    ↪ maxComputeWorkGroupCount from VkPhysicalDeviceLimits
pfUINT maxThreadsNum;    // Max number of threads from
    ↪ VkPhysicalDeviceLimits
pfUINT sharedMemorySizeStatic;    // Available for static
    ↪ allocation shared memory size, in bytes

```

```

pfUINT sharedMemorySize;    // Available for allocation shared
    ↪ memory size, in bytes
pfUINT sharedMemorySizePow2; // Power of 2 which is less or
    ↪ equal to sharedMemorySize, in bytes
pfUINT warpSize;           // Number of threads per warp/wavefront.
pfUINT halfThreads;        // Intel fix
pfUINT allocateTempBuffer; // Buffer allocated by app
    ↪ automatically if needed to reorder Four step algorithm.
    ↪ Parameter to check if it has been allocated
pfUINT reorderFourStep;    // Unshuffle Four step algorithm.
    ↪ Requires tempbuffer allocation (0 - off, 1 - on). Default
    ↪ 1.
pfINT maxCodeLength;       // Specify how big can be buffer used
    ↪ for code generation (in char). Default 1000000 chars.
pfINT maxTempLength;       // Specify how big can be buffer used
    ↪ for intermediate string sprintfs be (in char). Default 5000
    ↪ chars. If code segfaults for some reason - try increasing
    ↪ this number.
pfUINT autoCustomBluesteinPaddingPattern; // default value for
    ↪ useCustomBluesteinPaddingPattern
pfUINT useRaderUintLUT;    // allocate additional LUT to store
    ↪ g_pow
pfUINT vendorID;           // vendorID 0x10DE - NVIDIA, 0x8086 - Intel,
    ↪ 0x1002 - AMD, etc
#ifdef VKFFT_BACKEND==0 //Vulkan API
VkDeviceMemory tempBufferDeviceMemory; // Filled at app
    ↪ creation
VkCommandBuffer* commandBuffer;         // Filled at app execution
VkMemoryBarrier* memory_barrier;        // Filled at app creation
#elif VKFFT_BACKEND==1 //CUDA API
cudaEvent_t* stream_event;              // Filled at app creation
pfUINT streamCounter;                   // Filled at app creation
pfUINT streamID;                        // Filled at app creation
#elif VKFFT_BACKEND==2 //HIP API
hipEvent_t* stream_event;               // Filled at app creation
pfUINT streamCounter;                   // Filled at app creation
pfUINT streamID;                        // Filled at app creation
#elif VKFFT_BACKEND==3 //OpenCL API
cl_command_queue* commandQueue;         // Filled at app creation
#elif VKFFT_BACKEND==4
ze_command_list_handle_t* commandList;  // Filled at app
    ↪ creation
#elif VKFFT_BACKEND==5
MTL::CommandBuffer* commandBuffer;      // Filled at app execution

```



