

**Universidade do Minho**

Escola de Engenharia

Departamento de Informática

Jorge Francisco Teixeira Bastos da Mota

## **High Performance Fourier Transforms on GPUs**

December 2022



**Universidade do Minho**

Escola de Engenharia

Departamento de Informática

Jorge Francisco Teixeira Bastos da Mota

## **High Performance Fourier Transforms on GPUs**

Master dissertation

Integrated Master's in Informatics Engineering

Dissertation supervised by

**António Ramires**

December 2022

---

## COPYRIGHT AND TERMS OF USE FOR THIRD PARTY WORK

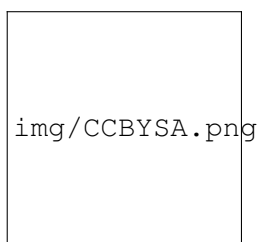
---

This dissertation reports on academic work that can be used by third parties as long as the internationally accepted standards and good practices are respected concerning copyright and related rights.

This work can thereafter be used under the terms established in the license below.

Readers needing authorization conditions not provided for in the indicated licensing should contact the author through the RepositóriUM of the University of Minho.

LICENSE GRANTED TO USERS OF THIS WORK:



**CC BY-SA**

<https://creativecommons.org/licenses/by-sa/4.0/>

---

## STATEMENT OF INTEGRITY

---

I hereby declare having conducted this academic work with integrity.

I confirm that I have not used plagiarism or any form of undue use of information or falsification of results along the process leading to its elaboration.

I further declare that I have fully acknowledged the Code of Ethical Conduct of the University of Minho.

---

## ABSTRACT

---

The Fast Fourier Transform is an algorithm or a family of algorithms indispensable for the computation of the Discrete Fourier Transform. Accordingly, these transforms are the core of many applications in several areas and are required to be computed efficiently in many scenarios.

The continuous evolution of GPUs has increased the popularity of parallelizable algorithm implementations on this type of hardware. Traditionally GPUs were associated to graphics background, however, with the popularization of the compute functionality of this hardware, most modern GPUs now have this capability, hence, algorithms now are more likely to be implemented in the general purpose compute pipeline of GPUs. As a result, many applications take advantage of compute programming in GPGPU capable frameworks such as GLSL, a high-level shading language recurrently used in the context of computer graphics.

In this dissertation we provide, refine and analyze GPU driven implementations of the family of FFT algorithms in GLSL, with the goal to provide programmers with efficient and simplified compute kernels for this transform, from the classic Cooley-Tukey algorithm to more suitable algorithms for the GPU. Accordingly, we also compare these same algorithms with different GPGPU frameworks with the goal to analyse their significance for different compute APIs. Finally, we demonstrate how all improvements discussed in this dissertation culminate in the performance improvement in a realtime rendering technique that heavily depends on this transform, as a case of study.

**KEYWORDS**     FFT, GPGPU, GLSL, cuFFT, analysis, performance, compute.

---

## RESUMO

---

A Transformada Rápida de Fourier é um algoritmo ou uma família de algoritmos indispensáveis para o cálculo da Transformada Discreta de Fourier. Assim, essas transformadas são o núcleo de muitas aplicações em diversas áreas e precisam ser calculadas de forma eficiente em muitos cenários.

A evolução contínua dos GPUs aumentou a popularidade das implementações de algoritmos paralelizáveis neste tipo de *hardware*. Tradicionalmente, os GPUs eram associadas ao fundo gráfico, no entanto, com a popularização da funcionalidade de *compute* desse hardware, os GPUs mais modernos agora têm essa capacidade, portanto, os algoritmos agora são mais propensos a serem implementados na *compute pipeline* de propósito geral dos GPUs. Como resultado, muitas aplicações aproveitam a programação em *compute* em *frameworks* compatíveis com GPGPU como GLSL, uma linguagem de *shading* de alto nível usada recorrentemente no contexto de computação gráfica.

Nesta dissertação fornecemos, refinamos e analisamos implementações em GPU da família de algoritmos FFT em GLSL, com o objetivo de fornecer aos programadores *compute kernels* eficientes e simplificados para esta transformada, desde o clássico algoritmo de Cooley-Tukey até algoritmos mais adequados para o GPU. Da mesma forma, também comparamos esses mesmos algoritmos com diferentes *frameworks* GPGPU com o objetivo de analisar a sua significância para diferentes APIs de *compute*. Por fim, demonstramos como todas as melhorias discutidas nesta dissertação culminam na melhoria de desempenho em uma técnica de renderização em tempo real que depende fortemente dessa transformação, como caso de estudo.

**PALAVRAS-CHAVE**     FFT, GPGPU, GLSL, cuFFT, análise, performance, compute.

---

## CONTENTS

---

<b>1</b>	<b>Introduction</b>	<b>5</b>
1.1	Contextualization	5
1.2	Motivation	5
1.3	Objectives	5
1.4	Document Organization	6
<b>2</b>	<b>The Fourier Transform</b>	<b>7</b>
2.1	Continuous Fourier Transform	8
2.2	Discrete Fourier Transform	9
2.2.1	Matrix multiplication	10
2.3	Fast Fourier Transform	12
2.3.1	Radix-2 Decimation-in-Time FFT	12
2.3.2	Radix-2 Decimation-in-Frequency FFT	16
<b>3</b>	<b>Algorithms analysis</b>	<b>19</b>
<b>4</b>	<b>Computation of the Fourier Transform</b>	<b>20</b>
4.1	Natural order Cooley-Tukey	20
4.2	Stockham algorithm	21
4.3	Radix-4 instead of Radix-2	23
<b>5</b>	<b>Implementation on the GPU</b>	<b>26</b>
5.1	GPU Programming model	26
5.2	2D Fourier Transform on the GPU	28
5.3	GLSL implementation	29
5.3.1	Cooley-Tukey	29
5.3.2	Radix-2 Stockham	34
5.3.3	Radix-4 Stockham	35
<b>6</b>	<b>Analysis and Comparison</b>	<b>38</b>
6.1	cuFFT	38
6.2	Implementation analysis in GLSL	39
6.3	Implementation analysis in CUDA	41

6.4	Case of study	45
6.4.1	Tensendorf waves	45
6.4.2	Results	47
<b>7</b>	<b>Conclusions and Future work</b>	<b>49</b>
	<b>Bibliography</b>	<b>50</b>
<b>I</b>	<b>Appendices</b>	
<b>A</b>	<b>GLSL FFT</b>	<b>53</b>
<b>B</b>	<b>cuFFT</b>	<b>67</b>
<b>C</b>	<b>CUDA FFT</b>	<b>69</b>



---

## LIST OF FIGURES

---

Figure 1	Time to frequency signal decomposition <b>Source:</b> NTiAudio	7
Figure 2	Radix-2 Decimation-in-Time FFT <b>Source:</b> Jones (2014)	13
Figure 3	Cooley-Tukey butterfly	13
Figure 4	Bit reverse permutation	14
Figure 5	Radix-2 Decimation-in-Frequency FFT <b>Source:</b> Jones (2014)	16
Figure 6	Gentleman-Sande butterfly	17
Figure 7	Chain of even and odd compositions over each stage for a natural order DIF	21
Figure 8	Illustration of a stage in radix-4 Stockham with each color representing a radix-4 butterfly computation	23
Figure 9	Radix-4 FFT butterfly structure <b>Source:</b> Marti-Puig and Bolano (2009)	24
Figure 10	CPU architecture compared with a GPU architecture	27
Figure 11	High level illustration of horizontal and vertical passes	28
Figure 12	Difference in invocation spaces for size 256 FFT to allow barrier synchronization between local threads	33
Figure 13	Frame time difference between using stage per pass approach and unique pass for the Radix-2 Cooley-Tukey algorithm	40
Figure 14	2D FFT benchmarks in milliseconds of out-of-place cuFFT and GLSL Radix-2 algorithms	42
Figure 15	2D FFT benchmarks of in milliseconds out-of-place cuFFT and GLSL Radix-4 Stockham algorithm	42
Figure 16	2D FFT benchmarks in milliseconds of GLSL Radix-2 Stockham with multiple butterflies per local threads	43
Figure 17	2D FFT benchmarks in milliseconds of CUDA and GLSL Radix-2 algorithms	44
Figure 18	2D FFT benchmarks in milliseconds of CUDA and GLSL Radix-4 Stockham algorithm	44
Figure 19	Tensendorf waves in Nau3D Engine	46
Figure 20	Time spent in the CPU and GPU for the size 512 FFT horizontal and vertical passes	46
Figure 21	Total time spent in the CPU and GPU for the size 1024 FFT horizontal and vertical passes	48

---

## LIST OF TABLES

---

Table 1	Dissertation schedule	6
Table 2	FFT algorithms benchmark. Results are <u>measured in milliseconds</u> for forward and inverse computation with varying input sizes	19

---

## LIST OF LISTINGS

---

5.1	Input buffer bindings . . . . .	29
5.2	Complex multiplication . . . . .	30
5.3	Invocation indices . . . . .	30
5.4	Uniform control variables . . . . .	30
5.5	FFT element index . . . . .	30
5.6	Unique pass structure for Cooley-Tukey . . . . .	34
5.7	Radix-2 Stockham DIF . . . . .	35
5.8	Radix-4 Stockham stage control variables . . . . .	35
5.9	Radix-4 Stockham butterfly . . . . .	36
5.10	Radix-4 Stockham dragonfly inverse arithmetic . . . . .	37
A.1	FFT Radix-2 Cooley-Tukey Horizontal stage pass, see <a href="#">Section 5.3.1</a> . . . . .	53
A.2	FFT Radix-2 Cooley-Tukey Vertical stage pass, see <a href="#">Section 5.3.1</a> . . . . .	54
A.3	FFT Radix-2 Cooley-Tukey Horizontal unique pass, see <a href="#">Section 5.3.1</a> . . . . .	56
A.4	FFT Radix-2 Cooley-Tukey Vertical unique pass, see <a href="#">Section 5.3.1</a> . . . . .	58
A.5	FFT Radix-2 Stockham Horizontal unique pass, see <a href="#">Section 5.3.2</a> . . . . .	60
A.6	FFT Radix-2 Stockham Vertical unique pass, see <a href="#">Section 5.3.2</a> . . . . .	61
A.7	FFT Radix-4 Stockham Horizontal unique pass, see <a href="#">Section 5.3.3</a> . . . . .	62
A.8	FFT Radix-4 Stockham Vertical unique pass, see <a href="#">Section 5.3.3</a> . . . . .	64
B.1	cuFFT, see <a href="#">Section 6.1</a> . . . . .	67
C.1	FFT Radix-2 Cooley-Tukey, see <a href="#">Section 6.3</a> . . . . .	69
C.2	FFT Radix-2 Stockham, see <a href="#">Section 6.3</a> . . . . .	72
C.3	FFT Radix-4 Stockham, see <a href="#">Section 6.3</a> . . . . .	75

---

## INTRODUCTION

---

### 1.1 CONTEXTUALIZATION

The Fast Fourier Transforms have been present in our surroundings for a long time, they're used extensively in digital signal processing including many other areas and they often need to be used in a realtime context, where the computations must be performed fast enough. Fast Fourier Transforms essentially are optimized algorithms to compute the Discrete Fourier Transform of some data, data that might be sampled from a signal, an oscillating object or even an image, which is transformed into the frequency domain allowing any kind of processing for a relatively low computational cost. Despite already existing pretty fast computations of the [FFT](#), many applications require the processing of several transforms so its necessary to manage the implementations properties and achieve the best speed.

### 1.2 MOTIVATION

The continuous progress of the evolution of GPUs has increased the popularity of parallelizable algorithm implementations on this type of hardware. Notably the [FFT](#) algorithms family is constantly present in Computer Graphics, it's usual to find inlined implementations in shader code which offer reliable Fast Fourier Transforms [Flügge \(2017\)](#), but lack tuning of settings for a more optimized versions of these computations. On the other hand there's already out there libraries that provide efficient implementations of [FFT](#) on the GPU and CPU like cuFFT [Nvidia](#), a library provided by NVIDIA exclusively for their GPU's implemented for CUDA, and FFTW [Frigo and Johnson \(2012\)](#), a library dedicated to computations of [FFT](#) on the CPU with SIMD instructions support.

Although this libraries can provide efficient transforms with specialized cases over a proper plan, in some applications its performance might be compromised for cases where, for example, the graphics pipeline needs to be synchronized with the computation of the Fourier Transform.

### 1.3 OBJECTIVES

The main objective of this dissertation is to provide efficient [FFT](#) alternatives in GLSL compared with dedicated tools for high performance of [FFT](#) computations like NVIDIA cuFFT library or FFTW, while analysing the intrinsic of a good Fast Fourier Transform implementation on the GPU and even make a one to one comparison of

implementations on different frameworks. To accomplish the main objective there are two stages taken in consideration, *Analysis of CUDA and GLSL kernels* to be well settled in their differences and to have a reference for the second stage *Analysis of application specific implementations* which will cluster the study's main objective and where we'll use as case of study applications with implementation of the FFT in the field of Computer Graphics that require realtime performance.

With constant progression of the research needed for this project, some steps of the work plan were refactored to meet the needs. The two main stages of the objectives stay the same but there are some adjustments to the schedule dates and steps as shown in Table 1.

- **Research Fast Fourier Transform;**
- **Study cuFFT**, understand internal optimizations and prepare specialized profiles;
- **Analysis of CUDA and GLSL kernels** for FFT raw computations;
- **Research of Application driven FFT**, specialized implementations on the context of the application;
- **Writing of pre-dissertation;**
- **Writing of dissertation.**

	2021			2022							
	Nov	Dez	Jan	Feb	Mar	Apr	May	Jun	Jul	Aug	Set
Research Fast Fourier Transform											
Study cuFFT											
Analysis of CUDA and GLSL kernels											
Research of Application driven FFT											
Writing of pre-dissertation											
Writing of dissertation											

**Table 1:** Dissertation schedule

## 1.4 DOCUMENT ORGANIZATION

This dissertation is organized in 3 chapters. Firstly, the Chapter 1 exposes an introduction to the subject of this dissertation with the respective background information and defines objectives including contextualization and this document organization section.

To give a state of the art overview of the theory and practice associated with Fourier Transforms, Chapter 2 covers most of basic understandings and algorithms needed for later chapters, this will only take simple approaches to each concept to give intuitive insight and empirical explanations without proving it formally.

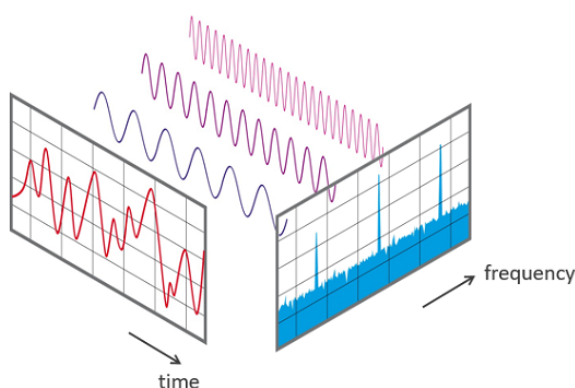
---

## THE FOURIER TRANSFORM

---

It's noticeable the presence of Fourier Transforms in a great variety of apparent unrelated fields of application, even the [FFT](#) is often called ubiquitous due to its effective nature of solving a great hand of problems for the most intended time complexity. Some applications include polynomial multiplication [Jia \(2014\)](#), numerical integration, time-domain interpolation, x-ray diffraction. Furthermore, it is present in several fields of study such as Applied Mechanics, Signal Processing, Sonics and Acoustics, Biomedical Engineering, Instrumentation, Radar, Numerical Methods, Electromagnetics, Computer Graphics and more [Brigham \(1988\)](#).

In **Signal Analysis** when representing a signal with amplitude as function of time, it can be translated to the frequency domain, a domain that consists of signals of sines and cosines waves of varied frequencies, as illustrated in [Figure 1](#), but to calculate the coefficients of those waves we use the Fourier Transform.



**Figure 1:** Time to frequency signal decomposition **Source:** [NTiAudio](#)

Since the sine and cosine waves are in simple wave forms they can then be manipulated with relative ease. This process is constantly present in communications since the transmission of data over wires and radio circuits through signals and most devices nowadays perform it frequently.

In this introductory chapter we present a rudimentary introduction to the Fourier Transform in [Section 2.1](#) and describe the discrete version of the Fourier Transform, which we focus more in this dissertation in [Section 2.2](#) and finalizing with the state of the art of the most popular algorithms in [Section 2.3](#).

## 2.1 CONTINUOUS FOURIER TRANSFORM

The **Fourier Transform** is a mathematical method to transform the domain referred to as *time* of a function, to the *frequency* domain, intuitively the Inverse Fourier Transform is the corresponding method to reverse that process and reconstruct the original function from the one in *frequency* domain representation.

Although there are many forms, the Fourier Transform key definition can be described as:

$$\begin{aligned} X(f) &= \int_{-\infty}^{+\infty} x(t)e^{-ift}dt \\ x(t) &= \frac{1}{2\pi} \int_{-\infty}^{+\infty} X(f)e^{-ift}df \end{aligned} \quad (1)$$

where:

- $X(f), \forall f \in \mathbb{R} \rightarrow$  function in *frequency* domain representation, also called the Fourier Transform of  $x(t)$ ;
- $x(t), \forall t \in \mathbb{R} \rightarrow$  function in *time* domain representation;
- $i \rightarrow$  imaginary unit  $i = \sqrt{-1}$ .

This formulation shows the usage of complex-valued domain, making the Fourier Transform range from real to complex values, one complex coefficient per frequency  $X : \mathbb{R} \rightarrow \mathbb{C}$

If we take into account Euler's formula ([Equation 2](#)), we can rewrite the Fourier Transform as represented in [Equation 3](#).

$$e^{ix} = \cos x + i \sin x \quad (2)$$

$$X(f) = \int_{-\infty}^{+\infty} x(t)(\cos(-ft) + i \sin(-ft))dt \quad (3)$$

Hence, we can break the Fourier Transform apart into two formulas that give each coefficient of the sine and cosine components as functions without dealing with complex numbers.

$$\begin{aligned} X(f) &= X_a(f) + iX_b(f) \\ X_a(f) &= \int_{-\infty}^{+\infty} x(t) \cos(ft)dt \\ X_b(f) &= \int_{-\infty}^{+\infty} x(t) \sin(ft)dt \end{aligned} \quad (4)$$

This model of the Fourier Transform applied to infinite domain functions is called **Continuous Fourier Transform**.

## 2.2 DISCRETE FOURIER TRANSFORM

The Fourier Transform of a finite sequence of equally-spaced samples of a function is the called the **Discrete Fourier Transform** (DFT). It converts a finite set of values in *time* domain to *frequency* domain representation. It is an important version of the Fourier Transform since it deals with a discrete amount of data, therefore, programmers use it to implement on computers.

$$X_k = \sum_{n=0}^{N-1} x_n \cdot e^{-\frac{i2\pi}{N}kn} \quad (5)$$

$$x_n = \frac{1}{N} \sum_{k=0}^{N-1} X_k \cdot e^{\frac{i2\pi}{N}kn} \quad (6)$$

Notably, the discrete version of the Fourier Transform has some obvious differences since it deals with a discrete time sequence. The first difference is that the sum covers all elements of the input values instead of integrating the infinite domain of the function, but we can also notice that the exponential, similar to the aforesaid, divides the values by  $N$  ( $N$  being the total number of elements in the sequence) due to the inability to look at frequency and time  $ft$  continuously we instead take the  $k$ 'th frequency over  $n$ .

We can expand this formula as:

$$X_k = x_0 + x_1 e^{\frac{i2\pi}{N}k} + \dots + x_{N-1} e^{\frac{i2\pi}{N}k(N-1)}$$

Having this sum simplified we then only need to resolve the complex exponential, and we can do that by replacing the  $e^{\frac{i2\pi}{N}kn}$  by the euler formula as mentioned before to reduce the maths to a simple sum of real and imaginary numbers.

$$X_k = x_0 + x_1(\cos b_1 + i \sin b_1) + \dots + x_{N-1}(\cos b_{N-1} + i \sin b_{N-1}) \quad (7)$$

$$\text{where } b_n = \frac{2\pi}{N}kn$$

Finally we'll be left with the result as a complex number

$$X_k = A_k + iB_k$$

**EXAMPLE** Let us now follow an example of calculation of the DFT for a sequence  $x$  with  $N$  number of elements.

$$x = [1 \quad 0.707 \quad 0 \quad -0.707 \quad -1 \quad -0.707 \quad 0 \quad 0.707]$$

$$N = 8$$



With this sequence we now want to transform it into the frequency domain, and for that we need to apply the Discrete Fourier Transform to each element  $x_n \rightarrow X_k$ , thus, for each  $k$ 'th element of  $X$  we apply the DFT for every element of  $x$ .

$$\begin{aligned} X_0 &= 1 \cdot e^{-\frac{i2\pi}{8} \cdot 0 \cdot 0} + 0.707 \cdot e^{-\frac{i2\pi}{8} \cdot 0 \cdot 1} + \dots + 0.707 \cdot e^{-\frac{i2\pi}{8} \cdot 0 \cdot 7} \\ &= (0 + 0i) \end{aligned}$$

$$\begin{aligned} X_1 &= 1 \cdot e^{-\frac{i2\pi}{8} \cdot 1 \cdot 0} + 0.707 \cdot e^{-\frac{i2\pi}{8} \cdot 1 \cdot 1} + \dots + 0.707 \cdot e^{-\frac{i2\pi}{8} \cdot 1 \cdot 7} \\ &= (4 + 0i) \end{aligned}$$

...

$$\begin{aligned} X_7 &= 1 \cdot e^{-\frac{i2\pi}{8} \cdot 7 \cdot 0} + 0.707 \cdot e^{-\frac{i2\pi}{8} \cdot 7 \cdot 1} + \dots + 0.707 \cdot e^{-\frac{i2\pi}{8} \cdot 7 \cdot 7} \\ &= (4 + 0i) \end{aligned}$$

And that will produce our complex-valued output in frequency domain, as simple as that.

$$X = \begin{bmatrix} 0i & 4 + 0i & 0i & 0i & 0i & 0i & 0i & 4 + 0i \end{bmatrix}$$

### 2.2.1 Matrix multiplication

The example shown above is done sequentially as if each frequency pin is computed individually, but there's a way to calculate the same result by using matrix multiplication [Rao and Yip \(2018\)](#). Since the operations are done equally without any extra step we can group all analysing function sinusoids ( $e^{-\frac{i2\pi}{N}kn}$ ), also referred to as twiddle factors.

$$W = \begin{bmatrix} \omega_N^{0 \cdot 0} & \omega_N^{1 \cdot 0} & \dots & \omega_N^{(N-1) \cdot 0} \\ \omega_N^{0 \cdot 1} & \omega_N^{1 \cdot 1} & \dots & \omega_N^{(N-1) \cdot 1} \\ \vdots & \vdots & \ddots & \vdots \\ \omega_N^{0 \cdot (N-1)} & \omega_N^{1 \cdot (N-1)} & \dots & \omega_N^{(N-1) \cdot (N-1)} \end{bmatrix} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ 1 & \omega & \dots & \omega^{(N-1)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & \omega^{(N-1)} & \dots & \omega^{(N-1) \cdot (N-1)} \end{bmatrix}$$

$$\text{where } \omega_N = e^{-\frac{i2\pi}{N}}$$

The substitution variable  $\omega$  allows us to avoid writing extensive exponents.

The symbol  $W$  represents the transformation matrix of the Discrete Fourier Transform, also called DFT matrix, and its inverse can be defined as.

$$W^{-1} = \frac{1}{N} \cdot \begin{bmatrix} 1 & 1 & \dots & 1 \\ 1 & \omega_N & \dots & \omega_N^{(N-1)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & \omega_N^{(N-1)} & \dots & \omega_N^{(N-1) \cdot (N-1)} \end{bmatrix}$$

$$\text{where } \omega_N = e^{-\frac{i2\pi}{N}}$$

By using this matrix multiplication form we can have a more efficient way to compute the DFT.

$$X = W \cdot x$$

$$x = W^{-1} \cdot X$$

Its also worth noting that normalizing the DFT and IDFT matrix be by  $\sqrt{N}$  instead of just normalizing the IDFT by  $N$ , will make  $W$  a unitary matrix [Horn and Johnson \(2012\)](#). However, this normalization by  $\sqrt{N}$  is not common in FFT implementations.

**EXAMPLE** Continuing the example [2.2](#), we can adapt the application of the DFT to the matrix multiplication form.

$$W = \begin{bmatrix} 1 & 1 & \dots & 1 \\ 1 & \omega_8 & \dots & \omega_8^7 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & \omega_8^7 & \dots & \omega_8^{49} \end{bmatrix}$$

$$\text{where } \omega_8 = e^{\frac{i2\pi}{8}}$$

$$X = W \cdot x = W \cdot \begin{bmatrix} 1 \\ 0.707 \\ \vdots \\ 0.707 \end{bmatrix} = \begin{bmatrix} 0 \\ 4 + 0i \\ \vdots \\ 4 + 0i \end{bmatrix}$$

It's conspicuous that the complexity time for each multiplication of every singular term of the sequence with the complex exponential value is  $O(N^2)$ , hence, the computation of the Discrete Fourier Transform rises exponentially as we use longer sequences. Therefore, over time new algorithms and techniques where developed to increase the performance of this transform due to its usefulness.

## 2.3 FAST FOURIER TRANSFORM

The **Fast Fourier Transform (FFT)** is a family of algorithms that compute the Discrete Fourier Transform (DFT) of a sequence, and its inverse, efficiently, since the direct usage of the DFT formulation is too slow for its applications. Thus, FFT algorithms exploit the DFT matrix structure by employing a divide-and-conquer approach (Chu and George (1999)) to segment its application.

Over time several variations of the algorithms were developed to improve the performance of the DFT and many aspects were rethought in the way we apply and produce the resulting transform.

There are many algorithms and approaches on the FFT family such as the well known Cooley-Tukey, known for its simplicity and effectiveness to compute any sequence with size as a power of two, but also Rader's algorithm Rader (1968) and Bluestein's algorithm Bluestein (1970) to deal with prime sized sequences, and even the Split-radix FFT Yavne (1968) that recursively expresses a DFT of length  $N$  in terms of one smaller DFT of length  $N/2$  and two smaller DFTs of length  $N/4$ .

The next two sections focus on the Cooley–Tukey algorithm, most specifically the radix-2 decimation-in-time (DIT) FFT and radix-2 decimation-in-frequency (DIF) FFT, both requiring the input sequence to have a power of two size. These two variations of the Cooley–Tukey algorithm represent the state of the art of what an individual in need to implement FFT will most likely be familiar with.

### 2.3.1 Radix-2 Decimation-in-Time FFT

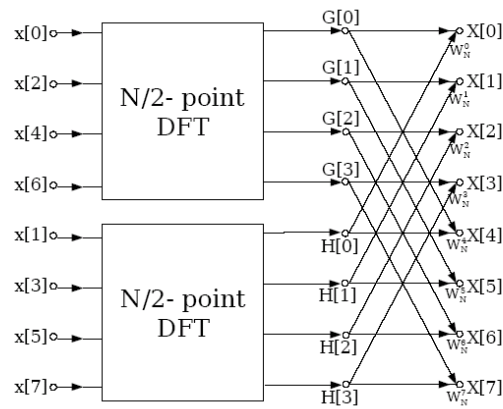
The Radix-2 Decimation-in-Time FFT algorithm rearranges the computation of a DFT of size  $N$  into two DFTs of size  $N/2$ , one as a sum over the even indexed elements and other as a sum over the odd indexed elements. Cooley and Tukey proved this possibility of dividing the DFT computation into two smaller DFT by exploiting the symmetry of this division, as presented in Equation 8. Hence, it is hinted the recursive definition of this algorithm on both DFT of size  $N/2$ .

$$\begin{aligned}
 X_k &= \sum_{n=0}^{N-1} x_n \cdot \omega_N^{kn} \\
 X_k &= \sum_{n=0}^{N/2-1} x_{2n} \cdot \omega_N^{k(2n)} + \sum_{n=0}^{N/2-1} x_{2n+1} \cdot \omega_N^{k(2n+1)} \\
 X_k &= \sum_{n=0}^{\frac{N}{2}-1} x_{2n} \cdot \omega_{N/2}^{kn} + \omega_N^k \sum_{n=0}^{\frac{N}{2}-1} x_{2n+1} \cdot \omega_{N/2}^{kn} \tag{8}
 \end{aligned}$$

where  $\omega_N = e^{\frac{j2\pi}{N}}$

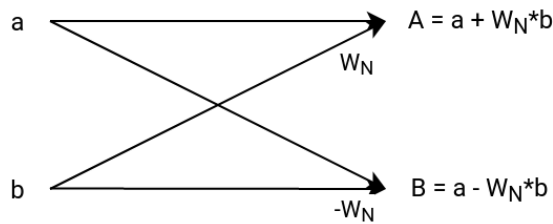
This formulation successfully segments the full sized DFT into two  $N/2$  sized DFT's of the even and odd indexed elements where the later is multiplied by a factor called twiddle  $\omega_N^k$ .

This algorithm is a Radix-2 Decimation-in-Time in the sense that the time values are regrouped in 2 subtransforms, and the decomposition reduces the time values to the frequency domain. Since the understanding of this algorithm can be applied recursively, the Figure 2 illustrates the basic behaviour and represents the  $N/2$  subtransforms with boxes that can be filled by the recursive application of this algorithm to produce the frequency domain sequence.



**Figure 2:** Radix-2 Decimation-in-Time FFT **Source:** Jones (2014)

Effectively, this smaller DFT's are recursively reduced by this algorithm until there's only the computation of a length-2 DFT. On each stage it is applied the Cooley-Tukey butterfly operation (Chu and George (1999)) with a shifted element according to the subtransform size, as illustrated in Figure 3.



**Figure 3:** Cooley-Tukey butterfly

The complexity work within the algorithm is distributed with the DIT approach which decomposes each DFT by 2 having  $\log(N)$  stages Smith (2007). There are  $N$  complex multiplications needed for each stage of the DIT decomposition, therefore, the multiplication complexity for a  $N$  sized DFT is reduced from  $O(N^2)$  to  $O(N \log(N))$ .

The splitting of the DFT into two smaller half sized DFTs causes the original input sequence to require a special reordering to pass the even and odd numbers into, and when this algorithm is applied recursively, this reordering is always needed, so in the end we need a special order for the input elements, fortunately, as noted by ? this

order corresponds to the bit reversed elements of the sequence, therefore, we need to bit reverse all elements of the input sequence.

The bit reversal of the input sequence corresponds to the permutation of swapping the elements position to its bit reversed index, as illustrated in [Figure 4](#). The *bit\_reverse* of an index depends directly on the indexing domain of the input sequence, therefore, it needs the size  $N$ , or more precisely the  $\log(N)$  value, to use as a reference to reverse the bit order while maintaining the value within the sequence range.



**Figure 4:** Bit reverse permutation

For example, for a sequence of size 16, we have some index with  $\log 16$  bits  $b_1b_2b_3b_4$ , which corresponds to the bit reversed index as  $b_4b_3b_2b_1$ .

Both the DIT and DIF FFT algorithms require this bit reversal permutation step, but in the case of the DIT we bit reverse the input sequence and on the DIF we apply it on the output, to result in a natural order sequence.

There are many implementations of the bit reversal, and since it is quite simple, any decent version can be used in regards of this FFT algorithm since it is not the main bottleneck. An algorithm such as [Algorithm 1](#) can be used for the *bit\_reverse* function or any other efficient alternatives ([Prado \(2004\)](#)).