```

MTL::ComputeCommandEncoder* commandEncoder;    // Filled at
    ↪ app execution
#endif
} VkFFTConfiguration;

```

4.3.1 Driver API parameters

In order to work, VkFFT needs some structures that are provided by the driver. They are backend API-dependent. VkFFT will return corresponding VkFFTResult if one of these structures are not provided (value equal to zero) unless it is stated that there is a default value assigned. VkFFT will not modify provided values directly.

Vulkan API will need the following information:

- `VkPhysicalDevice*` `physicalDevice` - Pointer to Vulkan physical device, obtained from `vkEnumeratePhysicalDevices()`
- `VkDevice*` `device` - Pointer to Vulkan device, created with `vkCreateDevice()`
- `VkQueue*` `queue` - Pointer to Vulkan queue, created with `vkGetDeviceQueue()`
- `VkCommandPool*` `commandPool` - Pointer to Vulkan command pool, created with `vkCreateCommandPool()`
- `VkFence*` `fence` - Pointer to Vulkan fence, created with `vkCreateFence()`
- `uint64_t` `isCompilerInitialized` - Specify if glslang compiler has been initialized before (0 - off, 1 - on). Default 0 - VkFFT will call `glslang_initialize_process()` at `initializeVkFFT()` and `glslang_finalize_process()` at `deleteVkFFT()` calls.

CUDA API will need the following information:

- `CUdevice*` `device` - Pointer to CUDA device, obtained from `cuDeviceGet()`
- `cudaStream_t*` `stream` - Pointer to streams (can be more than 1), where to execute the kernels. Default 0. Streams must be associated with the provided device. There is no real benefit in having more than one, however.
- `uint64_t` `num_streams` - Try to submit CUDA kernels in multiple streams for asynchronous execution. Default 0, set to ≥ 1 if you pass values in the stream pointer.

HIP API will need the following information:

- `hipDevice_t*` `device` - Pointer to HIP device, obtained from `hipDeviceGet()`
- `hipStream_t*` `stream` - Pointer to streams (can be more than 1), where to execute the kernels. Default 0. Streams must be associated with the provided device. There is no real benefit in having more than one, however.
- `uint64_t` `num_streams` - Try to submit HIP kernels in multiple streams for asynchronous execution. Default 0, set to ≥ 1 if you pass values in the stream pointer.

OpenCL API will need the following information:

- `cl_device_id*` device - Pointer to OpenCL device, obtained from `clGetDeviceIDs()`
- `cl_context*` context - Pointer to OpenCL context, obtained from `clCreateContext()`

Level Zero API will need the following information:

- `ze_device_handle_t*` device - Pointer to Level Zero device, obtained from `zeDeviceGet()`
- `ze_context_handle_t*` context - Pointer to Level Zero context, obtained from `zeContextGet()`
- `ze_command_queue_handle_t*` commandQueue - Pointer to Level Zero command queue with compute and copy capabilities, obtained from `zeCommandQueueCreate()`
- `uint32_t` commandQueueID - ID of the commandQueue with compute and copy capabilities

Metal API will need the following information:

- `MTL::Device*` device - Pointer to Metal device, obtained from `MTL::CopyAllDevices`
- `MTL::CommandQueue*` queue - Pointer to Metal queue, obtained from `device->newCommandQueue()`

4.3.2 Memory management parameters

There are five buffer types user can provide to `VkFFT`:

- the main buffer (buffer)
- temporary buffer used for calculations requiring out-of-place writes (tempBuffer)
- separate input buffer, from which initial read is performed (inputBuffer)
- separate output buffer, to which final write is performed (outputBuffer)
- kernel buffer, used for calculation of convolutions and cross-correlations (kernel)

These buffers must be passed by a pointer: in Vulkan API they are provided as `VkBuffer*`, in CUDA, HIP and Level Zero they are provided as `void**`, in OpenCL they are provided as `cl_mem*`, in Metal they are provided as `MTL::Buffer*`. Even though the underlying structure (`VkBuffer`, `void*`, `cl_mem`, `MTL::Buffer*`) is not a memory but just a number that the driver can use to access corresponding allocated memory on the GPU, passing them by a pointer allows for the user to query multiple GPU allocated buffers for `VkFFT` to use. Currently, it is only supported in Vulkan API - each of five buffer types can be made out of multiple separate memory allocations. For example, it is possible to combine multiple small unused at the point of FFT calculation buffers to form a tempBuffer. This option also allows Vulkan API to overcome the limit of 4GB for a single memory allocation - due to the fact that Vulkan can only use 32-bit numbers for addressing (other APIs support 64-bit addressing).

To use the buffers other than the main buffer, the user has to specify this in configuration at the application creation stage (set to zero by default, optional parameters):

- `uint64_t userTempBuffer` - enables manual temporary buffer allocation (otherwise it is managed by `VkFFT`)
- `uint64_t isInputFormatted` - specifies that initial read is performed from a separate buffer (`inputBuffer`)
- `uint64_t isOutputFormatted` - specifies that final write is performed to a separate buffer (`outputBuffer`)
- `uint64_t performConvolution` - enables convolution calculations, which requires precomputed kernel (`kernel`)

Buffer sizes (`bufferSize/tempBufferSize/inputBufferSize/outputBufferSize/kernelSize`) are provided as a `uint64_t` pointer to an array, where each element corresponds to the buffer size of the buffer with the same placement in the buffer array. Buffer sizes have to be provided in Vulkan API (due to the stricter memory management model and multiple buffer support) and are optional in other backends (they can be useful to determine when to switch for 64-bit addressing).

Buffer number (`bufferNum/tempBufferNum/inputBufferNum/outputBufferNum/kernelNum`) corresponds to how many elements are in the buffer and buffer size array. By default it is set to 1 and is not required to be provided by the user. Non-Vulkan backends currently don't support values other than default. Optional parameter.

Buffer offset (`bufferOffset/ tempBufferOffset/ inputBufferOffset/ outputBufferOffset/ kernelOffset`) specifies offset from the start of the buffer sequence. It must be specified in bytes and must be divisible by the number type size used in the corresponding array (otherwise, the offset will be truncated). It is provided as a single `uint64_t` value. Can be provided at launch time, if `specifyOffsetsAtLaunch` parameter is enabled during initialization call. Optional parameters.

User can provide custom dimension strides for `buffer/inputBuffer/outputBuffer` buffers - `uint64_t[VKFFT_MAX_FFT_DIMENSIONS]` array. Strides are specified in elements used in the array (not bytes). The first element corresponds to the stride between elements in the H direction, the second corresponds to the D direction and the third to C (or N, if the number of elements in C is 1). The first axis is assumed to be non-strided. Must be at least of the same size as default strides, otherwise the behavior is undefined. Optional parameters.

`uint64_t inverseReturnToInputBuffer` - an option that allows setting the final output buffer of the inverse transform to the same buffer, initial read of forward transform is performed from (`inputBuffer`, if `isInputFormatted` enabled). Optional parameter.

4.3.3 General FFT parameters

This section describes part of the configuration structure responsible for FFT specification.

`uint64_t FFTdim` - dimensionality of the transform (up to `VKFFT_MAX_FFT_DIMENSIONS`). Required parameter.

uint64_t size[VKFFT_MAX_FFT_DIMENSIONS] - WHD+ dimensions of the transform. Required parameter.

uint64_t numberBatches - N parameter of the transform. By default, it is set to 1. Optional parameter.

uint64_t performR2C - perform R2C/C2R decomposition. performDCT and performDST must be set to 0. Default 0, set to 1 to enable. Optional parameter.

uint64_t performDCT - perform DCT transformation. performR2C and performDST must be set to 0. Default 0, set to X for DCT-X (currently supported X: 1, 2, 3 and 4). Optional parameter.

uint64_t performDST - perform DST transformation. performR2C and performDCT must be set to 0. Default 0, set to X for DST-X (currently supported X: 1, 2, 3 and 4). Optional parameter.

uint64_t normalize - enabling this parameter will make the inverse transform divide the result by the FFT length. Default 0, set to 1 to enable. Optional parameter.

4.3.4 Precision parameters (and some things that can affect it):

uint64_t doublePrecision - perform calculations in double precision. Default 0, set to 1 to enable. In Vulkan/OpenCL/Level Zero your device must support double precision functionality. Metal API does not support double precision. Optional parameter.

uint64_t doublePrecisionFloatMemory - perform calculations in double precision, but all intermediate and final storage in float. Input/Output/main buffers must have single-precision layout. doublePrecision must be set to 0. This option increases precision, but not that much to be recommended for actual use. Default 0, set to 1 to enable. In Vulkan/OpenCL/Level Zero your device must support double precision functionality. Metal API does not support double precision. Optional parameter.

uint64_t quadDoubleDoublePrecision - perform calculations in double-double emulation of quad precision (0 - off, 1 - on). Optional parameter. Requires quadmath library (for now).

uint64_t quadDoubleDoublePrecisionDoubleMemory - perform calculations in double-double emulation of quad precision, while all memory storage is done in FP64. Optional parameter. Requires quadmath library (for now).

uint64_t halfPrecision - half-precision in VkFFT is implemented only as memory optimization. All calculations are done in single precision (similar way as doublePrecisionFloatMemory works for double and single precision). Default 0, set to 1 to enable. Optional parameter.

uint64_t halfPrecisionMemoryOnly - another way of performing half-precision in VkFFT, it will use half-precision only for initial and final memory storage in input/output buffer. Input/Output have to be allocated as half, buffer/tempBuffer have to be allocated as float (out-of-place mode only). Specify isInputFormatted and isOutputFormatted to use. So, for