---

**Algorithm 1:** Bit reverse

---

**Data:** Integer  $i$ **Result:** Bit reversed integer  $i$  $n \leftarrow 0$  **foreach**  $i = 0$  **to**  $\log(N) - 1$  **do**     $n \leftarrow n \ll 1$ ;     $n \leftarrow n | (i \& 1)$ ;     $i \leftarrow i \gg 1$ ;**end****return**  $i$ ;

---

In practice, [Algorithm 2](#) demonstrates the aforesaid as an iterative possible implementation. Although this algorithm is congruent with a code implementation, its worth noting that the input sequence can either have real or complex numbers, since the arithmetic is the same for both domains the only thing that needs to be specialized is the operator overloading in the inner most loop.

---

**Algorithm 2:** Radix-2 Decimation-in-Time Forward FFT

---

**Data:** Sequence  $in$  with size  $N$  power of 2**Result:** Sequence  $out$  with size  $N$  with the DFT of the input

/\* Bit reversal step

\*/

**foreach**  $i = 0$  **to**  $N - 1$  **do**     $out[\text{bit\_reverse}(i)] \leftarrow in[i]$ **end**

/\* FFT

\*/

**foreach**  $s = 1$  **to**  $\log(N)$  **do**     $m \leftarrow 2^s$ ;     $w_m \leftarrow \exp(-2\pi i / m)$ ;    **foreach**  $k = 0$  **to**  $N - 1$  **by**  $m$  **do**         $w \leftarrow 1$ ;        **foreach**  $j = 0$  **to**  $m/2$  **do**             $bw \leftarrow w \cdot out[k + j + m/2]$ ;             $a \leftarrow out[k + j]$ ;             $out[k + j] \leftarrow a + bw$ ;             $out[k + j + m/2] \leftarrow a - bw$ ;             $w \leftarrow w \cdot w_m$ ;        **end**    **end****end****return**  $out$ ;

---

### 2.3.2 Radix-2 Decimation-in-Frequency FFT

The Radix-2 Decimation-in-Frequency FFT algorithm is very similar to the DIT approach, its based on the same principle of divide-and-conquer but it rearranges the original Discrete Fourier Transform (DFT) into the computation of two transforms, one with the even indexed elements and other with the odd indexed elements; as in this simplified formulation [Equation 9](#).

$$X_{2k} = \sum_{n=0}^{\frac{N}{2}-1} (x_n + x_{n+\frac{N}{2}}) \cdot \omega_{N/2}^{kn} \quad (9)$$

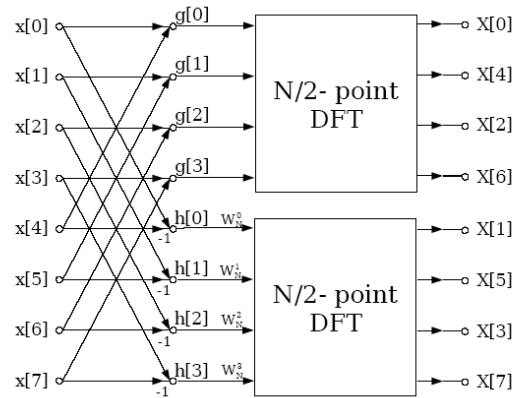
$$X_{2k+1} = \sum_{n=0}^{\frac{N}{2}-1} ((x_n - x_{n+\frac{N}{2}}) \cdot \omega_{N/2}^{kn}) \cdot \omega_N^n$$

where  $\omega_N = e^{\frac{j2\pi}{N}}$

The DFT divided into these two transforms from the full sized DFT By separating these two transforms from the full sized DFT we get two distinct

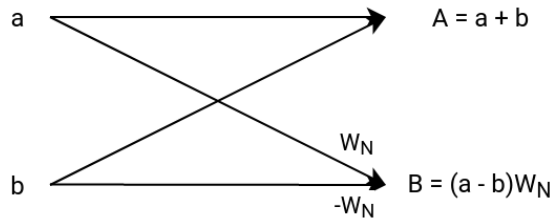
Notably, this formulation distinguishes the full sized DFT into two  $N/2$  sized DFT's of the even and odd indexed elements where the later is multiplied by a twiddle factor  $\omega_N^k$  with both outside the same context.

This algorithm is a Radix-2 Decimation-in-Frequency since the DFT is decimated into two distinct smaller DFT's and the frequency samples will be computed separately in different groups, as if the regrouping of the DFT's would reduce directly to the frequency domain. Since the understanding of this algorithm can be applied recursively, the [Figure 5](#) illustrates the this behaviour and represents the  $N/2$  subtransforms with boxes that can be filled by the recursive application of this algorithm to produce the frequency domain sequence. Additionally this illustration can be compared to [Figure 2](#) since both are symmetrically identical.



**Figure 5:** Radix-2 Decimation-in-Frequency FFT **Source:** [Jones \(2014\)](#)

Similarly to the DIT version, the DFT can be recursively reduced by the DIF algorithm until theres only the computation of a length-2 DFT. On each stage it is applied the Gentleman-Sande butterfly operation ([Chu and George \(1999\)](#)) with a shifted element according to the subtransform size, as illustrated in [Figure 6](#).



**Figure 6:** Gentleman-Sande butterfly

Since this algorithm has similarities with the DIT, its complexity also lives to this similarity, maintaining the same  $O(N \log(N))$  for number of multiplications, despite that, [Figure 6](#) and [Figure 3](#) might look different in number of arithmetic operations since the first has 1 addition, 1 subtraction, and 2 multiplications, and the second has 1 addition, 1 subtraction, and 1 multiplication, but effectively the  $W_N \cdot b$  can be reused and only computed once as seen in [Algorithm 2](#).

As mentioned in [Section 2.3.1](#) the bit reversal in DIF works a bit differently, this algorithm does the exact opposite of the DIT since it requires a natural order sequence and returns a bit reversed output, justifying why this step is applied after the algorithm.

In practice, [Algorithm 3](#) demonstrates the aforesaid with an iterative representation of a possible implementation. Although this algorithm is congruent with a code implementation, its worth noting that the input sequence can either have real or complex numbers, since the arithmetic is the same for both domains the only thing that needs to be specialized is the operator overloading in the inner most loop.



---

**Algorithm 3:** Radix-2 Decimation-in-Frequency Forward FFT

---

**Data:** Sequence *in* with size  $N$  power of 2**Result:** Sequence *out* with size  $N$  with the DFT of the input

```

/* FFT                                                    */
foreach  $s = 0$  to  $\log(N) - 1$  do
     $gs \leftarrow N \gg s$ ;
     $w_{gs} \leftarrow \exp(2\pi i / gs)$ ;
    foreach  $k = 0$  to  $N - 1$  by  $gs$  do
         $w \leftarrow 1$ ;
        foreach  $j = 0$  to  $gs/2 - 1$  do
             $a \leftarrow in[k + j + gs/2]$ ;
             $b \leftarrow in[k + j]$ ;
             $in[k + j] \leftarrow a + b$ ;
             $in[k + j + gs/2] \leftarrow (a - b) \cdot w$ ;
             $w \leftarrow w \cdot w_{gs}$ ;
        end
    end
end
/* Bit reversal step                                      */
foreach  $i = 0$  to  $N - 1$  do
     $out[bit\_reverse(i)] \leftarrow in[i]$ 
end
return out;

```

---

---

## ALGORITHMS ANALYSIS

---

To flavour this pre-dissertation report some work of benchmarking and analysis were done to compete with the theoretical explanations addressed on [Chapter 2](#). Hence, some implementations were tested to provide coherence to what has been studied, and algorithms such as [Algorithm 2](#) Radix-2 Decimation-in-Time [Algorithm 3](#) Radix-2 Decimation-in-Frequency were timed in [Table 2](#).

	Size 128	Size 256	Size 512	Size 1024	Size 2048
<b>DFT</b>	5.16593	17.2782	70.5689	293.104	1246.44
<b>FFT DIT</b>	0.169113	0.37668	0.86415	1.8793	4.47742
<b>FFT DIF</b>	0.159458	0.378722	0.881921	1.90661	4.13369
<b>Recursive FFT</b>	0.210895	0.485643	1.4421	2.32922	5.1178

**Table 2:** FFT algorithms benchmark. Results are measured in milliseconds for forward and inverse computation with varying input sizes

As we can see the Discrete Fourier Transform increases exponentially for higher sized sequences, as expected all FFT variants perform critically better than the original formulation.

One variant that wasn't exposed much in the above chapters is the Recursive FFT algorithm, which corresponds to a Decimation-in-Time approach with recursive reduction, therefore, this algorithm aggregates the divide-and-conquer method but with the disadvantage of recursive function overhead. This recursive look at the DIT approach can be easier to implement since the bit reversal step isn't explicitly applied before starting the FFT.

Finally, as expected the DIT and DIF algorithms overrule the other alternatives for every sized input.

---

## COMPUTATION OF THE FOURIER TRANSFORM

---

Nowadays there is a lot more to the computation of FFT's than just the basic Cooley-Tukey algorithm described in [Section 2.3.1](#) there are more algorithms, variations and improvements that enhance the computation in many aspects. Even the target hardware can change the restrictions on the computation of this primitive.

One could optimize the FFT in many ways, and in this section we introduce an intermediate step in the Cooley-Tukey algorithm to turn it into a natural order algorithm without explicit bit reversal permutation in [Section 4.1](#). Furthermore, this serves as base knowledge to introduce the Radix-2 Stockham algorithm in [Section 4.2](#). Finally we upgrade the Stockham algorithm to Radix-4 in [Section 4.3](#) and consequently the advantages and disadvantages of this.

### 4.1 NATURAL ORDER COOLEY-TUKEY

Despite existing already really fast solutions for index bit reversal ([Prado \(2004\)](#)), this *shuffle* step still weights the algorithms with extra overhead. The natural order Cooley-Tukey FFT is a modification of the Cooley-Tukey algorithm that allows the removal of this step by computing the butterfly and reordering the elements per stage ([OKAHISA](#)).

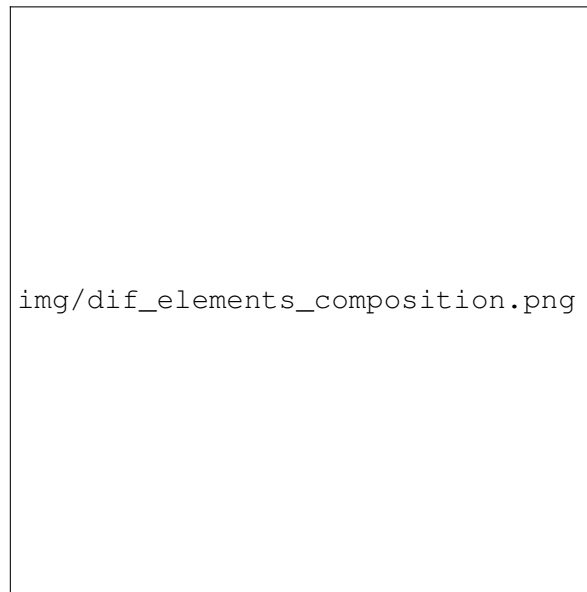
At the end of each stage the even and odd elements are composed in such a way that the elements will be in natural order at the end, as illustrated in [Figure 7](#). This composition follows the indexing scheme described in [Equation 10](#).

$$x[q + 2 * p] = y[q + p] \tag{10}$$

$$x[q + 2 * p + 1] = y[q + p + m] \tag{11}$$

Where  $x$  and  $y$  are alternated sequences for read or write over each stage,  $q$  corresponds to the sub FFT offset in this stage,  $p$  is the index of the sub FFT element shift for the current butterfly being computed, and finally  $m$  corresponds to the size of the sub FFT divided by 2.

Although this additional composition got rid of the bit reversal step in the Cooley-Tukey's algorithm, the computation is overloaded with more work for each stage, since first the algorithm computes the butterfly and then it reorders the elements every stage.



**Figure 7:** Chain of even and odd compositions over each stage for a natural order DIF

This reordering after the butterfly may seem unnecessary, therefore we can use a new sequence and alternate between the two each stage, hence, providing this composition when we write the butterfly results. This is where the Stockham algorithm comes in.

## 4.2 STOCKHAM ALGORITHM

As mentioned in the previous section, the Stockham algorithm eliminates the need to have the bit reversal permutation, it does this by taking advantage of a reordering of the elements (Govindaraju et al. (2008)) in between stages illustrated in Figure 7. The natural order elements are composed stage by stage and the butterfly computations stay the same, so this approach takes advantage of the Cooley-Tukey algorithm and turns it into a more suitable form for highly parallelizable hardware such as GPUs, making it a best fit for our implementation in any GPU programmable language.

Similarly to Section 4.1 the algorithm requires the usage alternated sequences for read write over each stage. With this, the read of an element which has been altered before read in a given stage. In the end the result will be in the last sequence we wrote to.

The Stockham algorithm is described in Algorithm 4 and this version may seem strictly different from the Cooley Tukey, specially the inner most loop, However, most of the logic stays the same, the indexing is a bit different and simpler for this version and the if statements are a consequence of using alternated ping pong sequences since there has to be branches for read and write of both arrays for this out-of-place algorithm.

**Algorithm 4:** Stockham Radix-2 Decimation-in-Frequency Forward FFT**Data:** Sequence *pingpong0* with size  $N$  power of 2**Result:** Sequence *out* with size  $N$  with the DFT of the input

---

```

foreach  $s = 0$  to  $\log(N) - 1$  do
     $gs \leftarrow N \gg s$ ;
     $stride \leftarrow 1 \ll s$ ;
    foreach  $i = 0$  to  $N - 1$  do
         $p \leftarrow i \text{ div } stride$ ;
         $q \leftarrow i \text{ mod } stride$ ;
         $w_p = \exp(-2\pi i / gs * p)$ ;
        if  $stage \text{ mod } 2 == 0$  then
             $a \leftarrow pingpong0[q + s * (p + 0)]$ ;
             $b \leftarrow pingpong0[q + s * (p + gs/2)]$ ;
            /* Perform butterfly */
             $pingpong1[q + s * (2 * p + 0)] = a + b$ ;
             $pingpong1[q + s * (2 * p + 1)] = (a - b) * w_p$ ;
        else
             $a \leftarrow pingpong1[q + s * (p + 0)]$ ;
             $b \leftarrow pingpong1[q + s * (p + gs/2)]$ ;
            /* Perform butterfly */
             $pingpong0[q + s * (2 * p + 0)] = a + b$ ;
             $pingpong0[q + s * (2 * p + 1)] = (a - b) * w_p$ ;
        end
    end
end
if  $\log(N) \text{ mod } 2 == 0$  then
    return pingpong1;
else
    return pingpong0;
end

```

---

This algorithm description will give us a solid code base to use as reference when implementing it on the GPU, mainly due to the way we are indexing and the usage of alternating read write sequences.

### 4.3 RADIX-4 INSTEAD OF RADIX-2

We've discussed about multiple radix-2 approaches , however, we can explore a wide range of alternatives when we get into higher radices other than just radix-2. These FFT algorithms can use higher radix for better performance and even mixed radix Singleton (1969) for wider range of input sequence sizes.

The radix-4 is a good improvement over radix-2 since getting to higher radix than 4 applies size constraints that are not suitable for our case of study. Rearranging the Stockham algorithm to radix-4 upgrades the computation of the butterflies while reducing the number of stages which will be crucial in later sections.

Theoretically, radix-4 formulation can be twice as fast as a radix-2 Hussain et al. (2010) since it only takes half the stages with more complexity in the butterflies which are sometimes called dragonflies. Additionally, we can use less multiplications with better factorizations Marti-Puig and Bolano (2009).

The radix-4 Stockham splits the FFT of size  $N$  into four subtransforms of size  $N/4$  each stage, therefore this algorithm only features  $\log(N)/2$  stages and requires the size to be multiple of 4. Since this is a natural order algorithm the computation of the dragonfly includes the reordering of the elements in natural order every stage, as illustrated in Figure 8.



**Figure 8:** Illustration of a stage in radix-4 Stockham with each color representing a radix-4 butterfly computation

The forward dragonfly for this algorithm is presented in Figure 9 and its computation involves 4 elements.

With each stage the dragonflies are calculated and the elements are reordered around, therefore, in the end we get Algorithm 5 and its inverse Algorithm 6.



**Figure 9:** Radix-4 FFT butterfly structure **Source:** Marti-Puig and Bolano (2009)

---

**Algorithm 5:** Stockham Radix-4 Decimation-in-Frequency Forward FFT

---

**Data:** Sequence *pingpong0* with size  $N$  power of 4

**Result:** Sequence *out* with size  $N$  with the DFT of the input

**foreach**  $stage = 0$  **to**  $\log(N)/2 - 1$  **do**

$n \leftarrow 1 \ll ((\log(N)/2) - stage) * 2$ ;

$s \leftarrow 1 \ll (stage * 2)$ ;

$n0 \leftarrow 0$ ;

$n1 \leftarrow n/4$ ;

$n2 \leftarrow n/2$ ;

$n3 \leftarrow n1 + n2$ ;

**foreach**  $i = 0$  **to**  $N - 1$  **do**

$p \leftarrow i \text{ div } s$ ;

$q \leftarrow i \text{ mod } s$ ;

$w_{1p} = \exp(-2\pi i / n * p)$ ;

$w_{2p} = w_{1p} * w_{1p}$ ;

$w_{3p} = w_{1p} * w_{2p}$ ;

**if**  $stage \text{ mod } 2 == 0$  **then**

$a \leftarrow pingpong0[q + s * (p + n0)]$ ;

$b \leftarrow pingpong0[q + s * (p + n1)]$ ;

$c \leftarrow pingpong0[q + s * (p + n2)]$ ;

$d \leftarrow pingpong0[q + s * (p + n3)]$ ;

/\* Perform dragonfly

\*/

$pingpong1[q + s * (4 * p + 0)] = a + c + b + d$ ;

$pingpong1[q + s * (4 * p + 1)] = w_{1p} * ((a - c) - (b - d) * \sqrt{-1})$ ;

$pingpong1[q + s * (4 * p + 2)] = w_{2p} * ((a + c) - (b + d))$ ;

$pingpong1[q + s * (4 * p + 3)] = w_{3p} * ((a - c) + (b - d) * \sqrt{-1})$ ;

**else**

**Algorithm 6:** Stockham Radix-4 Decimation-in-Time Inverse FFT**Data:** Sequence *pingpong0* with size  $N$  power of 4**Result:** Sequence *out* with size  $N$  with the DFT of the input

```

foreach stage = 0 to  $\log(N)/2 - 1$  do
     $n \leftarrow 1 \ll ((\log(N)/2) - \textit{stage} * 2);$ 
     $s \leftarrow 1 \ll (\textit{stage} * 2);$ 
     $n0 \leftarrow 0;$ 
     $n1 \leftarrow n/4;$ 
     $n2 \leftarrow n/2;$ 
     $n3 \leftarrow n1 + n2;$ 
    foreach  $i = 0$  to  $N - 1$  do
         $p \leftarrow i \text{ div } s;$ 
         $q \leftarrow i \text{ mod } s;$ 
         $w_{1p} = \exp(2\pi i / n * p);$ 
         $w_{2p} = w_{1p} * w_{1p};$ 
         $w_{3p} = w_{1p} * w_{2p};$ 
        if  $\textit{stage} \bmod 2 == 0$  then
             $a \leftarrow \textit{pingpong0}[q + s * (p + n0)];$ 
             $b \leftarrow \textit{pingpong0}[q + s * (p + n1)];$ 
             $c \leftarrow \textit{pingpong0}[q + s * (p + n2)];$ 
             $d \leftarrow \textit{pingpong0}[q + s * (p + n3)];$ 
            /* Perform dragonfly */
             $\textit{pingpong1}[q + s * (4 * p + 0)] = a + c + b + d;$ 
             $\textit{pingpong1}[q + s * (4 * p + 1)] = w_{1p} * ((a - c) + (b - d) * \sqrt{-1});$ 
             $\textit{pingpong1}[q + s * (4 * p + 2)] = w_{2p} * ((a + c) - (b + d));$ 
             $\textit{pingpong1}[q + s * (4 * p + 3)] = w_{3p} * ((a - c) - (b - d) * \sqrt{-1});$ 
        else
             $a \leftarrow \textit{pingpong1}[q + s * (p + n0)];$ 
             $b \leftarrow \textit{pingpong1}[q + s * (p + n1)];$ 
             $c \leftarrow \textit{pingpong1}[q + s * (p + n2)];$ 
             $d \leftarrow \textit{pingpong1}[q + s * (p + n3)];$ 
            /* Perform dragonfly */
             $\textit{pingpong0}[q + s * (4 * p + 0)] = a + c + b + d;$ 
             $\textit{pingpong0}[q + s * (4 * p + 1)] = w_{1p} * ((a - c) - (b - d) * \sqrt{-1});$ 
             $\textit{pingpong0}[q + s * (4 * p + 2)] = w_{2p} * ((a + c) - (b + d));$ 
             $\textit{pingpong0}[q + s * (4 * p + 3)] = w_{3p} * ((a - c) + (b - d) * \sqrt{-1});$ 
        end
    end
end
if  $\log(N) \bmod 2 == 0$  then
    return pingpong1;
else
    return pingpong0;
end

```



---

## IMPLEMENTATION ON THE GPU

---

To be able to analyse the implementation of FFT on the GPU, we need to provide some background on what kind of hardware we're dealing with, what constraints are relevant and what we can or can't do, after, we proceed with a description of how the Fourier Transform suits the GPU programming model and finally we'll move forward to a detailed analysis of the implementation of FFT algorithms in a compute shader in GLSL.

### 5.1 GPU PROGRAMMING MODEL

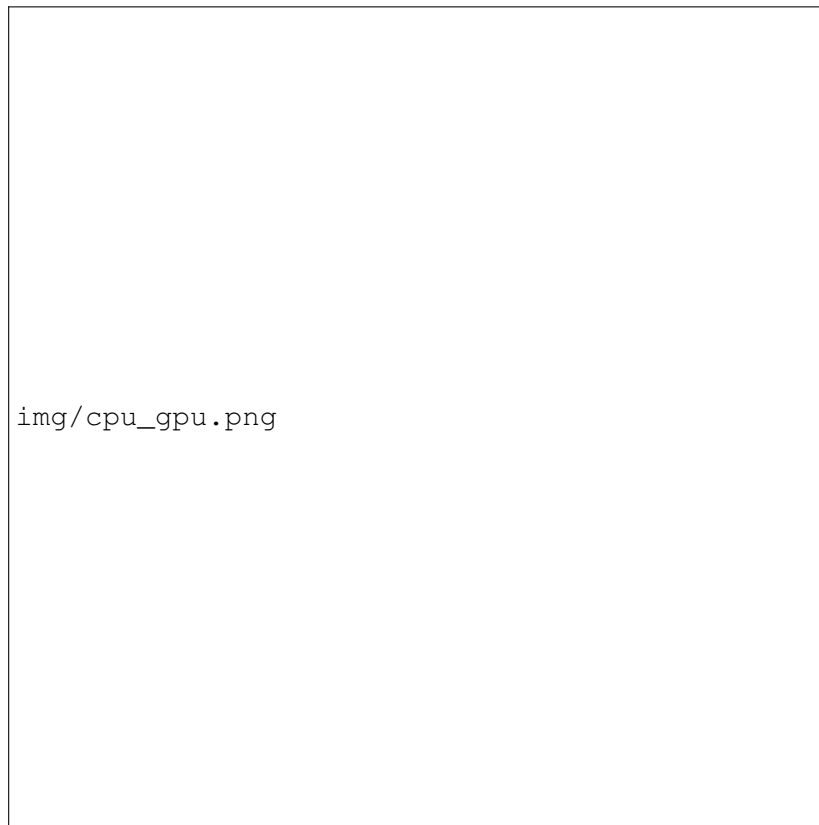
The Graphics Processing Unit (GPU) provides much higher instruction throughput and memory bandwidth than the CPU within a similar price and power envelope having a much faster performance growth curve. Many applications leverage these higher capabilities to run faster on the GPU than on the CPU. While the hardware may change some components and architecture on different models, fundamentally most modern GPUs adopt a specific kind of single instruction multiple data (SIMD) stream architecture which is single instruction multiple threads (SIMT) that is an extension of the SIMD paradigm with large scale multi-threading, streaming memory and dynamic scheduling.

A SIMT stream architecture forwards an instruction to multiple threads with dedicated memory for each instance, allowing data-level parallelism to compute more quickly and effectively. A modern CPU uses a von Neumann architecture, so it executes instructions sequentially one at a time and updates the memory progressively. However, a stream architecture processor works in a slightly different way, they contain multiple streams of simpler processors with shared memory just like the illustration in [Figure 10](#). These processors execute programs called kernels that receive a finite set of input fragments to produce another set of output fragments in parallel [Fernando et al. \(2004\)](#).

With the introduction of programmable General Purpose Graphics Processing Unit (GPGPU) the industry required a lot more parallel algorithms.

platforms allow developers to ignore the language barrier that exists between the CPU and the GPU and, instead, focus on higher-level computing concepts.

Developers and researchers often attend more specific scenarios and they configure multiple GPUs with distributed workflows in sync with each other to take better advantage of this hardware for more intensive problems [Heldens et al. \(2022\)](#).



**Figure 10:** CPU architecture compared with a GPU architecture

The GPU programming model although sophisticated, at its core it provides thread groups hierarchies, shared memories and GPU barriers for synchronization of threads, such abstractions provide data and thread level parallelism for the developer to utilize.

In this programming model we introduce the concept of granularity, that refers to the amount of computation relatively to the transfer of data. It is used to describe types of parallelism as fine-grained or coarse-grained. Fine-grained parallelism describes a small multi threaded tasks in kernel size and execution time with frequent small data transfers between the processors. On the opposite side there's coarse-grained parallelism that describes larger amounts of computation followed by larger infrequent data transfers [NVIDIA](#).

This concept is important to reflect on the association of the GPU programming model with further details on the provided implementations within the next chapters.

Most GPGPU programming frameworks such as CUDA provide coarse-grained data and task parallelism with nested fine-grained data and thread parallelism. This type of task hierarchy architecture promotes the partitioning of the problems to independently solvable subgroups of smaller problems which fit into smaller blocks of threads that can be solved cooperatively in parallel. This relevant information provides us more insight of how the work groups for FFT should be dispatched.

## 5.2 2D FOURIER TRANSFORM ON THE GPU

Computing a 2D Fourier Transform requires a two dimensional input sequence to produce another for the output. However, we could use any usable arbitrary values in the the real world we often use this 2D FFT on images and it isn't immediately obvious how this should be applied since the target data source might have three color channels, therefore, three different values for this multidimensional sequence element that can be used, and not even in floating point complex domain. Precisely, it needs to be adapted to the application use case, we may use it as greyscale image if want a derivation of the luminance via quantized RGB signals of the image (?), use only the values of one channel, or compute multiple FFT's for each channel values. Either way we will preferably at the end have a two dimensional buffer with floating point complex values prepared to be used.

The 2D FFT is computed by performing single dimension FFT's for every row and then for every columns after that ?, so we can divide its application to a horizontal and vertical pass. This describes the way 2D FFT are computed but it is independent of the 1D FFT implementation chosen, so there's freedom to use any type of algorithm.



**Figure 11:** High level illustration of horizontal and vertical passes

At the end of a forward FFT vertical pass the result will be a 2D complex buffer with the frequency domain values of the original image as illustrated in [Figure 11](#).

### 5.3 GLSL IMPLEMENTATION

The implementations were made using GLSL, a high-level shader language for graphics API's such as OpenGL, and we used the compute pipeline from OpenGL to integrate a FFT implementation using compute shaders, which are general purpose programmable shaders.

Since there are many aspects that may impact on the performance the implementation was an iterative process that required researching and testing to compose an optimal solution. One of the main targets was to keep the code generalized so that it could be used as base for other implementations using FFTs with ease.

The next sections go in detail about the way every major iteration evolved into the next one and why there was the need to do it, starting by the Cooley-Tukey algorithm (Section 5.3.1) then progressively improving to the Stockham algorithm (Section 5.3.2 and Section 5.3.3), all this while implementing good GPGPU programming strategies.

#### 5.3.1 Cooley-Tukey

The GPU implementation took as a starting point was with the DIT Cooley-Tukey algorithm, since it is the most popular one with time complexity of  $O(N \log(N))$ , and it is based on the iterative version adapted for parallel processors.

Since this implementation is highly parallel there's the need to separate the reads and writes for each processor into two different buffers in memory due lack of order between processors, therefore, the declaration of two complex pingpong buffers.

```
layout (binding = 0, rg32f) uniform image2D pingpong0;
layout (binding = 1, rg32f) uniform image2D pingpong1;
```

**Listing 5.1:** Input buffer bindings

The read write control for this buffers can be achieved with a flag variable `pingpong`.

Initially this algorithm will work with a pass per stage approach, where a kernel is dispatched every stage and since there's a synchronization step at the end of the pass its granted that all the work groups have finished off writing to the buffer when a pass ends. This case holds up for both the horizontal and vertical FFT steps.

As a result, each kernel has the opportunity to work within every segment of the image, so the local threads can be dispatched with two dimensions, so each work group will have a total of 32 local threads, a reference number used in this implementations for setting up local threads in a work group, this may vary depending on the GPU for optimal performance but 32 is a good number to fill in the thread warp size os most GPUs.

So an example dispatch group for this implementation could be  $(fft\_width/8, fft\_height/8)$  work groups since 8 is the number of threads in the  $y$  axis

By using GLSL there is some advantages on the complex values operations, since the addition and subtraction for vector types already have operators overloading which function the same as in the complex domain. However,

the multiplication works a bit differently so we need to provide an auxiliary function to abstract and support this operator.

```
vec2 complex_mult(vec2 v0, vec2 v1) {
    return vec2(v0.x * v1.x - v0.y * v1.y,
               v0.x * v1.y + v0.y * v1.x);
}
```

**Listing 5.2:** Complex multiplication

Due to the adoption of a different programming paradigm the FFT segment iteration loop doesn't exist such as in [Algorithm 2](#), instead the processors identifiers are used to fetch the index of the butterflies they're gonna work on based on the work groups and threads dispatch setup, and this holds up for any implementation using compute shaders.

```
int line = int(gl_GlobalInvocationID.x);
int column = int(gl_GlobalInvocationID.y);
```

**Listing 5.3:** Invocation indices

Since this first approach is a dynamic implementation that invokes a pass per stage some stage control variables need to be feed into the shader in order to compute the correct butterfly index or control the butterfly process.

```
uniform int pingpong;
uniform int log_width;
uniform int stage;
uniform int fft_dir;
```

**Listing 5.4:** Uniform control variables

Effectively, we use these shader uniform input variables and obtain actual index we're gonna use on the 1D FFT of the image.

```
int group_size = 2 << stage;
int shift = 1 << stage;

int idx = (line % shift) + group_size * (line / shift);
```

**Listing 5.5:** FFT element index

To calculate the twiddle factor we use Euler's formula such as in [Equation 2](#) and resort to the control variable `fft_dir` to flip the twiddle factor for the inverse if we want to reuse this shader.

```
vec2 euler(float angle) {
    return vec2(cos(angle), sin(angle));
}
```

```

void main() {
    // ...
    vec2 w = euler(fft_dir * 2 * (M_PI / group_size) * ((idx % group_size) %
    shift));
    // ...
}

```

Now with the computed twiddle factor we may proceed to compute and store the Cooley-Tukey FFT butterfly, and that's where we need to control the reads and writes for each stage. Since the `pingpong` variable is toggled every pass invocation, it is used to choose with what image we will use to lookup the elements and which one to store into, and that is achieved with an if statement, just like it is used in [Algorithm 4](#).