example, intermediate storage between axes FFTs in the multidimensional case will be done in single precision, as opposed to half-precision in the base halfPrecision case. halfPrecision must be set to 1. Default 0, set to 1 to enable. Optional parameter.

int64_t useLUT - switches from calculating sines and cosines (via special function units in single precision or as a polynomial approximation in double precision) to using precomputed Look-Up Tables. -1 - off, 0 - auto, 1 - on. Default 0 in single precision, 1 in double precision. Set to 1 by default for Intel GPUs. If you have issues with single-precision accuracy on your GPU, try enabling this parameter (mobile GPUs may be affected). Optional parameter.

int64_t useLUT_4step - switches from calculating sincos to using precomputed LUT tables for intermediate roots of 1 in the Four-step FFT algorithm. (-1 - off, 0 - auto, 1 - on). Configured by initialization routine. Optional parameter.

4.3.5 Advanced parameters (code will work fine without using them)

uint64_t omitDimension[VKFFT_MAX_FFT_DIMENSIONS] - parameter, that disables the FFT calculation for a particular axis (WHD). Note, that omitted dimensions still need to be included in FFTdim and size. This parameter simply works as a switch during execution - by not executing the particular dimension code. It doesn't work with the non-strided axis (W) of R2C/C2R mode. It doesn't work with convolution calculations. Default 0, set to 1 to enable. Optional parameter.

uint64_t useUint64 - forces VkFFT to use 64-bit addressing in generated kernels. It is automatically enabled if the estimated buffer size is more than 4GB. Doesn't work with the Vulkan backend. By default, it is set to 0. Optional parameter.

uint64_t coalescedMemory - number of bytes to coalesce per one transaction. For Nvidia and AMD is equal to 32, Intel is equal to 64. Going to work regardless, but if specified by the user correctly, the performance will be higher. Default 64 for other GPUs. For half-precision should be multiplied by two. Should be a power of two. Optional parameter.

uint64_t numSharedBanks - configure the number of shared banks on the target GPU. Default 32. Minor performance boost as it solves shared memory conflicts for the power of two systems. Optional parameter.

uint64_t aimThreads - try to aim all kernels at this amount of threads. Gains/losses are not predictable, just a parameter to play with (it is not guaranteed that the target kernel will use that many threads). Default 128. Optional parameter.

uint64_t useUint64 - forces 64-bit addressing in generated kernels. Should be enabled automatically for systems spanning more than 4GB, but it is better to have an option to force it as a failsafe. Doesn't work in Vulkan API (use multiple buffer binding). Default 0, set to 1 to enable. Optional parameter.

uint64_t performBandwidthBoost - try to reduce coalesced number by a factor of X to get bigger sequence in one upload for strided axes. Default: -1(inf) for DCT and DST, 2 for Bluestein's algorithm (or -1 if DCT and DST), 0 otherwise

uint64_t disableMergeSequencesR2C - disable the optimization that performs merging of two real sequences to reduce calculations (in R2C/C2R and R2R). If enabled, calculations will be performed by simply setting the imaginary component to zero. Default 0, set to 1 to enable. Optional parameter.

uint64_t disableReorderFourStep - disables unshuffling of the Four Step FFT algorithm (last transposition of data). With this option enabled, tempBuffer will not be needed (unless it is required by Bluestein's multi-upload FFT algorithm). Default 0, set to 1 to enable. Automatically enabled for convolution calculations and Bluestein's algorithm. Optional parameter.

uint64_t makeForwardPlanOnly - generate code only for forward FFT. Default 0, set to 1 to enable. Mutually exclusive with makeInversePlanOnly. Optional parameter.

uint64_t makeInversePlanOnly - generate code only for inverse FFT. Default 0, set to 1 to enable. Mutually exclusive with makeForwardPlan. Optional parameter.

uint64_t considerAllAxesStrided - will create a plan for a non-strided axis similar to a strided axis (used with disableReorderFourStep to get the same layout for Bluestein kernel). Default 0, set to 1 to enable. Optional parameter.

uint64_t keepShaderCode - debugging option, will keep shader code and print all executed shaders during the plan execution in order. Default 0, set to 1 to enable. Optional parameter.

uint64_t printMemoryLayout - debugging option, will print order of buffers used in kernels. Default 0, set to 1 to enable. Optional parameter.

uint64_t saveApplicationToString - will save all compiled binaries to VkFFTApplication.saveApplicationString (will be allocated by VkFFT, deallocated with deleteVkFFT call). VkFFTApplication.applicationStringSize will contain size of binary in bytes. Currently disabled in Metal backend. Default 0, set to 1 to enable. Optional parameter.

uint64_t loadApplicationFromString - will load all binaries from loadApplicationString instead of recompiling them (loadApplicationString must be allocated by user, must contain what saveApplicationToString call generated previously in VkFFTApplication.saveApplicationString). Currently disabled in Metal backend. Default 0, set to 1 to enable. Optional parameter. Mutually exclusive with saveApplicationToString

void* loadApplicationString - memory binary array through which user can load VkFFT binaries, must be provided by user if loadApplicationFromString = 1. Use rb/wb flags to load/save.

uint64_t disableSetLocale - disables all VkFFT attempts to set locale to C - user must ensure that VkFFT has C locale during the plan initialization. This option is needed for multithreading. Default 0.

uint64_t groupedBatch[VKFFT_MAX_FFT_DIMENSIONS] - try to force this many FFTs to be performed by one threadblock for each dimension. Optional parameter.

4.3.6 Rader control parameters

`uint64_t fixMinRaderPrimeMult` - start direct multiplication Rader's algorithm for radix primes from this number. This means that `VkFFT` will inline custom Rader kernels if sequence is divisible by these primes. Default is 17, as `VkFFT` has kernels for 2-13. If you make it less than 13, `VkFFT` will switch from these kernels to Rader.

`uint64_t fixMaxRaderPrimeMult` - switch from Mult Rader's algorithm for radix primes from this number. Current limitation for Rader is $\text{maxThreadNum}/2+1$, realistically you would want to switch somewhere on 30-100 range. Default is vendor-specific (currently ~40)

`uint64_t fixMinRaderPrimeFFT` - start FFT convolution version of Rader for radix primes from this number. Better than direct multiplication version for almost all primes (except small ones, like 17-23 on some GPUs). Must be bigger or equal to `fixMinRaderPrimeMult`. Default 29 on AMD and 17 on other GPUs.

`uint64_t fixMaxRaderPrimeFFT` - switch to Bluestein's algorithm for radix primes from this number. Switch may happen earlier if prime can't fit in shared memory. Default is 16384, which is bigger than most current GPU's shared memory.

4.3.7 Bluestein control parameters

If the sequence can not be decomposed as a multiplication of primes up to 13, FFT is performed as a convolution. The sequence to pad to with the best performance is usually device-dependent. `VkFFT` uses parameters manually tuned for all sequences between 2 and 4096 for both double and single precision on Nvidia A100 (Nvidia profile) and AMD MI250 (default profile). To control this process, `VkFFT` allows for the following parameters specification:

`uint64_t fixMaxRadixBluestein` - controls the padding of sequences in Bluestein convolution. If specified, padded sequence will be made of up to `fixMaxRadixBluestein` primes. Default: 2 for CUDA and Vulkan/OpenCL/HIP up to 1048576 combined dimension FFT system, 7 for Vulkan/OpenCL/HIP past after. Min = 2, Max = 13.

`uint64_t forceBluesteinSequenceSize` - force the sequence size to pad to in Bluestein's algorithm. Must be at least $2*N-1$ and decomposable with primes 2-13.

`uint64_t useCustomBluesteinPaddingPattern` - force the sequence sizes to pad to in Bluestein's algorithm, but on a range. This number specifies the number of elements in `primeSizes` and in `paddedSizes` arrays. `primeSizes` - array of non-decomposable as radix scheme sizes - 17, 23, 31 etc. `paddedSizes` - array of lengths to pad to. `paddedSizes[i]` will be the padding size for all non-decomposable sequences from `primeSizes[i]` to `primeSizes[i+1]` (will use default scheme after last one) - 42, 60, 64 for `primeSizes` before and 37+ will use default scheme (for example). Default is vendor and API-based specified in `autoCustomBluesteinPaddingPattern`.

`uint64_t* primeSizes` - described in `useCustomBluesteinPaddingPattern`

`uint64_t* paddedSizes` - described in `useCustomBluesteinPaddingPattern`

4.3.8 Zero padding parameters

`uint64_t performZeropadding[VKFFT_MAX_FFT_DIMENSIONS]` - do not read/write some data/perform computations if some part of the sequence is known to have zeros. Set separately for each axis (WHD). If enabled, all 1D sequences in this direction will be considered padded (independent of other zero-padded axes). Default 0, set to 1 to enable. Optional parameter.

`uint64_t fft_zeropad_left[VKFFT_MAX_FFT_DIMENSIONS]` - specify start boundary of zero block in the system for each axis. Default 0, set to the value between 0 and `size[X]-1`. Optional parameter.

`uint64_t fft_zeropad_right[VKFFT_MAX_FFT_DIMENSIONS]` - specify end boundary of zero block in the system for each axis. Default 0, set to the value between `fft_zeropad_left[X]` and `size[X]-1`. Optional parameter.

`uint64_t frequencyZeroPadding` - enables zero padding of the frequency domain, so the first read of inverse FFT will consider the parts of the system from `fft_zeropad_left` to `fft_zeropad_right` as zero. Default 0 - spatial zero padding, set to 1 to enable. Optional parameter.

4.3.9 Convolution parameters