The butterfly computation itself is simply calculated as a Cooley-Tukey DIT butterfly as illustrated in [Figure 3](#).

```

if (pingpong == 0) {
    // Read
    a = imageLoad(pingpong0, ivec2(idx, column)).rg;
    b = imageLoad(pingpong0, ivec2(idx + shift, column)).rg;

    // Compute and store
    vec2 raux = a + complex_mult(w, b);
    imageStore(pingpong1, ivec2(idx, column), vec4(raux, 0, 0));
    raux = a - complex_mult(w, b);
    imageStore(pingpong1, ivec2(idx + shift, column), vec4(raux, 0, 0));
}
else {
    // Read
    a = imageLoad(pingpong1, ivec2(idx, column)).rg;
    b = imageLoad(pingpong1, ivec2(idx + shift, column)).rg;

    // Compute and store
    vec2 raux = a + complex_mult(w, b);
    imageStore(pingpong0, ivec2(idx, column), vec4(raux, 0, 0));
    raux = a - complex_mult(w, b);
    imageStore(pingpong0, ivec2(idx + shift, column), vec4(raux, 0, 0));
}

```

Finally there's only one step missing, the bit reversal of indices to have a natural order result. A `bit_reverse` function could be easily defined, However, there's already a GLSL alternative which is `bitfieldReverse` together with `bitfieldExtract` ?.

Despite not having a noticeable performance hit, it is good practice to use GLSL predefined functions and operators since they might be optimised for that specific device hardware and using these functions instead of a

handmade implementation reduces the kernel size significantly about approximately 400 bytes (evaluated using *glslang*) since they are reusable functions.

```
int bit_reverse(int k) {
    uint br = bitfieldReverse(k);
    return int(bitfieldExtract(br, 32 - log_width, log_width));
}
```

This auxiliary function will be conditionally used inside the branching if statement for alternating read writes to be only applied on the first stage of transform both for the possible values of `pingpong == 0` and `pingpong == 1`, since the vertical pass might start reading on the first pingpong buffer. This is not the case for the horizontal pass, its ensured that the first stage will always read from `pingpong0` reforcing that the bit reverse branching when `pingpong == 0` is disposable.

```
if (pingpong == 0) {
    if (stage == 0) {
        a = imageLoad(pingpong0, ivec2(bit_reverse(idx), column)).rg;
        b = imageLoad(pingpong0, ivec2(bit_reverse(idx + shift), column)).rg;
    }
    else {
        a = imageLoad(pingpong0, ivec2(idx, column)).rg;
        b = imageLoad(pingpong0, ivec2(idx + shift, column)).rg;
    }

    // ... Compute and store results
}
else {
    a = imageLoad(pingpong1, ivec2(idx, column)).rg;
    b = imageLoad(pingpong1, ivec2(idx + shift, column)).rg;

    // ... Compute and store results
}
```

With all this steps aggregated, the shader for the horizontal pass of this FFT DIT Cooley-Tukey implementation is presented in [Listing A.1](#), see [Appendix A](#).

For the vertical pass shader there is However, one extra multiplication by `mult_factor` in the butterfly results for the last stage when the pass is inverse. This corresponds to the normalization of the multidimensional transform similar to the normalization in the inverse DFT in ??, this is demonstrated in [Listing A.2](#).

We can already see in ?? the above code a lot of branching that might be undesirable on the GPU ...

*All stages in one pass*

The previous implementation demonstrates a generic 2D FFT that may be reused for multiple FFT sizes. However, this comes at a cost of efficiency when it comes to the synchronization of the stage itself, moreover it requires use multiple uniform variables for control of the FFT that are transferred between CPU and GPU when there are updates in between stages.

The ideal solution here is to port this implementation synchronization step to be kernel-wise and this detail changes how the code is structured and how it may be dispatched.

At the moment there isn't a way to trivially ensure the reads and writes of all work groups ? without interrupting at the end of the stage, but there are functions that make use of barriers that synchronize all the threads within a work group. However, there is a lot of literature on behalf of GPU compute execution synchronization we'll make use of the GLSL predefined barrier synchronization functions.

The current grouping of threads doesn't allow the use of these barrier synchronization since the computation of one dimensional FFT is distributed between multiple work groups, so using a call to `barrier()` inside the kernel wouldn't fix the race conditions of several segments of the image. We could however, change the setup of these work groups in such way that each 1D FFT threads fit in one work group as illustrated in [Figure 12](#).



**Figure 12:** Difference in invocation spaces for size 256 FFT to allow barrier synchronization between local threads

By restricting the invocation space to be only one dimensional we grant the possibility to use the barrier correctly but at the cost of resizability, the work group local size must now be restricted to half the size of the FFT. This constraint may seem like a penalty on the size of the work group if the FFT size is very high, but this isn't the case for the sizes we tested (see [Figure 16](#) in [Section 6.2](#)), the GPU is an extremely parallel hardware, therefore, it's better to dispatch smaller programs with more instances than overloading local threads in a work group to decrease its.

When fitting all stages in one kernel, the algorithm stays the same but we provide all information needed, such as the size and log size of the FFT for the compiler to unroll the loop, since the for loop feature in GLSL requires to



be used with constant expressions that allow the loop to be inlined in the compiled code.

```
#define FFT_SIZE 256
#define LOG_SIZE 8

layout (local_size_x = FFT_SIZE/2, local_size_y = 1) in;

layout (binding = 0, rg32f) uniform image2D pingpong0;
layout (binding = 1, rg32f) uniform image2D pingpong1;
uniform int fft_dir;
// ... Auxiliar functions

void main() {
    // ...
    int pingpong = 0;

    for(int stage = 0; stage < LOG_SIZE; ++stage) {
        int group_size = 2 << stage;
        int shift = 1 << stage;

        vec2 a, b;
        int idx = (line % shift) + group_size * (line / shift);
        vec2 w = euler(fft_dir * 2 * (M_PI / group_size) * ((idx % group_size) %
shift));
        // ... Perform butterflies

        pingpong = (pingpong + 1) % 2;
        barrier();
    }
}
```

**Listing 5.6:** Unique pass structure for Cooley-Tukey

Aside from the variable `fft_dir` that allows the re usage of this pass for the inverse, all needed data for the FFT lives inside the code, hence, there are more opportunities for optimization when the kernel is compiled since most data is within the code itself. For full code for horizontal and vertical passes see [Appendix A](#).

### 5.3.2 Radix-2 Stockham

Since we wanted to get rid of the bit reversal permutation in the shader code, we implemented the Stockham algorithm (see [Section 4.2](#)). What differs from the previous code is mainly the reordering of the elements when writing the results of the butterfly. It is also important to note that this algorithm is in DIF instead of DIT like the previous ones.

When performing the butterfly in [Listing 5.7](#) we read the elements from the image according to the thread identifier, however we store it in such way that the output has its elements in natural order (see [Algorithm 4](#)).

```
for(int stage = 0; stage < LOG_SIZE; ++stage) {
    int n = 1 << (LOG_SIZE - stage); // group_size
    int m = n >> 1; // shift
    int s = 1 << stage;

    int p = line / s;
    int q = line % s;
    vec2 wp = euler(fft_dir * 2 * (M_PI / n) * p);

    if(pingpong == 0) {
        vec2 a = imageLoad(pingpong0, ivec2(q + s*(p + 0), column)).rg;
        vec2 b = imageLoad(pingpong0, ivec2(q + s*(p + m), column)).rg;

        vec2 res = (a + b);
        imageStore(pingpong1, ivec2(q + s*(2*p + 0), column), vec4(res,0,0));
        res = complex_mult(wp, (a - b));
        imageStore(pingpong1, ivec2(q + s*(2*p + 1), column), vec4(res,0,0));
    }
    else {
        // ... Read pingpong1, write pingpong0
    }

    // ... Sync
}
```

**Listing 5.7:** Radix-2 Stockham DIF

Consequently, we got rid of the conditional read with bit reversed index, and the compute shader code for the horizontal [Listing A.5](#) and vertical passes [Listing A.6](#) got simpler as a result of this.

### 5.3.3 Radix-4 Stockham

In [Section 4.3](#) we introduced how higher radix factorizations could improve the performance with the cost of size constraints, hence, here we change the code with ease to implement the radix-4 factorization described previously.

First we update the stage control variables and compute the multiple twiddle factors used in the radix-4 butterfly. Note that now the for loop only iterates  $\log(N)/2$  times and the number of local threads is reduced by half.

```
#define FFT_SIZE 256
#define LOG_SIZE 8 // log2(FFT_SIZE)
#define HALF_LOG_SIZE 4 // log2(FFT_SIZE) / 2

layout (local_size_x = FFT_SIZE/4, local_size_y = 1) in;
```

```

// ...

void main() {
    // ...

    for(int stage = 0; stage < HALF_LOG_SIZE; ++stage) {
        int n = 1 << (HALF_LOG_SIZE - stage)*2;
        int s = 1 << stage*2;

        int n0 = 0;
        int n1 = n/4;
        int n2 = n/2;
        int n3 = n1 + n2;

        int p = line / s;
        int q = line % s;

        vec2 w1p = euler(2*(M_PI / n) * p * fft_dir);
        vec2 w2p = complex_mult(w1p,w1p);
        vec2 w3p = complex_mult(w1p,w2p);

        // ... Radix-4 butterfly
    }
}

```

**Listing 5.8:** Radix-4 Stockham stage control variables

Then, we compute the radix-4 butterfly and store the results in natural order, as described in [Figure 9](#).

```

if(pingpong == 0) {
    vec2 a = imageLoad(pingpong0, ivec2(q + s*(p + n0), column)).rg;
    vec2 b = imageLoad(pingpong0, ivec2(q + s*(p + n1), column)).rg;
    vec2 c = imageLoad(pingpong0, ivec2(q + s*(p + n2), column)).rg;
    vec2 d = imageLoad(pingpong0, ivec2(q + s*(p + n3), column)).rg;

    vec2 apc = a + c;
    vec2 amc = a - c;
    vec2 bpd = b + d;
    vec2 jbmd = complex_mult(vec2(0,1), b - d);

    imageStore(pingpong1, ivec2(q + s*(4*p + 0), column), vec4(apc + bpd, 0,0));
    imageStore(pingpong1, ivec2(q + s*(4*p + 1), column), vec4(complex_mult(w1p,
    amc - jbmd), 0,0));
    imageStore(pingpong1, ivec2(q + s*(4*p + 2), column), vec4(complex_mult(w2p,
    apc - bpd ), 0,0));
}

```

```

    imageStore(pingpong1, ivec2(q + s*(4*p + 3), column), vec4(complex_mult(w3p,
    amc + jbmd), 0,0));
}
else {
    // ...
}

```

**Listing 5.9:** Radix-4 Stockham butterfly

We take advantage of branchless programming to flip the signal of the butterfly to avoid unnecessary if statements for the inverse.

```

imageStore(pingpong1, ivec2(q + s*(4*p + 0), column), vec4(apc + bpd, 0,0));
imageStore(pingpong1, ivec2(q + s*(4*p + 1), column), vec4(complex_mult(w1p, amc
+ jbmd*fft_dir), 0,0));
imageStore(pingpong1, ivec2(q + s*(4*p + 2), column), vec4(complex_mult(w2p, apc
- bpd ), 0,0));
imageStore(pingpong1, ivec2(q + s*(4*p + 3), column), vec4(complex_mult(w3p, amc
- jbmd*fft_dir), 0,0));

```

**Listing 5.10:** Radix-4 Stockham dragonfly inverse arithmetic

With this changes we now have the radix-4 Stockham shader complete (see [Listing A.7](#) and [Listing A.8](#)).

---

## ANALYSIS AND COMPARISON

---

Finally on this chapter an evaluation of the explored implementations and improvements is provided followed by an empirical analysis based on the results and tests done.

To establish a reference point on the results provided, we used `cuFFT` which is the `CUDA` Fast Fourier Transform library from NVIDIA.

Additionally, to deepen the analysis, this chapter also delivers an equivalent comparison of the researched algorithms applied to a different compute framework such as `CUDA`.

### 6.1 CUFFT

The `cuFFT` library is the NVIDIA framework designed to provide high performance `FFT` exclusively on its own GPUs that supports a wide range of `FFT` inputs and settings that compute `FFT`s efficiently on NVIDIA GPUs.

The `cuFFT` library is acknowledged as one of the most efficient `FFT` GPU framework for the flexibility it provides and it is "*de-facto a standard GPU implementation for developers using CUDA*" (?). Furthermore, it offers all kinds of settings needed for most use cases, such as multidimensional transforms, complex and real-valued input and output, support for half, single and double floating point precision, execution of multiple transforms simultaneously and finally since all this is implemented using `CUDA` we can take advantage of streamed execution, enabling asynchronous computation and data movement.

Unfortunately, as mentioned before the main downside of `cuFFT` is the unavailability of this library for GPUs from other vendors.

The `cuFFT` library uses algorithms highly optimized for input sizes that can be written in the form  $2^a \times 3^b \times 5^c \times 7^d$ , so it factorizes the input size to allow arbitrary sized `FFT` sequences. Furthermore, sizes with lower prime factors have intuitively better performance.

To use the results of the `cuFFT` library as a reference point, we need to establish equivalent conditions to that of the GLSL implementations:

- Out-of-place 2D `FFT`, input buffer is different from the output buffer;
- Base 2 input sizes, such as 128, 256, 512 and 1024;

- Complex to complex FFT, input and output buffer are complex valued;
- The benchmarks are average milliseconds of multiple executions, However, the first dispatch is not taking into account since takes extra time to setup things on the GPU.

The results of the cuFFT out-of-place benchmarks can be found in [Figure 14](#) and [Figure 15](#).

## 6.2 IMPLEMENTATION ANALYSIS IN GLSL

The implementations discussed in [Section 5.3](#), as said before, were studied and benchmarks were made to come to a conclusion about the advantages and disadvantages of using each one and how do they perform. With this in mind we prepared an interactive test environment using Nau 3D engine ? and profiled it using an internal pass profiler.

All benchmarks, the ones in this section and [Section 6.3](#) were tested with the following hardware and software configuration:

- **CPU:** Intel(R) Core(TM) i7-8750H @ 2.20GHz;
- **GPU:** NVIDIA GeForce GTX 1050 Ti Max-Q;
- **NVIDIA driver:** 511.65;
- **CUDA version:** V11.6.124;
- **GLSL version:** 4.60.

In [Section 5.3.1](#) we discussed how the implementation would benefit by having a unique pass that synchronized by stage instead of dispatching multiple stage passes, accordingly we can clearly notice this difference in [Figure 13](#).

The results in [Figure 13](#) show us how a real time application would behave by adopting both strategies. The pass per stage approach introduces a lot of runtime overhead in between stages, since it needs to update the stage for the next iteration dispatch, on the other hand the unique pass kernel is highly optimized for its own size so most calculation are inlined by the GLSL compiler and the synchronization is kept inside the GPU until the kernel is done executing.

As we can see in [Figure 14](#) the GLSL radix-2 implementation of the Stockham algorithm has overall better performance comparing it to the Cooley-Tukey version. This happens due to the removal of the data reordering process of the bit reversal done in the Cooley-Tukey version only on the first stage. Effectively, this change improves the results consistently within the test size ranges, however, it is worth noting that the Stockham algorithm has worst data access locality than the Cooley-Tukey algorithm since it accesses data arbitrarily within the FFT instead of performing the sorting right away, therefore, these results may not hold this way for larger sizes.

Even more important, the results for the radix-4 version of the Stockham algorithm outperform radix-2 as demonstrated in [Figure 15](#), it drastically improved the performance halfway close to the cuFFT.