`uint64_t performConvolution` - main parameter that enables convolutions in the application. If enabled, you must specify kernel buffer, number of kernel buffers and kernel sizes (in Vulkan API). Disables reordering of the Four Step FFT algorithm. Default 0, set to 1 to enable. Optional parameter.

`uint64_t conjugateConvolution` - default 0, set to 1 to enable enables conjugation of the sequence FFT is currently done on, 2 to enable conjugation of the convolution kernel. Optional parameter.

`uint64_t crossPowerSpectrumNormalization` - normalize the FFT * kernel multiplication in frequency domain. Default 0, set to 1 to enable. Optional parameter.

`uint64_t coordinateFeatures` - max coordinate (C), or dimension of the features vector. In matrix convolution - the size of the vector. The main purpose is to support Matrix-Vector convolutions. Use `numberBatches` parameter in tasks, not requiring two separate coordinate-like enumerations of data. Default 1. Optional parameter.

`uint64_t matrixConvolution` - set to 2 to perform 2x2, set to 3 to perform 3x3 matrix-vector convolution. Matrix-vector convolution is a form of point-wise multiplication in the Fourier space, used by the convolution theorem, where multiplication takes the form of Matrix-vector multiplication. Overrides `coordinateFeatures` during execution. Default 0. Optional parameter.

`uint64_t symmetricKernel` - specify if kernel in 2x2 or 3x3 matrix convolution is symmetric. You need to store data as xx, xy, yy (upper-triangular) if enabled and as xx, xy, yx, yy (along rows then along columns, from left to right) if disabled. Default 0, set to 1 to enable. Optional parameter.

uint64_t numberKernels - specify how many kernels were initialized before performing one input/multiple output convolutions. Overwrites numberBatches (N). Only used in convolution step and the following inverse transforms. Default 1. Optional parameter.

uint64_t kernelConvolution - specify if this application is used to create kernel for convolution, so it has the same properties/memory layout. performConvolution has to be set to 0 for the kernel creation. Default 0, set to 1 to enable. Optional parameter, but it is a required parameter for kernel generation.

4.3.10 Register overutilization

Only works in C2C mode, without convolution support. Enabled in Vulkan, OpenCL and Level Zero APIs only (it works in other APIs, but worse, does not work in Metal). Experimental feature.

uint64_t registerBoost - specify if the register file size is bigger than shared memory and can be used to extend it X times (on Nvidia 256KB register file can be used instead of 32KB of shared memory, set this constant to 4 to emulate 128KB of shared memory). Default 1 - no overutilization. In Vulkan, OpenCL and Level Zero it is set to 4 on Nvidia GPUs, to 2 if the driver shows 64KB or more of shared memory on AMD, to 2 if the driver shows less than 64KB of shared memory on AMD, to 1 if the driver shows 64KB or more of shared memory on Intel, to 2 if the driver shows less than 64KB of shared memory on Intel. Optional parameter.

uint64_t registerBoostNonPow2 - specify if register overutilization should be used on non-power of 2 sequences. Default 0, set to 1 to enable. Optional parameter.

uint64_t registerBoost4Step - specify if register file overutilization should be used in big sequences ($>2^{14}$), same definition as registerBoost. Default 1. Optional parameter.

4.3.11 Extra advanced parameters (filled automatically)

uint64_t computeCapabilityMajor - CUDA/HIP compute capability of the device

uint64_t computeCapabilityMinor - CUDA/HIP compute capability of the device

uint64_t maxComputeWorkGroupCount[3] - how many workgroups can be launched at one dispatch. Automatically derived from the driver, can be artificially lowered. Then VkFFT will perform a logical split and extension of the number of workgroups to cover the required range.

uint64_t maxComputeWorkGroupSize[3] - max dimensions of the workgroup. Automatically derived from the driver. Can be modified if there are some issues with the driver (as there were with ROCm 4.0, when it returned 1024 for maxComputeWorkGroupSize and actually supported only up to 256 threads).

uint64_t maxThreadsNum - max number of threads per block. Similar to maxComputeWorkGroupSize, but aggregated. Automatically derived from the driver.

uint64_t sharedMemorySizeStatic - available for static allocation shared memory size, in bytes. Automatically derived from the driver. Can be controlled by the user, if desired.

uint64_t sharedMemorySize - available for allocation shared memory size, in bytes. VkFFT uses dynamic shared memory in CUDA/HIP as it allows for bigger allocations. Automatically derived from the driver. Can be controlled by the user, if desired.

uint64_t sharedMemorySizePow2 - the power of 2 which is less or equal to sharedMemorySize, in bytes. Automatically computed.

uint64_t warpSize - number of threads per warp/wavefront. Automatically derived from the driver, but can be modified (can increase performance, though unpredictable as defaults have good values). Must be a power of two.

uint64_t halfThreads - Intel GPU fix, tries to reduce the amount of dispatched threads in half to solve performance degradation in the Four Step FFT algorithm. Default 0 for other GPUs, try enabling it if performance degrades in the Four Step FFT algorithm for your GPU as well.

int64_t maxCodeLength - specify how big can the buffer used for code generation be (in char). Default 1000000 chars.

int64_t maxTempLength - specify how big can the buffer used for intermediate string printf's be (in char). Default 5000 chars. If code segfaults for some reason - try increasing this number.

uint64_t autoCustomBluesteinPaddingPattern - default value for useCustomBluesteinPaddingPattern

uint64_t useRaderUintLUT - allocate additional LUT to store g_pow

uint64_t vendorID - vendorID 0x10DE - NVIDIA, 0x8086 - Intel, 0x1002 - AMD, etc.

5 VkFFT Benchmark/Precision Suite and utils_VkFFT helper routines

The only licensed (MIT) part of the VkFFT repository is the VkFFT header file - core library. Other files are either external helper libraries (half, glslang, with their respective licenses) or unlicensed code that is intended for simple copy-pasting (benchmark_scripts, utils_VkFFT.h). It is the easiest way to understand how to use VkFFT by taking the provided scripts and tinker them to the particular task. The current version of the benchmark and precision verification suite has the following codes available:

- user_benchmark_VkFFT - generalization of the main configuration parameters that can be used to launch simplest in-place transforms for the most important supported functionality
- Sample 0 - FFT + iFFT C2C benchmark 1D batched in single precision
- Sample 1 - FFT + iFFT C2C benchmark 1D batched in double precision
- Sample 2 - FFT + iFFT C2C benchmark 1D batched in half precision
- Sample 3 - FFT + iFFT C2C multidimensional benchmark in single precision
- Sample 4 - FFT + iFFT C2C multidimensional benchmark in single precision, native zeropadding
- Sample 5 - FFT + iFFT C2C benchmark 1D batched in single precision, no reshuffling
- Sample 6 - FFT + iFFT R2C / C2R benchmark, in-place.
- Sample 7 - FFT + iFFT C2C Bluestein benchmark in single precision
- Sample 8 - FFT + iFFT C2C Bluestein benchmark in double precision
- Sample 9 - FFT + iFFT C2C benchmark 1D batched in double-double emulation of quad precision
- Sample 10 - multiple buffers (4 by default) split version of benchmark 0
- Sample 11 - VkFFT / xFFT / FFTW C2C precision test in single precision (xFFT can be cuFFT or rocFFT)
- Sample 12 - VkFFT / xFFT / FFTW C2C precision test in double precision (xFFT can be cuFFT or rocFFT)
- Sample 13 - VkFFT / cuFFT / FFTW C2C precision test in half precision
- Sample 14 - VkFFT / FFTW C2C radix 3 / 5 / 7 / 11 / 13 / Bluestein precision test in single precision
- Sample 15 - VkFFT / xFFT / FFTW R2C+C2R precision test in single precision, out-of-place. (xFFT can be cuFFT or rocFFT)
- Sample 16 - VkFFT / FFTW R2R DCT-I, II, III and IV precision test in single precision

- Sample 17 - VkFFT / FFTW R2R DCT-I, II, III and IV precision test in double precision
- Sample 18 - VkFFT / FFTW C2C radix 3 / 5 / 7 / 11 / 13 / Bluestein precision test in double precision
- Sample 19 - VkFFT / FFTW C2C precision test in double-double emulation of quad precision
- Sample 50 - convolution example with identity kernel
- Sample 51 - zero padding convolution example with identity kernel
- Sample 52 - batched convolution example with identity kernel
- Sample 100 - VkFFT FFT + iFFT R2R DCT multidimensional benchmark in single precision
- Sample 101 - VkFFT FFT + iFFT R2R DCT multidimensional benchmark in double precision
- Sample 1000 - FFT + iFFT C2C benchmark 1D batched in single precision: all supported systems from 2 to 4096
- Sample 1001 - FFT + iFFT C2C benchmark 1D batched in double precision: all supported systems from 2 to 4096
- Sample 1002 - FFT + iFFT C2C benchmark 1D batched in half precision: all supported systems from 2 to 4096
- Sample 1003 - FFT + iFFT C2C benchmark 1D batched in single precision: all supported systems from 2 to 4096
- Sample 1004 - FFT + iFFT C2C benchmark 1D batched in double-double emulation of quad precision: all supported systems from 2 to 4096

5.1 `utils_VkFFT` helper routines

Launching even the simplest Vulkan application can be a non-trivial task. To help with this, `utils_VkFFT` contains the routines that can help to create the simplest Vulkan application, allocate memory, record command buffers and launch them. Code has some comments explaining what is going on at each step. It also has some useful struct defines (like `vkGPU`) that keep the most important handles used in Vulkan Compute. This section may be expanded in the future to the proper step-by-step guide on Vulkan Compute simple application creation. I also encourage to check <https://github.com/DTolm/VulkanComputeSamples-Transposition> repository for another example of a compute algorithm (matrix transposition) implemented with Vulkan API.

`utils_VkFFT` also has a routine that prints the list of available devices.

`vkGPU` struct has the following definition:

```

typedef struct {
    #if(VKFFT_BACKEND==0) //Vulkan API
    VkInstance instance; //a connection between the application and
        ↪ the Vulkan library
    VkPhysicalDevice physicalDevice; //a handle for the graphics
        ↪ card used in the application
    VkPhysicalDeviceProperties physicalDeviceProperties; //basic
        ↪ device properties
    VkPhysicalDeviceMemoryProperties
        ↪ physicalDeviceMemoryProperties; //basic memory properties
        ↪ of the device
    VkDevice device; //a logical device, interacting with physical
        ↪ device
    VkDebugUtilsMessengerEXT debugMessenger; //extension for
        ↪ debugging
    uint64_t queueFamilyIndex; //if multiple queues are available,
        ↪ specify the used one
    VkQueue queue; //a place, where all operations are submitted
    VkCommandPool commandPool; //an opaque objects that command
        ↪ buffer memory is allocated from
    VkFence fence; //a vkGPU->fence used to synchronize dispatches
    std::vector<const char*> enabledDeviceExtensions;
    uint64_t enableValidationLayers;
    #elif(VKFFT_BACKEND==1) //CUDA API
    CUdevice device;
    CUcontext context;
    #elif(VKFFT_BACKEND==2) //HIP API
    hipDevice_t device;
    hipCtx_t context;
    #elif(VKFFT_BACKEND==3) //OpenCL API
    cl_platform_id platform;
    cl_device_id device;
    cl_context context;
    cl_command_queue commandQueue;
    #elif(VKFFT_BACKEND==4) //Level Zero API
    ze_driver_handle_t driver;
    ze_device_handle_t device;
    ze_context_handle_t context;
    ze_command_queue_handle_t commandQueue;
    uint32_t commandQueueID;
    #elif(VKFFT_BACKEND==5) //Metal API
    MTL::Device* device;
    MTL::CommandQueue* queue;

```

```
#endif
uint64_t device_id; //an id of a device, reported by
    ↪ devices_list call
} VkGPU;
```

6 VkFFT Code Examples

This section will provide some simple pseudocode for VkFFT usage, which will once again outline important steps required to launch FFT with VkFFT. More information (and fully working code) can be found in this folder of the VkFFT repository:

/benchmark_samples/vkFFT_scripts/src/

6.1 Driver initializations

Before launching VkFFT, do not forget to do all necessary driver initializations. The following code specifies them for all the supported backends, though the final implementation may be different depending on the particular user's configuration.

```
/*if(VKFFT_BACKEND==0) //Vulkan API*/
VkResult res = VK_SUCCESS;
//create instance - a connection between the application and
↳ the Vulkan library
res = createInstance(vkGPU, sample_id);
if (res != 0) {
    //printf("Instance creation failed, error code: %" PRIu64
    ↳ "\n", res);
    return VKFFT_ERROR_FAILED_TO_CREATE_INSTANCE;
}
//set up the debugging messenger
res = setupDebugMessenger(vkGPU);
if (res != 0) {
    //printf("Debug messenger creation failed, error code: %"
    ↳ PRIu64 "\n", res);
    return VKFFT_ERROR_FAILED_TO_SETUP_DEBUG_MESSENGER;
}
//check if there are GPUs that support Vulkan and select one
res = findPhysicalDevice(vkGPU);
if (res != 0) {
    //printf("Physical device not found, error code: %" PRIu64
    ↳ "\n", res);
    return VKFFT_ERROR_FAILED_TO_FIND_PHYSICAL_DEVICE;
}
//create logical device representation
res = createDevice(vkGPU, sample_id);
if (res != 0) {
    //printf("Device creation failed, error code: %" PRIu64 "\n",
    ↳ res);
    return VKFFT_ERROR_FAILED_TO_CREATE_DEVICE;
}
```

```

//create fence for synchronization
res = createFence(vkGPU);
if (res != 0) {
    //printf("Fence creation failed, error code: %" PRIu64 "\n",
    ↪ res);
    return VKFFT_ERROR_FAILED_TO_CREATE_FENCE;
}
//create a place, command buffer memory is allocated from
res = createCommandPool(vkGPU);
if (res != 0) {
    //printf("Fence creation failed, error code: %" PRIu64
    ↪ "\n", res);
    return VKFFT_ERROR_FAILED_TO_CREATE_COMMAND_POOL;
}
vkGetPhysicalDeviceProperties(vkGPU->physicalDevice,
    ↪ &vkGPU->physicalDeviceProperties);
vkGetPhysicalDeviceMemoryProperties(vkGPU->physicalDevice,
    ↪ &vkGPU->physicalDeviceMemoryProperties);
glslang_initialize_process();
//compiler can be initialized before VkFFT

#elif(VKFFT_BACKEND==1) //CUDA API
CUresult res = CUDA_SUCCESS;
cudaError_t res2 = cudaSuccess;
res = cuInit(0);
if (res != CUDA_SUCCESS) return
    ↪ VKFFT_ERROR_FAILED_TO_INITIALIZE;
res2 = cudaSetDevice((int)vkGPU->device_id);
if (res2 != cudaSuccess) return
    ↪ VKFFT_ERROR_FAILED_TO_SET_DEVICE_ID;
res = cuDeviceGet(&vkGPU->device, (int)vkGPU->device_id);
if (res != CUDA_SUCCESS) return
    ↪ VKFFT_ERROR_FAILED_TO_GET_DEVICE;
res = cuCtxCreate(&vkGPU->context, 0, (int)vkGPU->device);
if (res != CUDA_SUCCESS) return
    ↪ VKFFT_ERROR_FAILED_TO_CREATE_CONTEXT;
#elif(VKFFT_BACKEND==2) //HIP API
hipError_t res = hipSuccess;
res = hipInit(0);
if (res != hipSuccess) return VKFFT_ERROR_FAILED_TO_INITIALIZE;
res = hipSetDevice((int)vkGPU->device_id);
if (res != hipSuccess) return
    ↪ VKFFT_ERROR_FAILED_TO_SET_DEVICE_ID;
res = hipDeviceGet(&vkGPU->device, (int)vkGPU->device_id);

```

```

if (res != hipSuccess) return VKFFT_ERROR_FAILED_TO_GET_DEVICE;
res = hipCtxCreate(&vkGPU->context, 0, (int)vkGPU->device);
if (res != hipSuccess) return
    ↪ VKFFT_ERROR_FAILED_TO_CREATE_CONTEXT;
#elif(VKFFT_BACKEND==3) //OpenCL API
cl_int res = CL_SUCCESS;
cl_uint numPlatforms;
res = clGetPlatformIDs(0, 0, &numPlatforms);
if (res != CL_SUCCESS) return VKFFT_ERROR_FAILED_TO_INITIALIZE;
cl_platform_id* platforms =
    ↪ (cl_platform_id*)malloc(sizeof(cl_platform_id) *
    ↪ numPlatforms);
if (!platforms) return VKFFT_ERROR_MALLOC_FAILED;
res = clGetPlatformIDs(numPlatforms, platforms, 0);
if (res != CL_SUCCESS) return VKFFT_ERROR_FAILED_TO_INITIALIZE;
uint64_t k = 0;
for (uint64_t j = 0; j < numPlatforms; j++) {
    cl_uint numDevices;
    res = clGetDeviceIDs(platforms[j], CL_DEVICE_TYPE_ALL, 0,
        ↪ 0, &numDevices);
    cl_device_id* deviceList =
        ↪ (cl_device_id*)malloc(sizeof(cl_device_id) *
        ↪ numDevices);
    if (!deviceList) return VKFFT_ERROR_MALLOC_FAILED;
    res = clGetDeviceIDs(platforms[j], CL_DEVICE_TYPE_ALL,
        ↪ numDevices, deviceList, 0);
    if (res != CL_SUCCESS) return
        ↪ VKFFT_ERROR_FAILED_TO_GET_DEVICE;
    for (uint64_t i = 0; i < numDevices; i++) {
        if (k == vkGPU->device_id) {
            vkGPU->platform = platforms[j];
            vkGPU->device = deviceList[i];
            vkGPU->context = clCreateContext(NULL, 1,
                ↪ &vkGPU->device, NULL, NULL, &res);
            if (res != CL_SUCCESS) return
                ↪ VKFFT_ERROR_FAILED_TO_CREATE_CONTEXT;
            cl_command_queue commandQueue =
                ↪ clCreateCommandQueue(vkGPU->context,
                ↪ vkGPU->device, 0, &res);
            if (res != CL_SUCCESS) return
                ↪ VKFFT_ERROR_FAILED_TO_CREATE_COMMAND_QUEUE;
            vkGPU->commandQueue = commandQueue;
            i=numDevices;
            j=numPlatforms;
        }
    }
}

```

```

        }
        else {
            k++;
        }
    }
    free(deviceList);
}
free(platforms);
#elif(VKFFT_BACKEND==4)
    ze_result_t res = ZE_RESULT_SUCCESS;
    res = zeInit(0);
    if (res != ZE_RESULT_SUCCESS) return
        ↪ VKFFT_ERROR_FAILED_TO_INITIALIZE;
    uint32_t numDrivers = 0;
    res = zeDriverGet(&numDrivers, 0);
    if (res != ZE_RESULT_SUCCESS) return
        ↪ VKFFT_ERROR_FAILED_TO_INITIALIZE;
    ze_driver_handle_t* drivers =
        ↪ (ze_driver_handle_t*)malloc(numDrivers *
        ↪ sizeof(ze_driver_handle_t));
    if (!drivers) return VKFFT_ERROR_MALLOC_FAILED;
    res = zeDriverGet(&numDrivers, drivers);
    if (res != ZE_RESULT_SUCCESS) return
        ↪ VKFFT_ERROR_FAILED_TO_INITIALIZE;
    uint64_t k = 0;
    for (uint64_t j = 0; j < numDrivers; j++) {
        uint32_t numDevices = 0;
        res = zeDeviceGet(drivers[j], &numDevices, nullptr);
        if (res != ZE_RESULT_SUCCESS) return
            ↪ VKFFT_ERROR_FAILED_TO_GET_DEVICE;
        ze_device_handle_t* deviceList =
            ↪ (ze_device_handle_t*)malloc(numDevices *
            ↪ sizeof(ze_device_handle_t));
        if (!deviceList) return VKFFT_ERROR_MALLOC_FAILED;
        res = zeDeviceGet(drivers[j], &numDevices, deviceList);
        if (res != ZE_RESULT_SUCCESS) return
            ↪ VKFFT_ERROR_FAILED_TO_GET_DEVICE;
        for (uint64_t i = 0; i < numDevices; i++) {
            if (k == vkGPU->device_id) {
                vkGPU->driver = drivers[j];
                vkGPU->device = deviceList[i];
                ze_context_desc_t contextDescription = {};
                contextDescription.stype =
                    ↪ ZE_STRUCTURE_TYPE_CONTEXT_DESC;
            }
        }
    }
}

```



```

res = zeContextCreate(vkGPU->driver,
    ↪ &contextDescription, &vkGPU->context);
if (res != ZE_RESULT_SUCCESS) return
    ↪ VKFFT_ERROR_FAILED_TO_CREATE_CONTEXT;

uint32_t queueGroupCount = 0;
res = zeDeviceGetCommandQueueGroupProperties(
    ↪ vkGPU->device, &queueGroupCount,
    ↪ 0);
if (res != ZE_RESULT_SUCCESS) return
    ↪ VKFFT_ERROR_FAILED_TO_CREATE_COMMAND_QUEUE;

ze_command_queue_group_properties_t*
    ↪ cmdqueueGroupProperties =
    ↪ (ze_command_queue_group_properties_t*)
    ↪ malloc(queueGroupCount *
    ↪ sizeof(ze_command_queue_group_properties_t));
if (!cmdqueueGroupProperties) return
    ↪ VKFFT_ERROR_MALLOC_FAILED;
res = zeDeviceGetCommandQueueGroupProperties(
    ↪ vkGPU->device, &queueGroupCount,
    ↪ cmdqueueGroupProperties);
if (res != ZE_RESULT_SUCCESS) return
    ↪ VKFFT_ERROR_FAILED_TO_CREATE_COMMAND_QUEUE;

uint32_t commandQueueID = -1;
for (uint32_t i = 0; i < queueGroupCount; ++i)
    ↪ {
        if ((cmdqueueGroupProperties[i].flags &&
            ↪ ZE_COMMAND_QUEUE_GROUP_PROPERTY_FLAG_COMPUTE)
            ↪ && (cmdqueueGroupProperties[i].flags
            ↪ &&
            ↪ ZE_COMMAND_QUEUE_GROUP_PROPERTY_FLAG_COPY))
            ↪ {
                commandQueueID = i;
                break;
            }
    }
if (commandQueueID == -1) return
    ↪ VKFFT_ERROR_FAILED_TO_CREATE_COMMAND_QUEUE;
vkGPU->commandQueueID = commandQueueID;
ze_command_queue_desc_t
    ↪ commandQueueDescription = {};

```

```

        commandQueueDescription.stype =
            ↪ ZE_STRUCTURE_TYPE_COMMAND_QUEUE_DESC;
        commandQueueDescription.ordinal =
            ↪ commandQueueID;
        commandQueueDescription.priority =
            ↪ ZE_COMMAND_QUEUE_PRIORITY_NORMAL;
        commandQueueDescription.mode =
            ↪ ZE_COMMAND_QUEUE_MODE_DEFAULT;
        res = zeCommandQueueCreate(vkGPU->context,
            ↪ vkGPU->device, &commandQueueDescription,
            ↪ &vkGPU->commandQueue);
        if (res != ZE_RESULT_SUCCESS) return
            ↪ VKFFT_ERROR_FAILED_TO_CREATE_COMMAND_QUEUE;
        free(cmdqueueGroupProperties);
        i=numDevices;
        j=numPlatforms;
    }
    else {
        k++;
    }
}

free(deviceList);
}
free(drivers);
#elif(VKFFT_BACKEND==5)
NS::Array* devices = MTL::CopyAllDevices();
MTL::Device* device =
    ↪ (MTL::Device*)devices->object(vkGPU->device_id);
vkGPU->device = device;
MTL::CommandQueue* queue = device->newCommandQueue();
vkGPU->queue = queue;
#endif

```

6.2 Simple FFT application example: 1D (one dimensional) C2C (complex to complex) FP32 (single precision) FFT

This example performs the simplest case of FFT. It shows all the necessary fields that the user must fill during the configuration and the submission process. Other samples will build on this one, as driver parameters initialization and code execution commands are the same for all configurations (except for the launch parameters that can be configured after application creation).