**Figure 13:** Frame time difference between using stage per pass approach and unique pass for the Radix-2 Cooley-Tukey algorithm

By choosing a radix-4 approach the number of stages reduces to half but with a lot more complex operations per dragonfly on each stage, although the work complexity remains the same, there are much less barrier synchronization events.

For each stage there's a synchronization barrier for each local thread inside a work group, so less stages means less sync points, hence, compensating the 70% kernel size increase of the radix-4 version .

It is also worth noting, in our benchmarks we tested the possibility to integrate multiple butterflies per threads in hope that there was a threshold on the number of local threads that would compensate the workload per kernel for less threads in the synchronization step, however, this didn't happen and for all implementations introducing multiple butterflies per processor and reducing the number of threads. This approach didn't improve the performance overall, so the results of this testing are presented in [Figure 16](#) for the Radix-2 Stockham algorithm. We can conclude that the compute shaders can handle well high threads work groups, therefore less work per kernel is preferred instead of lesser threads per work group.

### 6.3 IMPLEMENTATION ANALYSIS IN CUDA

To study how these implementations behave on other GPGPU frameworks we also provide a study of these implementations using CUDA. Since the objective is to analyze the impact on a different GPGPU API the conditions of the execution of the implementations are identical, therefore, invocation spaces are the same and a CUDA surface object is used as the ping pong buffers for read and write access.

We analyze the performance difference for the Radix-2 Cooley-Tukey ([Listing C.1](#)), Radix-2 Stockham ([Listing C.3](#)) and Radix-4 Stockham (??).

When comparing the Radix-2 CUDA and GLSL implementations plotted in [Figure 17](#), we are able to notice interesting results which are only influenced by the difference in GPGPU platforms. The Cooley-Tukey version in CUDA behaves similarly to the GLSL version however for higher FFT sizes such as 1024 the GLSL implementation outperforms it with a gentle performance increase curve.

The GLSL implementations have slightly better benchmarks and the Radix-2 Stockham version outperforms the implementations written in CUDA.

We can note that the CUDA implementation of Radix-2 Stockham has worst performance when sizes are above 256 than the Cooley-Tukey version which doesn't happen in GLSL.

The sparse read and write access of the ping pong buffers in CUDA implementations have a noticeable impact on the performance of these implementations.

In [Figure 18](#) we have plotted the results for the Radix-2 and Radix-4 Stockham algorithm, were we find that the Radix-4 version have overall better performance.

With the presented results, preferably in GLSL a programmer might want to implement this Radix-2 Stockham algorithm on the GPU, and if the target input buffer has a size multiple of 4, then the Radix-4 version should be used instead. However, this logic doesn't not apply the same way for CUDA implementations.





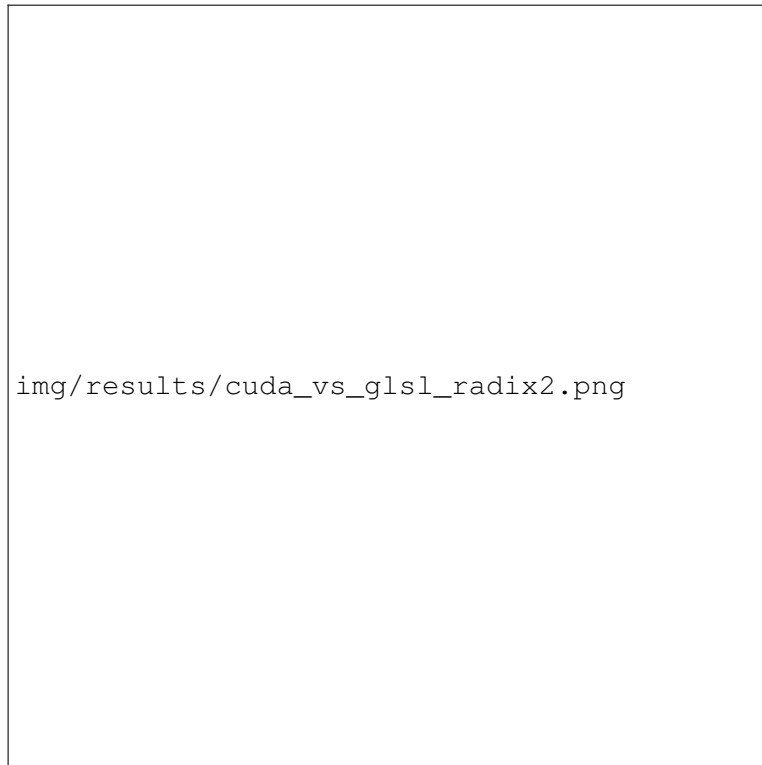
**Figure 14:** 2D FFT benchmarks in milliseconds of out-of-place cuFFT and GLSL Radix-2 algorithms



**Figure 15:** 2D FFT benchmarks of in milliseconds out-of-place cuFFT and GLSL Radix-4 Stockham algorithm



**Figure 16:** 2D FFT benchmarks in milliseconds of GLSL Radix-2 Stockham with multiple butterflies per local threads



**Figure 17:** 2D FFT benchmarks in milliseconds of CUDA and GLSL Radix-2 algorithms



**Figure 18:** 2D FFT benchmarks in milliseconds of CUDA and GLSL Radix-4 Stockham algorithm

## 6.4 CASE OF STUDY

Based on the findings of [Section 6.3](#) we measured how good implementations in GLSL are, by analyzing the performance of a test application that converted a 2D image to frequency domain and then reversed it to its original look. However, with the goal of highlighting the importance of the performance increase of these algorithms, in this section we provide an overview of the implementations impact within a more real scenario by using a ocean rendering technique demo that heavily relies on the usage of FFT ([Figure 19](#)).

In this section we brief the Tensendorf waves demo in [Section 6.4.1](#) where we describe in which way FFTs are relevant for the implementation of this rendering technique, and how we improved an existing implementation for Nau3D that used a pass per stage Cooley-Tukey implementation. After that we present results and how the FFT implementation improves the demo performance and by how much in [Section 6.4.2](#).

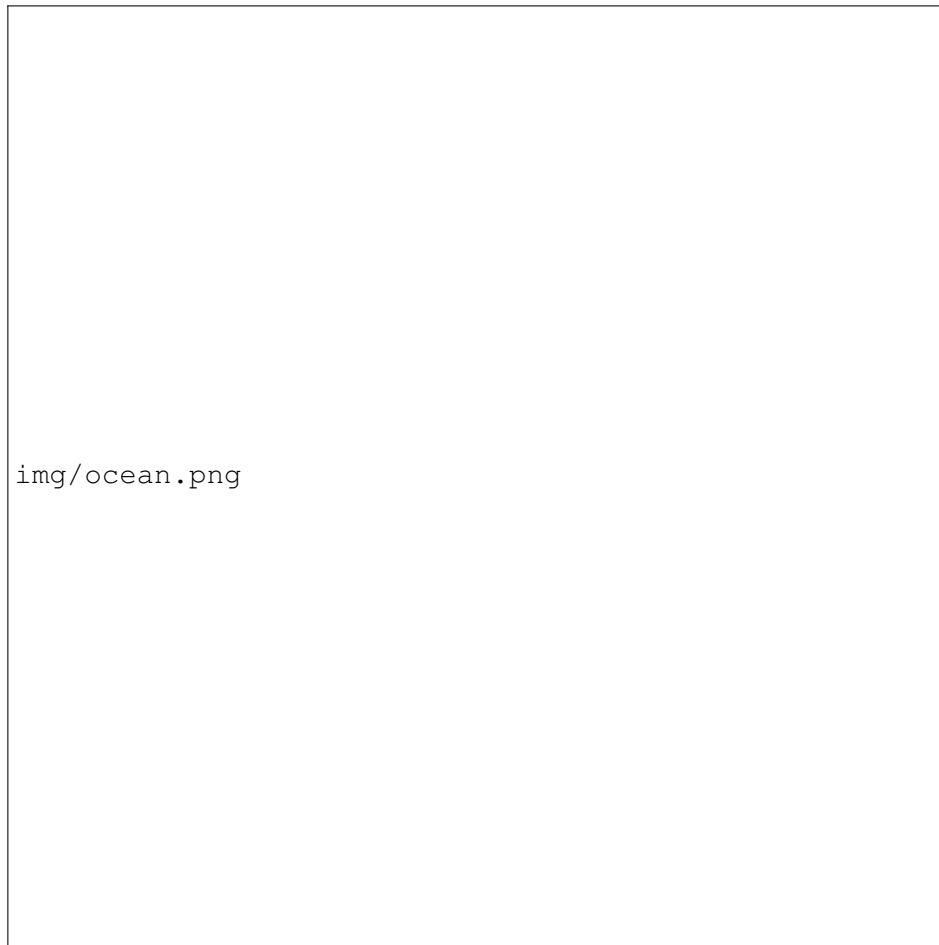
### 6.4.1 Tensendorf waves

The rendering of oceans demo we used as a starting point was a real time implementation in Nau3D of the popular article published by ?. In this demo there are two main stages, the generation of the height map and the actual rendering, the FFTs come into place in the generation of the height map since we need to generate it and the additional vectors used for shading. In total, there are 4 2D FFTs computed each frame, that translates to  $8 * FFT\_SIZE$  1D FFTs in total.

Regarding the results in [Section 6.2](#), we first changed the pipeline from a pass per stage radix-2 Cooley-Tukey algorithm to implement radix-2 Stockham with synchronization within the kernel execution for the horizontal and vertical passes. Each pass computes multiple FFTs at a time and we take advantage of the same kernel to compute all 4 required FFTs at the same time. Within each `image2D` we store 2 complex values in a `vec4` so we also take advantage of SIMD operations optimize the additional cost.

Although the FFTs take a big role in this demo, it also renders the ocean waves that have around 2 million vertices, so the performance does not only depend on the FFTs computation.

For this demo we used 512 as the width for the 2D complex valued buffers, this width allows the waves to have good enough quality, however, ? mentions we can increase the FFT size to 1024 or 2048 for better wave quality if needed. Accordingly, we also tested this demo with width 1024 for more detailed waves and used radix-4 Stockham to achieve best performance.



**Figure 19:** Tensendorf waves in Nau3D Engine



**Figure 20:** Time spent in the CPU and GPU for the size 512 FFT horizontal and vertical passes

#### 6.4.2 Results

In [Figure 20](#) we can note the performance difference the radix-2 Stockham gives when performing the horizontal and vertical pass comparing it to the radix-2 Cooley-Tukey with a pass per stage, overall there is a huge improvement just by removing the stage update phase that is done on the CPU side.

When running the application its notable the performance difference, and the frame rate is improved up to 20%.

Testing the demo for higher quality waves the radix-4 Stockham performance stands out in [Figure 21](#), as predicted. When running the application, the difference in performance is noticeable, with the frame rate using Stockham radix-2 being improved by up to 20%, while the radix-4 delivers a frame rate of up to 60%.

These results proved the necessity of implementing adequate algorithms and appropriate GPGPU programming practices.



**Figure 21:** Total time spent in the CPU and GPU for the size 1024 FFT horizontal and vertical passes

---

## CONCLUSIONS AND FUTURE WORK

---



---

## BIBLIOGRAPHY

---

- Leo Bluestein. A linear filtering approach to the computation of discrete fourier transform. *IEEE Transactions on Audio and Electroacoustics*, 18(4):451–455, 1970.
- E Oran Brigham. *The fast Fourier transform and its applications*. Prentice-Hall, Inc., 1988.
- Eleanor Chu and Alan George. *Inside the FFT black box: serial and parallel fast Fourier transform algorithms*. CRC press, 1999.
- Randima Fernando et al. *GPU gems: programming techniques, tips, and tricks for real-time graphics*, volume 590. Addison-Wesley Reading, 2004.
- Fynn-Jorin Flügge. Realtime gpgpu fft ocean water simulation. Technical report, 2017.
- Matteo Frigo and Steven G Johnson. Fftw: Fastest fourier transform in the west. *Astrophysics Source Code Library*, pages ascl–1201, 2012.
- Naga K Govindaraju, Brandon Lloyd, Yuri Dotsenko, Burton Smith, and John Manferdelli. High performance discrete fourier transforms on graphics processors. In *SC’08: Proceedings of the 2008 ACM/IEEE conference on Supercomputing*, pages 1–12. Ieee, 2008.
- Stijn Heldens, Pieter Hijma, Ben Van Werkhoven, Jason Maassen, and Rob V Van Nieuwpoort. Lightning: Scaling the gpu programming model beyond a single gpu. In *2022 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, pages 492–503. IEEE, 2022.
- Roger A Horn and Charles R Johnson. *Matrix analysis*. Cambridge university press, 2012.
- Waqar Hussain, Fabio Garzia, and Jari Nurmi. Evaluation of radix-2 and radix-4 fft processing on a reconfigurable platform. In *13th IEEE Symposium on Design and Diagnostics of Electronic Circuits and Systems*, pages 249–254. IEEE, 2010.
- Yan-Bin Jia. Polynomial multiplication and fast fourier transform. *Com S*, 477:577, 2014.
- Douglas L Jones. *Digital signal processing: A user’s guide*. 2014.
- Pere Marti-Puig and Ramon Reig Bolano. Radix-4 fft algorithms with ordered input and output data. In *2009 16th International Conference on Digital Signal Processing*, pages 1–6. IEEE, 2009.
- NTiAudio. Fast Fourier Transformation FFT - Basics. <https://www.nti-audio.com/en/support/know-how/fast-fourier-transform-fft>. Accessed at 2022-01-28.

- NVIDIA. Cuda c++ programming guide. <https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#scalable-programming-model>. Accessed: 2022-07-27.
- CUDA Nvidia. Cufft library (2018). URL developer. [nvidia.com/cuFFT](https://developer.nvidia.com/cuFFT).
- OK Ojisan (Takuya OKAHISA). Introduction to the stockham fft. <http://www.pikara.ne.jp/okojisan/otfft-en/stockham1.html>. Accessed: 2022-07-22.
- J Prado. A new fast bit-reversal permutation algorithm based on a symmetry. *IEEE Signal Processing Letters*, 11 (12):933–936, 2004.
- Charles M Rader. Discrete fourier transforms when the number of data samples is prime. *Proceedings of the IEEE*, 56(6):1107–1108, 1968.
- Kamisetty Ramam Rao and Patrick C Yip. *The transform and data compression handbook*. CRC press, 2018.
- RC Singleton. An algorithm for computing the mixed radix fast fourier transform. *IEEE Transactions on audio and electroacoustics*, 17(2):93–103, 1969.
- Julius Orion Smith. *Mathematics of the discrete Fourier transform (DFT): with audio applications*. Julius Smith, 2007.
- R Yavne. An economical method for calculating the discrete fourier transform. In *Proceedings of the December 9-11, 1968, fall joint computer conference, part I*, pages 115–125, 1968.