```

//zero-initialize configuration + FFT application
VkFFTConfiguration configuration = {};
VkFFTApplication app = {};

configuration.FFTdim = 1; //FFT dimension
configuration.size[0] = Nx; //FFT size
uint64_t bufferSize = (uint64_t)sizeof(float) * 2 *
    ↪ configuration.size[0];

//Device management + code submission
#if(VKFFT_BACKEND==5)
configuration.device = vkGPU->device;
#else
configuration.device = &vkGPU->device;
#endif

#if(VKFFT_BACKEND==0) //Vulkan API
configuration.queue = &vkGPU->queue;
configuration.fence = &vkGPU->fence;
configuration.commandPool = &vkGPU->commandPool;
configuration.physicalDevice = &vkGPU->physicalDevice;
    ↪
configuration.isCompilerInitialized = isCompilerInitialized;
    ↪ //glslang compiler can be initialized before VkFFT plan
    ↪ creation. if not, VkFFT will create and destroy one after
    ↪ initialization
#elif(VKFFT_BACKEND==3) //OpenCL API
configuration.context = &vkGPU->context;
#elif(VKFFT_BACKEND==4)
configuration.context = &vkGPU->context;
configuration.commandQueue = &vkGPU->commandQueue;
configuration.commandQueueID = vkGPU->commandQueueID;
#elif(VKFFT_BACKEND==5)
configuration.queue = vkGPU->queue;
#endif

allocateBuffer(buffer, bufferSize); //Pseudocode for buffer
    ↪ allocation, differs between APIs
transferDataFromCPU(buffer, cpu_buffer); //Pseudocode for data
    ↪ transfer from CPU to GPU, differs between APIs

#if(VKFFT_BACKEND==0) //Vulkan API needs bufferSize at
    ↪ initialization

```

```

configuration.bufferSize = &bufferSize;
#endif

VkFFTResult resFFT = initializeVkFFT(&app, configuration);

VkFFTLaunchParams launchParams = {};
launchParams.buffer = &buffer;
#ifdef VKFFT_BACKEND==0) //Vulkan API
launchParams.commandBuffer = &commandBuffer;
#elif(VKFFT_BACKEND==3) //OpenCL API
launchParams.commandQueue = &commandQueue;
#elif(VKFFT_BACKEND==4) //Level Zero API
launchParams->commandList = &commandList;
#elif(VKFFT_BACKEND==5) //Metal API
launchParams->commandBuffer = commandBuffer;
launchParams->commandEncoder = commandEncoder;
#endif
resFFT = VkFFTAppend(app, -1, &launchParams);

//add synchronization relevant to your API - vkWaitFor-
↳ Fences/cudaDeviceSynchronize/hipDeviceSynchronize/clFinish
transferDataToCPU(cpu_buffer, buffer); //Pseudocode for data
↳ transfer from GPU to CPU, differs between APIs

freeBuffer(buffer, bufferSize); //Pseudocode for buffer
↳ deallocation, differs between APIs

deleteVkFFT(&app);

```

6.3 Advanced FFT application example: ND, C2C/R2C/R2R, different precisions, batched FFT

This example shows how to configure the main parameters of interest in the VkFFT library: multidimensional case, different types of transforms, different precision, perform batched transforms.

In the code below X, Y and Z are the dimensions of FFT, B - number of batches, R2C - real to complex mode 0 or 1 (on/off), DCT - 0, 1, 2, 3 or 4 (off/DCT type), P - precision (0 - single, 1 - double, 2 - half).

```

//zero-initialize configuration + FFT application
VkFFTConfiguration configuration = {};
VkFFTApplication app = {};

configuration.FFTdim = 1; //FFT dimension

```

```

configuration.size[0] = X;
configuration.size[1] = Y;
configuration.size[2] = Z;
if (Y > 1) configuration.FFTdim++;
if (Z > 1) configuration.FFTdim++;
configuration.numberBatches = B;
configuration.performR2C = R2C;
configuration.performDCT = DCT;
if (P == 1) configuration.doublePrecision = 1;
if (P == 2) configuration.halfPrecision = 1;
if (P == 3) configuration.quadDoubleDoublePrecision = 1;

uint64_t bufferSize = 0;

if (R2C) {
    bufferSize = (uint64_t)(storageComplexSize / 2) *
        ↪ (configuration.size[0] + 2) * configuration.size[1] *
        ↪ configuration.size[2] * configuration.numberBatches;
}
else {
    if (DCT) {
        bufferSize = (uint64_t)(storageComplexSize / 2) *
            ↪ configuration.size[0] * configuration.size[1] *
            ↪ configuration.size[2] *
            ↪ configuration.numberBatches;
    }
    else {
        bufferSize = (uint64_t)storageComplexSize *
            ↪ configuration.size[0] * configuration.size[1] *
            ↪ configuration.size[2] *
            ↪ configuration.numberBatches;
    }
} // storageComplexSize - 4/8/16 for FP16/FP32/FP64
  ↪ respectively.

//Device management + code submission - code is identical to
  ↪ the previous example

```

6.4 Advanced FFT application example: out-of-place R2C FFT with custom strides

In this example, VkFFT is configured to calculate a 3D out-of-place R2C FFT of a system with custom strides. VkFFT reads data from the inputBuffer and produces the result in the buffer.

```

//zero-initialize configuration + FFT application
VkFFTConfiguration configuration = {};
VkFFTApplication app = {};

configuration.FFTdim = 3; //FFT dimension
configuration.size[0] = Nx;
configuration.size[1] = Ny;
configuration.size[2] = Nz;

configuration.performR2C = 1;

//out-of-place - we need to specify that input buffer is
↳ separate from the main buffer
configuration.isInputFormatted = 1;
configuration.inputBufferStride[0] = configuration.size[0];
configuration.inputBufferStride[1] =
↳ configuration.inputBufferStride[0] *
↳ configuration.size[1];
configuration.inputBufferStride[2] =
↳ configuration.inputBufferStride[1] *
↳ configuration.size[2];

configuration.bufferStride[0] = (uint64_t)
↳ (configuration.size[0] / 2) + 1;
configuration.bufferStride[1] = configuration.bufferStride[0]
↳ * configuration.size[1];
configuration.bufferStride[2] = configuration.bufferStride[1]*
↳ configuration.size[2];

uint64_t inputBufferSize = (uint64_t)sizeof(float) *
↳ configuration.size[0] * configuration.size[1] *
↳ configuration.size[2];

uint64_t bufferSize = (uint64_t)sizeof(float) * 2 *
↳ (configuration.size[0]/2+1) * configuration.size[1] *
↳ configuration.size[2];

//Device management + code submission - code is identical to
↳ the first example, except that you need to allocate two
↳ buffers (and provide them in the launch configuration).

```

6.5 Advanced FFT application example: 3D FFT with innermost batching

In this example, VkFFT is configured to calculate a 3D in-place C2C FFT of a system with innermost (fast axis) batching.

```
//zero-initialize configuration + FFT application
VkFFTConfiguration configuration = {};
VkFFTApplication app = {};

configuration.FFTdim = 4; //FFT dimension, 3+1
configuration.size[0] = Nbatch;
configuration.size[1] = Nx;
configuration.size[2] = Ny;
configuration.size[3] = Nz;

configuration.omitDimension[0] = 1; // disable FFT over the
    ↪ batching dimension

uint64_t bufferSize = (uint64_t)sizeof(float) *
    ↪ configuration.size[0] * configuration.size[1] *
    ↪ configuration.size[2] * configuration.size[3];

//Device management + code submission - code is identical to
    ↪ the first example, except that you need to allocate two
    ↪ buffers (and provide them in the launch configuration).
```

6.6 Advanced FFT application example: 3D zero-padded FFT

In this example, VkFFT is configured to calculate a 3D FFT of a system. The meaningful data is located in the first octant of the buffer, the rest is padded with zeros. This configuration removes the circular part of the convolution and allows modelling of open systems.

```
//zero-initialize configuration + FFT application
VkFFTConfiguration configuration = {};
VkFFTApplication app = {};

configuration.FFTdim = 3; //FFT dimension
configuration.size[0] = Nx;
configuration.size[1] = Ny;
configuration.size[2] = Nz;
```

```

configuration.performZeropadding[0] = 1; //Perform padding
↳ with zeros on GPU. Still need to properly align input data
↳ (no need to fill padding area with meaningful data) but
↳ this will increase performance due to the lower amount of
↳ the memory reads/writes and omitting sequences only
↳ consisting of zeros.
configuration.performZeropadding[1] = 1;
configuration.performZeropadding[2] = 1;
configuration.fft_zeropad_left[0] =
↳ (uint64_t)ceil(configuration.size[0] / 2.0);
configuration.fft_zeropad_right[0] = configuration.size[0];
configuration.fft_zeropad_left[1] =
↳ (uint64_t)ceil(configuration.size[1] / 2.0);
configuration.fft_zeropad_right[1] = configuration.size[1];
configuration.fft_zeropad_left[2] =
↳ (uint64_t)ceil(configuration.size[2] / 2.0);
configuration.fft_zeropad_right[2] = configuration.size[2];

uint64_t bufferSize = (uint64_t)storageComplexSize *
↳ configuration.size[0] * configuration.size[1] *
↳ configuration.size[2];

//Device management + code submission - code is identical to
↳ the first example

```

6.7 Convolution application example: 3x3 matrix-vector convolution in 1D

In this example, VkFFT is configured to calculate a kernel, represented by a 3x3 matrix and a system, represented by a 3D vector. Their convolution is a matrix-vector multiplication in the frequency domain.