Part I

APPENDICES



---

## GLSL FFT

---

```
#version 440

#define M_PI 3.1415926535897932384626433832795

layout (local_size_x = 4, local_size_y = 8) in;

layout (binding = 0, rg32f) uniform image2D pingpong0;
layout (binding = 0, rg32f) uniform image2D pingpong1;

uniform int pingpong;
uniform int log_width;
uniform int stage;
uniform int fft_dir;

vec2 complex_mult(vec2 v0, vec2 v1) {
    return vec2(v0.x * v1.x - v0.y * v1.y,
                v0.x * v1.y + v0.y * v1.x);
}

int bit_reverse(int k) {
    uint br = bitfieldReverse(k);
    return int(bitfieldExtract(br, 32 - log_width, log_width));
}

vec2 euler(float angle) {
    return vec2(cos(angle), sin(angle));
}

void main() {
    int line = int(gl_GlobalInvocationID.x);
    int column = int(gl_GlobalInvocationID.y);

    int group_size = 2 << stage;
    int shift = 1 << stage;
```

```

vec2 a, b;

    int idx = (line % shift) + group_size * (line / shift);
    vec2 w = euler(fft_dir * 2 * (M_PI / group_size) * ((idx % group_size) %
shift));

    if (pingpong == 0) {
        if (stage == 0) {
            a = imageLoad(pingpong0, ivec2(bit_reverse(idx), column)).rg;
            b = imageLoad(pingpong0, ivec2(bit_reverse(idx + shift), column)).rg
;
        }
        else {
            a = imageLoad(pingpong0, ivec2(idx, column)).rg;
            b = imageLoad(pingpong0, ivec2(idx + shift, column)).rg;
        }

        vec2 raux = a + complex_mult(w, b);
        imageStore(pingpong1, ivec2(idx, column), vec4(raux, 0, 0));

        raux = a - complex_mult(w, b);
        imageStore(pingpong1, ivec2(idx + shift, column), vec4(raux, 0, 0));
    }
    else {
        a = imageLoad(pingpong1, ivec2(idx, column)).rg;
        b = imageLoad(pingpong1, ivec2(idx + shift, column)).rg;

        vec2 raux = a + complex_mult(w, b);
        imageStore(pingpong0, ivec2(idx, column), vec4(raux,0,0));

        raux = a - complex_mult(w, b);
        imageStore(pingpong0, ivec2(idx + shift, column), vec4(raux,0,0));
    }
}

```

**Listing A.1:** FFT Radix-2 Cooley-Tukey Horizontal stage pass, see [Section 5.3.1](#)

```

#version 440

#define M_PI 3.1415926535897932384626433832795

layout (local_size_x = 8, local_size_y = 4) in;

layout (binding = 0, rg32f) uniform image2D pingpong0;
layout (binding = 1, rg32f) uniform image2D pingpong1;

uniform int pingpong;

```

```

uniform int log_width;
uniform int stage;
uniform int fft_dir;

int iter = 1 << log_width;
int shift = (1 << stage);

vec2 complex_mult(vec2 v0, vec2 v1) {
    return vec2(v0.x * v1.x - v0.y * v1.y,
                v0.x * v1.y + v0.y * v1.x);
}

int bit_reverse(int k) {
    uint br = bitfieldReverse(k);
    return int(bitfieldExtract(br, 32 - log_width, log_width));
}

vec2 euler(float angle) {
    return vec2(cos(angle), sin(angle));
}

void main() {
    int line = int(gl_GlobalInvocationID.x);
    int column = int(gl_GlobalInvocationID.y);

    int group_size = 2 << stage;
    int shift = 1 << stage;

    vec2 a, b;

    int idx = (column % shift) + group_size * (column / shift);
    vec2 w = euler(fft_dir * 2 * (M_PI / group_size) * ((idx % group_size) %
    shift));

    float mult_factor = 1.0;
    if ((stage == log_width - 1) && fft_dir == 1) {
        mult_factor = 1.0 / (iter*iter);
    }

    if (pingpong == 0) {
        if (stage == 0) {
            a = imageLoad(pingpong0, ivec2(line, bit_reverse(idx))).rg;
            b = imageLoad(pingpong0, ivec2(line, bit_reverse(idx + shift))).rg;
        }
        else {
            a = imageLoad(pingpong0, ivec2(line, idx)).rg;
            b = imageLoad(pingpong0, ivec2(line, idx + shift)).rg;
        }
    }
}

```

```

    }

    vec2 raux = (a + complex_mult(w, b)) * mult_factor;
    imageStore(pingpong1, ivec2(line, idx), vec4(raux,0,0));

    raux = (a - complex_mult(w, b)) * mult_factor;
    imageStore(pingpong1, ivec2(line, idx + shift), vec4(raux,0,0));
}
else {
    if (stage == 0) {
        a = imageLoad(pingpong1, ivec2(line, bit_reverse(idx))).rg;
        b = imageLoad(pingpong1, ivec2(line, bit_reverse(idx + shift))).rg;
    }
    else {
        a = imageLoad(pingpong1, ivec2(line, idx)).rg;
        b = imageLoad(pingpong1, ivec2(line, idx + shift)).rg;
    }

    vec2 raux = (a + complex_mult(w, b)) * mult_factor;
    imageStore(pingpong0, ivec2(line, idx), vec4(raux,0,0));

    raux = (a - complex_mult(w, b)) * mult_factor;
    imageStore(pingpong0, ivec2(line, idx + shift), vec4(raux,0,0));
}
}

```

**Listing A.2:** FFT Radix-2 Cooley-Tukey Vertical stage pass, see [Section 5.3.1](#)

```

#version 440

#define M_PI 3.1415926535897932384626433832795
#define FFT_SIZE 256
#define LOG_SIZE 8

layout (local_size_x = FFT_SIZE/2, local_size_y = 1) in;

layout (binding = 0, rg32f) uniform image2D pingpong0;
layout (binding = 1, rg32f) uniform image2D pingpong1;

uniform int fft_dir;

vec2 complex_mult(vec2 v0, vec2 v1) {
    return vec2(v0.x * v1.x - v0.y * v1.y,
               v0.x * v1.y + v0.y * v1.x);
}

int bit_reverse(int k) {

```

```

uint br = bitfieldReverse(k);
return int(bitfieldExtract(br, 32 - LOG_SIZE, LOG_SIZE));
}

vec2 euler(float angle) {
    return vec2(cos(angle), sin(angle));
}

void main() {
    int line = int(gl_GlobalInvocationID.x);
    int column = int(gl_WorkGroupID.y);
    int pingpong = 0;

    for(int stage = 0; stage < LOG_SIZE; ++stage) {
        int group_size = 2 << stage;
        int shift = 1 << stage;

        vec2 a, b;
        int idx = (line % shift) + group_size * (line / shift);
        vec2 w = euler(fft_dir * 2 * (M_PI / group_size) * ((idx % group_size) %
shift));

        // alternate between textures
        if (pingpong == 0) {
            if (stage == 0) {
                a = imageLoad(pingpong0, ivec2(bit_reverse(idx), column)).rg;
                b = imageLoad(pingpong0, ivec2(bit_reverse(idx + shift), column))
.rg;
            }
            else {
                a = imageLoad(pingpong0, ivec2(idx, column)).rg;
                b = imageLoad(pingpong0, ivec2(idx + shift, column)).rg;
            }

            vec2 raux = a + complex_mult(w, b);
            imageStore(pingpong1, ivec2(idx, column), vec4(raux, 0, 0));
            raux = a - complex_mult(w, b);
            imageStore(pingpong1, ivec2(idx + shift, column), vec4(raux, 0, 0));
        }
        else {
            a = imageLoad(pingpong1, ivec2(idx, column)).rg;
            b = imageLoad(pingpong1, ivec2(idx + shift, column)).rg;

            vec2 raux = a + complex_mult(w, b);
            imageStore(pingpong0, ivec2(idx, column), vec4(raux,0,0));
            raux = a - complex_mult(w, b);
            imageStore(pingpong0, ivec2(idx + shift, column), vec4(raux,0,0));
        }
    }
}

```



```

    }

    pingpong = (pingpong + 1) % 2;
    barrier();
}
}

```

**Listing A.3:** FFT Radix-2 Cooley-Tukey Horizontal unique pass, see [Section 5.3.1](#)

```

#version 440

#define M_PI 3.1415926535897932384626433832795
#define FFT_SIZE 256
#define LOG_SIZE 8

layout (local_size_x = FFT_SIZE/2, local_size_y = 1) in;

layout (binding = 0, rg32f) uniform image2D pingpong0;
layout (binding = 1, rg32f) uniform image2D pingpong1;

uniform int fft_dir;

vec2 complex_mult(vec2 v0, vec2 v1) {
    return vec2(v0.x * v1.x - v0.y * v1.y,
               v0.x * v1.y + v0.y * v1.x);
}

int bit_reverse(int k) {
    uint br = bitfieldReverse(k);
    return int(bitfieldExtract(br, 32 - LOG_SIZE, LOG_SIZE));
}

vec2 euler(float angle) {
    return vec2(cos(angle), sin(angle));
}

void main() {
    int line = int(gl_WorkGroupID.y);
    int column = int(gl_GlobalInvocationID.x);
    int pingpong = LOG_SIZE % 2;

    for(int stage = 0; stage < LOG_SIZE; ++stage) {
        int group_size = 2 << stage;
        int shift = 1 << stage;

        vec2 a, b;
        int idx = (column % shift) + group_size * (column / shift);
    }
}

```

```

    vec2 w = euler(fft_dir * 2 * (M_PI / group_size) * ((idx % group_size) %
shift));

    float mult_factor = 1.0;
    if ((stage == LOG_SIZE - 1) && fft_dir == 1) {
        mult_factor = 1.0 / (FFT_SIZE*FFT_SIZE);
    }

    if (pingpong == 0) {
        if (stage == 0) {
            a = imageLoad(pingpong0, ivec2(line, bit_reverse(idx))).rg;
            b = imageLoad(pingpong0, ivec2(line, bit_reverse(idx + shift))).
rg;
        }
        else {
            a = imageLoad(pingpong0, ivec2(line, idx)).rg;
            b = imageLoad(pingpong0, ivec2(line, idx + shift)).rg;
        }

        vec2 raux = (a + complex_mult(w, b)) * mult_factor;
        imageStore(pingpong1, ivec2(line, idx), vec4(raux,0,0));
        raux = (a - complex_mult(w, b)) * mult_factor;
        imageStore(pingpong1, ivec2(line, idx + shift), vec4(raux,0,0));

    }
    else {
        if (stage == 0) {
            a = imageLoad(pingpong1, ivec2(line, bit_reverse(idx))).rg;
            b = imageLoad(pingpong1, ivec2(line, bit_reverse(idx + shift))).
rg;
        }
        else {
            a = imageLoad(pingpong1, ivec2(line, idx)).rg;
            b = imageLoad(pingpong1, ivec2(line, idx + shift)).rg;
        }

        vec2 raux = (a + complex_mult(w, b)) * mult_factor;
        imageStore(pingpong0, ivec2(line, idx), vec4(raux,0,0));
        raux = (a - complex_mult(w, b)) * mult_factor;
        imageStore(pingpong0, ivec2(line, idx + shift), vec4(raux,0,0));

    }

    pingpong = ((pingpong + 1) % 2);
    barrier();
}
}

```

**Listing A.4:** FFT Radix-2 Cooley-Tukey Vertical unique pass, see [Section 5.3.1](#)

```

#version 440

#define M_PI 3.1415926535897932384626433832795
#define FFT_SIZE 256
#define LOG_SIZE 8

layout (local_size_x = (FFT_SIZE/2)/NUM_BUTTERFLIES, local_size_y = 1) in;

layout (binding = 0, rg32f) uniform image2D pingpong0;
layout (binding = 1, rg32f) uniform image2D pingpong1;

uniform int fft_dir;

vec2 complex_mult(vec2 v0, vec2 v1) {
    return vec2(v0.x * v1.x - v0.y * v1.y,
                v0.x * v1.y + v0.y * v1.x);
}

vec2 euler(float angle) {
    return vec2(cos(angle), sin(angle));
}

void main() {
    int line = int(gl_GlobalInvocationID.x);
    int column = int(gl_WorkGroupID.y);
    int pingpong = 0;

    for(int stage = 0; stage < LOG_SIZE; ++stage) {
        int n = 1 << (LOG_SIZE - stage);
        int m = n >> 1;
        int s = 1 << stage;

        int p = line / s;
        int q = line % s;

        vec2 wp = euler(fft_dir * 2 * (M_PI / n) * p);
        if(pingpong == 0) {
            vec2 a = imageLoad(pingpong0, ivec2(q + s*(p + 0), column)).rg;
            vec2 b = imageLoad(pingpong0, ivec2(q + s*(p + m), column)).rg;

            vec2 res = (a + b);
            imageStore(pingpong1, ivec2(q + s*(2*p + 0), column), vec4(res,0,0));
            res = complex_mult(wp, (a - b));

```

```

        imageStore(pingpong1, ivec2(q + s*(2*p + 1), column), vec4(res,0,0));
    }
    else {
        vec2 a = imageLoad(pingpong1, ivec2(q + s*(p + 0), column)).rg;
        vec2 b = imageLoad(pingpong1, ivec2(q + s*(p + m), column)).rg;

        vec2 res = (a + b);
        imageStore(pingpong0, ivec2(q + s*(2*p + 0), column), vec4(res,0,0));
        res = complex_mult(wp, (a - b));
        imageStore(pingpong0, ivec2(q + s*(2*p + 1), column), vec4(res,0,0));
    }

    pingpong = (pingpong + 1) % 2;
    barrier();
}
}

```

**Listing A.5:** FFT Radix-2 Stockham Horizontal unique pass, see [Section 5.3.2](#)

```

#version 440

#define M_PI 3.1415926535897932384626433832795
#define FFT_SIZE 256
#define LOG_SIZE 8

layout (local_size_x = FFT_SIZE/2, local_size_y = 1) in;

layout (binding = 0, rg32f) uniform image2D pingpong0;
layout (binding = 1, rg32f) uniform image2D pingpong1;

uniform int fft_dir;

vec2 complex_mult(vec2 v0, vec2 v1) {
    return vec2(v0.x * v1.x - v0.y * v1.y,
               v0.x * v1.y + v0.y * v1.x);
}

vec2 euler(float angle) {
    return vec2(cos(angle), sin(angle));
}

void main() {
    int line = int(gl_WorkGroupID.y);
    int column = int(gl_GlobalInvocationID.x);
    int pingpong = LOG_SIZE % 2;

    for(int stage = 0; stage < LOG_SIZE; ++stage) {

```

```

    int n = 1 << (LOG_SIZE - stage);
    int m = n >> 1;
    int s = 1 << stage;

    float mult_factor = 1.0;
    if ((stage == LOG_SIZE-1) && fft_dir == 1) {
        mult_factor = 1.0 / (FFT_SIZE*FFT_SIZE) ;
    }

    int p = column / s;
    int q = column % s;

    vec2 wp = euler(fft_dir * 2 * (M_PI / n) * p);
    if(pingpong == 0) {
        vec2 a = imageLoad(pingpong0, ivec2(line, q + s*(p + 0))).rg;
        vec2 b = imageLoad(pingpong0, ivec2(line, q + s*(p + m))).rg;

        vec2 res = (a + b) * mult_factor;
        imageStore(pingpong1, ivec2(line, q + s*(2*p + 0)), vec4(res,0,0));
        res = complex_mult(wp, (a - b)) * mult_factor;
        imageStore(pingpong1, ivec2(line, q + s*(2*p + 1)), vec4(res,0,0));
    }
    else {
        vec2 a = imageLoad(pingpong1, ivec2(line, q + s*(p + 0))).rg;
        vec2 b = imageLoad(pingpong1, ivec2(line, q + s*(p + m))).rg;

        vec2 res = (a + b) * mult_factor;
        imageStore(pingpong0, ivec2(line, q + s*(2*p + 0)), vec4(res,0,0));
        res = complex_mult(wp, (a - b)) * mult_factor;
        imageStore(pingpong0, ivec2(line, q + s*(2*p + 1)), vec4(res,0,0));
    }

    pingpong = (pingpong + 1) % 2;
    barrier();
}
}

```

**Listing A.6:** FFT Radix-2 Stockham Vertical unique pass, see [Section 5.3.2](#)

```

#version 440

#define M_PI 3.1415926535897932384626433832795
#define FFT_SIZE 256
#define LOG_SIZE 8 // log2(FFT_SIZE)
#define HALF_LOG_SIZE 4 // log2(FFT_SIZE) / 2

layout (local_size_x = FFT_SIZE/4, local_size_y = 1) in;