```

//zero-initialize configuration + FFT application, we need two
↳ - one for kernel calculation
VkFFTConfiguration kernel_configuration = {};
VkFFTConfiguration convolution_configuration = {};
VkFFTApplication app_kernel = {};
VkFFTApplication app_convolution = {};

kernel_configuration.FFTdim = 1; //FFT dimension
kernel_configuration.size[0] = Nx; //FFT size

```



```

uint64_t bufferSize = (uint64_t)sizeof(float) * 2 *
    ↪ kernel_configuration.size[0];

//configure kernel
kernel_configuration.kernelConvolution = 1; //specify if this
    ↪ plan is used to create kernel for convolution
kernel_configuration.coordinateFeatures = 9; //Specify
    ↪ dimensionality of the input feature vector (default 1).
    ↪ Each component is stored not as a vector, but as a separate
    ↪ system and padded on it's own according to other options
    ↪ (i.e. for x*y system of 3-vector, first x*y elements
    ↪ correspond to the first dimension, then goes x*y for the
    ↪ second, etc).
//coordinateFeatures number is an important constant for
    ↪ convolution. If we perform 1x1 convolution, it is equal to
    ↪ number of features, but matrixConvolution should be equal
    ↪ to 1. For matrix convolution, it must be equal to
    ↪ matrixConvolution parameter. If we perform 2x2 convolution,
    ↪ it is equal to 3 for symmetric kernel (stored as xx, xy,
    ↪ yy) and 4 for nonsymmetric (stored as xx, xy, yx, yy).
    ↪ Similarly, 6 (stored as xx, xy, xz, yy, yz, zz) and 9
    ↪ (stored as xx, xy, xz, yx, yy, yz, zx, zy, zz) for 3x3
    ↪ convolutions.
kernel_configuration.normalize = 1;

//Initialize app_kernel and perform a single forward FFT like
    ↪ in examples before. You pass kernel as a buffer for the
    ↪ preparation stage.

convolution_configuration = kernel_configuration;
convolution_configuration.kernelConvolution = 0;
convolution_configuration.performConvolution = 1;
convolution_configuration.symmetricKernel = 0; //Specify if
    ↪ convolution kernel is symmetric. In this case we only pass
    ↪ upper triangle part of it in the form of: (xx, xy, yy) for
    ↪ 2d and (xx, xy, xz, yy, yz, zz) for 3d.
convolution_configuration.matrixConvolution = 3; //we do matrix
    ↪ convolution, so kernel is 9 numbers (3x3), but vector
    ↪ dimension is 3
convolution_configuration.coordinateFeatures = 3; //equal to
    ↪ matrixConvolution size

```

```
//Initialize app_convolution and perform a single forward FFT
↪ like in examples before. You pass kernel as kernel and
↪ system to be convolved with it as buffer
```

6.8 Convolution application example: R2C cross-correlation between two sets of N images

In this example, VkFFT is configured to calculate a kernel, represented by three 2D vectors (RGB values of a pixel) and a system, also represented by three 2D vectors. There are N kernels and N systems. Their cross-correlation is a conjugate convolution in the frequency domain. Images are usually stored as real, not complex numbers, so code uses R2C optimization as well.

```
//zero-initialize configuration + FFT application, we need two
↪ - one for kernel calculation
VkFFTConfiguration kernel_configuration = {};
VkFFTConfiguration convolution_configuration = {};
VkFFTApplication app_kernel = {};
VkFFTApplication app_convolution = {};

kernel_configuration.FFTdim = 2; //FFT dimension
kernel_configuration.size[0] = Nx;
kernel_configuration.size[1] = Ny;
kernel_configuration.coordinateFeatures = 3;
kernel_configuration.numberBatches = N;
kernel_configuration.performR2C = 1;
kernel_configuration.normalize = 1;

uint64_t bufferSize = (uint64_t)sizeof(float) * 2 *
↪ (kernel_configuration.size[0]/2+1) *
↪ kernel_configuration.size[1] *
↪ kernel_configuration.coordinateFeatures *
↪ kernel_configuration.numberBatches;

kernel_configuration.kernelConvolution = 1; //specify if this
↪ plan is used to create kernel for convolution

//Initialize app_kernel and perform a single forward FFT like
↪ in examples before. Pad in-place R2C system like this:

for (uint64_t n = 0; n < kernel_configuration.numberBatches;
↪ n++) {
```

```

for (uint64_t c = 0; c <
    ↪ kernel_configuration.coordinateFeatures; c++) {
    for (uint64_t j = 0; j < kernel_configuration.size[1];
        ↪ j++) {
        for (uint64_t i = 0; i <
            ↪ kernel_configuration.size[0]; i++) {
            kernel_padded_GPU[i + j * 2 *
                ↪ (kernel_configuration.size[0]/2 + 1) + c *
                ↪ 2 * (kernel_configuration.size[0]/2 + 1) *
                ↪ kernel_configuration.size[1] + n * 2 *
                ↪ (kernel_configuration.size[0]/2 + 1) *
                ↪ kernel_configuration.size[1] *
                ↪ kernel_configuration.coordinateFeatures] =
                ↪ kernel_input[i + j *
                ↪ kernel_configuration.size[0] + c *
                ↪ kernel_configuration.size[0] *
                ↪ kernel_configuration.size[1] + n *
                ↪ kernel_configuration.size[0] *
                ↪ kernel_configuration.size[1] *
                ↪ kernel_configuration.coordinateFeatures];
        }
    }
}

convolution_configuration = kernel_configuration;
convolution_configuration.kernelConvolution = 0;
convolution_configuration.performConvolution = 1;
convolution_configuration.conjugateConvolution = 1;

//Initialize app_convolution and perform a single forward FFT
    ↪ like in examples before. Pad the system in the same way as
    ↪ the kernel

```

6.9 Simple FFT application binary reuse application

This example shows how to save/load binaries generated by VkFFT. This can reduce time taken by initializeVkFFT call by removing RTC components from it. Be sure that rest of the configuration stays the same to reuse the binary. Use rb/wb flags to load/save. This does not currently work in Metal backend.

```

VkFFTConfiguration configuration = {};
VkFFTApplication app = {};

//configuration is initialized like in other examples

```

```

configuration.saveApplicationToString = 1;
//configuration.loadApplicationFromString = 1; //choose one to
↪ save/load binary file

if (configuration.loadApplicationFromString) {
    FILE* kernelCache;
    uint64_t str_len;
    kernelCache = fopen("VkFFT_binary", "rb");
    fseek(kernelCache, 0, SEEK_END);
    str_len = ftell(kernelCache);
    fseek(kernelCache, 0, SEEK_SET);
    configuration.loadApplicationString = malloc(str_len);
    fread(configuration.loadApplicationString, str_len, 1,
        ↪ kernelCache);
    fclose(kernelCache);
}

resFFT = initializeVkFFT(&app, configuration);
if (resFFT != VKFFT_SUCCESS) return resFFT;

if (configuration.loadApplicationFromString)
    free(configuration.loadApplicationString);

if (configuration.saveApplicationToString) {
    FILE* kernelCache;
    kernelCache = fopen("VkFFT_binary", "wb");
    fwrite(app.saveApplicationString,
        ↪ app.applicationStringSize, 1, kernelCache);
    fclose(kernelCache);
}

//application is launched like in other examples

```