```

```

layout (binding = 0, rg32f) uniform image2D pingpong0;
layout (binding = 1, rg32f) uniform image2D pingpong1;

uniform int fft_dir;

vec2 complex_mult(vec2 v0, vec2 v1) {
    return vec2(v0.x * v1.x - v0.y * v1.y,
                v0.x * v1.y + v0.y * v1.x);
}

vec2 euler(float angle) {
    return vec2(cos(angle), sin(angle));
}

void main() {
    int line = int(gl_GlobalInvocationID.x);
    int column = int(gl_WorkGroupID.y);
    int pingpong = 0;

    for(int stage = 0; stage < HALF_LOG_SIZE; ++stage) {
        int n = 1 << (HALF_LOG_SIZE - stage)*2;
        int s = 1 << stage*2;

        int n0 = 0;
        int n1 = n/4;
        int n2 = n/2;
        int n3 = n1 + n2;

        int p = line / s;
        int q = line % s;

        vec2 w1p = euler(2*(M_PI / n) * p * fft_dir);
        vec2 w2p = complex_mult(w1p, w1p);
        vec2 w3p = complex_mult(w1p, w2p);

        if(pingpong == 0) {
            vec2 a = imageLoad(pingpong0, ivec2(q + s*(p + n0), column)).rg;
            vec2 b = imageLoad(pingpong0, ivec2(q + s*(p + n1), column)).rg;
            vec2 c = imageLoad(pingpong0, ivec2(q + s*(p + n2), column)).rg;
            vec2 d = imageLoad(pingpong0, ivec2(q + s*(p + n3), column)).rg;

            vec2 apc = a + c;
            vec2 amc = a - c;
            vec2 bpd = b + d;
            vec2 jbmd = complex_mult(vec2(0,1), b - d);

```

```

        imageStore(pingpong1, ivec2(q + s*(4*p + 0), column), vec4(apc + bpd,
0,0));
        imageStore(pingpong1, ivec2(q + s*(4*p + 1), column), vec4(
complex_mult(w1p, amc + jbmd*fft_dir), 0,0));
        imageStore(pingpong1, ivec2(q + s*(4*p + 2), column), vec4(
complex_mult(w2p, apc - bpd ), 0,0));
        imageStore(pingpong1, ivec2(q + s*(4*p + 3), column), vec4(
complex_mult(w3p, amc - jbmd*fft_dir), 0,0));
    }
    else {
        vec2 a = imageLoad(pingpong1, ivec2(q + s*(p + n0), column)).rg;
        vec2 b = imageLoad(pingpong1, ivec2(q + s*(p + n1), column)).rg;
        vec2 c = imageLoad(pingpong1, ivec2(q + s*(p + n2), column)).rg;
        vec2 d = imageLoad(pingpong1, ivec2(q + s*(p + n3), column)).rg;

        vec2 apc = a + c;
        vec2 amc = a - c;
        vec2 bpd = b + d;
        vec2 jbmd = complex_mult(vec2(0,1), b - d);

        imageStore(pingpong0, ivec2(q + s*(4*p + 0), column), vec4(apc + bpd,
0,0));
        imageStore(pingpong0, ivec2(q + s*(4*p + 1), column), vec4(
complex_mult(w1p, amc + jbmd*fft_dir), 0,0));
        imageStore(pingpong0, ivec2(q + s*(4*p + 2), column), vec4(
complex_mult(w2p, apc - bpd ), 0,0));
        imageStore(pingpong0, ivec2(q + s*(4*p + 3), column), vec4(
complex_mult(w3p, amc - jbmd*fft_dir), 0,0));
    }

    pingpong = (pingpong + 1) % 2;
    barrier();
}
}

```

**Listing A.7:** FFT Radix-4 Stockham Horizontal unique pass, see [Section 5.3.3](#)

```

#version 440

#define M_PI 3.1415926535897932384626433832795
#define FFT_SIZE 256
#define LOG_SIZE 8 // log2(FFT_SIZE)
#define HALF_LOG_SIZE 4 // log2(FFT_SIZE) / 2

layout (local_size_x = (FFT_SIZE/4)/NUM_BUTTERFLIES, local_size_y = 1) in;

```

```

layout (binding = 0, rg32f) uniform image2D pingpong0;
layout (binding = 1, rg32f) uniform image2D pingpong1;

uniform int fft_dir;

vec2 complex_mult(vec2 v0, vec2 v1) {
    return vec2(v0.x * v1.x - v0.y * v1.y,
               v0.x * v1.y + v0.y * v1.x);
}

vec2 euler(float angle) {
    return vec2(cos(angle), sin(angle));
}

void main() {
    int line = int(gl_WorkGroupID.y);
    int column = int(gl_GlobalInvocationID.x);
    int pingpong = HALF_LOG_SIZE % 2;

    for(int stage = 0; stage < HALF_LOG_SIZE; ++stage) {
        int group_size = 2 << stage;
        int shift = 1 << stage;

        int n = 1 << ((HALF_LOG_SIZE - stage)*2);
        int s = 1 << (stage*2);

        int n0 = 0;
        int n1 = n/4;
        int n2 = n/2;
        int n3 = n1 + n2;

        float mult_factor = 1.0;
        if((stage == HALF_LOG_SIZE - 1) && fft_dir == 1) {
            mult_factor = 1.0 / (FFT_SIZE*FFT_SIZE) ;
        }

        int p = column / s;
        int q = column % s;

        vec2 w1p = euler(2*(M_PI / n) * p * fft_dir);
        vec2 w2p = complex_mult(w1p,w1p);
        vec2 w3p = complex_mult(w1p,w2p);

        if(pingpong == 0) {
            vec2 a = imageLoad(pingpong0, ivec2(line, q + s*(p + n0))).rg;
            vec2 b = imageLoad(pingpong0, ivec2(line, q + s*(p + n1))).rg;
            vec2 c = imageLoad(pingpong0, ivec2(line, q + s*(p + n2))).rg;

```



```

        vec2 d = imageLoad(pingpong0, ivec2(line, q + s*(p + n3))).rg;

        vec2 apc = a + c;
        vec2 amc = a - c;
        vec2 bpd = b + d;
        vec2 jbmd = complex_mult(vec2(0,1), b - d);

        imageStore(pingpong1, ivec2(line, q + s*(4*p + 0)), vec4(mult_factor
* (apc + bpd), 0,0));
        imageStore(pingpong1, ivec2(line, q + s*(4*p + 1)), vec4(mult_factor
* (complex_mult(w1p, amc + jbmd*fft_dir), 0,0));
        imageStore(pingpong1, ivec2(line, q + s*(4*p + 2)), vec4(mult_factor
* (complex_mult(w2p, apc - bpd ), 0,0));
        imageStore(pingpong1, ivec2(line, q + s*(4*p + 3)), vec4(mult_factor
* (complex_mult(w3p, amc - jbmd*fft_dir), 0,0));
    }
    else {
        vec2 a = imageLoad(pingpong1, ivec2(line, q + s*(p + n0))).rg;
        vec2 b = imageLoad(pingpong1, ivec2(line, q + s*(p + n1))).rg;
        vec2 c = imageLoad(pingpong1, ivec2(line, q + s*(p + n2))).rg;
        vec2 d = imageLoad(pingpong1, ivec2(line, q + s*(p + n3))).rg;

        vec2 apc = a + c;
        vec2 amc = a - c;
        vec2 bpd = b + d;
        vec2 jbmd = complex_mult(vec2(0,1), b - d);

        imageStore(pingpong0, ivec2(line, q + s*(4*p + 0)), vec4(mult_factor
* (apc + bpd), 0,0));
        imageStore(pingpong0, ivec2(line, q + s*(4*p + 1)), vec4(mult_factor
* (complex_mult(w1p, amc + jbmd*fft_dir), 0,0));
        imageStore(pingpong0, ivec2(line, q + s*(4*p + 2)), vec4(mult_factor
* (complex_mult(w2p, apc - bpd ), 0,0));
        imageStore(pingpong0, ivec2(line, q + s*(4*p + 3)), vec4(mult_factor
* (complex_mult(w3p, amc - jbmd*fft_dir), 0,0));
    }

    pingpong = (pingpong + 1) % 2;
    barrier();
}
}

```

**Listing A.8:** FFT Radix-4 Stockham Vertical unique pass, see [Section 5.3.3](#)

---

## CUFFT

---

```
#include <stdio>
#include <cufft.h>
#include <cuda.h>

#define FFT_SIZE 256

#define CU_ERR_CHECK_MSG(err, msg) { \
    if(err != cudaSuccess) { \
        fprintf(stderr, msg); \
        exit(1); \
    } \
}

#define CU_CHECK_MSG(res, msg) { \
    if(res != CUFFT_SUCCESS) { \
        fprintf(stderr, msg); \
        exit(1); \
    } \
}

int main() {
    const size_t data_size = sizeof(cufftComplex)*FFT_SIZE*FFT_SIZE;
    cufftComplex* data = reinterpret_cast<cufftComplex*>(malloc(data_size));
    cufftComplex* gpu_data_in;
    cufftComplex* gpu_data_out;
    cudaError_t err;

    // Initializing input sequence
    for(size_t i = 0; i < FFT_SIZE*FFT_SIZE; ++i) {
        data[i].x = i;
        data[i].y = 0.00;
    }

    // Allocate Input GPU buffer
    err = cudaMalloc(&gpu_data_in, data_size);
```

```

CU_ERR_CHECK_MSG(err, "Cuda error: Failed to allocate\n");

// Allocate Output GPU buffer
err = cudaMalloc(&gpu_data_out, data_size);
CU_ERR_CHECK_MSG(err, "Cuda error: Failed to allocate\n");

// Copy data to GPU buffer
err = cudaMemcpy(gpu_data_in, data, data_size, cudaMemcpyHostToDevice);
CU_ERR_CHECK_MSG(err, "Cuda error: Failed to copy buffer to GPU\n");

// Setup cufft plan
cufftHandle plan;
cufftResult_t res;
res = cufftPlan2d(&plan, FFT_SIZE, FFT_SIZE, CUFFT_C2C);
CU_CHECK_MSG(res, "cuFFT error: Plan creation failed\n");

// Execute Forward 2D FFT
res = cufftExecC2C(plan, gpu_data_in, gpu_data_out, CUFFT_FORWARD);
CU_CHECK_MSG(res, "cuFFT error: ExecC2C Forward failed\n");

// Await end of execution
err = cudaDeviceSynchronize();
CU_ERR_CHECK_MSG(err, "Cuda error: Failed to synchronize\n");

// Execute Inverse 2D FFT
res = cufftExecC2C(plan, gpu_data_in, gpu_data_out, CUFFT_FORWARD);
CU_CHECK_MSG(res, "cuFFT error: ExecC2C Forward failed\n");

// Await end of execution
err = cudaDeviceSynchronize();
CU_ERR_CHECK_MSG(err, "Cuda error: Failed to synchronize\n");

// Retrieve computed FFT buffer
err = cudaMemcpy(data, gpu_data_in, data_size, cudaMemcpyDeviceToHost);
CU_ERR_CHECK_MSG(err, "Cuda error: Failed to copy buffer to GPU\n");

// Destroy Cuda and cuFFT resources
cufftDestroy(plan);
cudaFree(gpu_data_in);

return 0;
}

```

**Listing B.1:** cuFFT, see [Section 6.1](#)

---

## CUDA FFT

---

```
__device__ __forceinline__
float2 complex_mult(float2 v0, float2 v1) {
    return float2{
        v0.x * v1.x - v0.y * v1.y,
        v0.x * v1.y + v0.y * v1.x
    };
}

__device__ __forceinline__
float2 complex_add(float2 v0, float2 v1) {
    return float2{v0.x + v1.x, v0.y + v1.y };
}

__device__ __forceinline__
float2 complex_sub(float2 v0, float2 v1) {
    return float2{v0.x - v1.x, v0.y - v1.y };
}

__device__ __forceinline__
float2 euler(float angle) {
    return float2{cos(angle), sin(angle)};
}

__device__ __forceinline__
unsigned int bit_reverse(unsigned int k) {
    return __brev(k) >> (32-LOG_SIZE);
}

__global__
void stockham_fft_horizontal(cudaSurfaceObject_t pingpong0, cudaSurfaceObject_t
    pingpong1, float fft_dir) {
    int line = threadIdx.x;
    int column = blockIdx.x;
    int pingpong = 0;

    for(int stage = 0; stage < LOG_SIZE; ++stage) {
```

```

    int group_size = 2 << stage;
    int shift = 1 << stage;

    float2 a, b;
    int idx = (line % shift) + group_size * (line / shift);
    float2 w = euler(fft_dir * 2 * (M_PI / group_size) * ((idx % group_size)
% shift));

    if(pingpong == 0) {
        if(stage == 0) {
            a = surf2Dread<float2>(pingpong0, bit_reverse(idx)*sizeof(float2)
, column);
            b = surf2Dread<float2>(pingpong0, bit_reverse(idx+shift)*sizeof(
float2), column);
        }
        else {
            a = surf2Dread<float2>(pingpong0, (idx)*sizeof(float2), column);
            b = surf2Dread<float2>(pingpong0, (idx+shift)*sizeof(float2),
column);
        }

        surf2Dwrite<float2>(a + complex_mult(w, b), pingpong1, (idx)*sizeof(
float2), column);
        surf2Dwrite<float2>(a - complex_mult(w, b), pingpong1, (idx+shift)*
sizeof(float2), column);
    }
    else {
        if(stage == 0) {
            a = surf2Dread<float2>(pingpong1, bit_reverse(idx)*sizeof(float2)
, column);
            b = surf2Dread<float2>(pingpong1, bit_reverse(idx+shift)*sizeof(
float2), column);
        }
        else {
            a = surf2Dread<float2>(pingpong1, (idx)*sizeof(float2), column);
            b = surf2Dread<float2>(pingpong1, (idx+shift)*sizeof(float2),
column);
        }

        surf2Dwrite<float2>(a + complex_mult(w, b), pingpong0, (idx)*sizeof(
float2), column);
        surf2Dwrite<float2>(a - complex_mult(w, b), pingpong0, (idx+shift)*
sizeof(float2), column);
    }

    pingpong = ((pingpong + 1) % 2);

```

```

        __syncthreads();
    }
}

__global__
void stockham_fft_vertical(cudaSurfaceObject_t pingpong0, cudaSurfaceObject_t
pingpong1, float fft_dir) {
    int line = blockIdx.x;
    int column = threadIdx.x;
    int pingpong = LOG_SIZE % 2;

    for(int stage = 0; stage < LOG_SIZE; ++stage) {
        int group_size = 2 << stage;
        int shift = 1 << stage;

        float2 a, b;
        int idx = (column % shift) + group_size * (column / shift);
        float2 w = euler(fft_dir * 2 * (M_PI / group_size) * ((idx % group_size)
% shift));

        float mult_factor = 1.0;
        if ((stage == LOG_SIZE - 1) && fft_dir == 1) {
            mult_factor = 1.0 / (FFT_SIZE*FFT_SIZE) ;
        }

        if(pingpong == 0) {
            if(stage == 0) {
                a = surf2Dread<float2>(pingpong0, line*sizeof(float2),
bit_reverse(idx));
                b = surf2Dread<float2>(pingpong0, line*sizeof(float2),
bit_reverse(idx+shift));
            }
            else {
                a = surf2Dread<float2>(pingpong0, line*sizeof(float2), idx);
                b = surf2Dread<float2>(pingpong0, line*sizeof(float2), idx+shift)
;
            }

            surf2Dwrite<float2>((a + complex_mult(w, b)) * mult_factor, pingpong1
, line*sizeof(float2), idx);
            surf2Dwrite<float2>((a - complex_mult(w, b)) * mult_factor, pingpong1
, line*sizeof(float2), idx+shift);
        }
        else {
            if(stage == 0) {

```

```

        a = surf2Dread<float2>(pingpong1, line*sizeof(float2),
bit_reverse(idx));
        b = surf2Dread<float2>(pingpong1, line*sizeof(float2),
bit_reverse(idx+shift));
    }
    else {
        a = surf2Dread<float2>(pingpong1, line*sizeof(float2), idx);
        b = surf2Dread<float2>(pingpong1, line*sizeof(float2), idx+shift)
;
    }

    surf2Dwrite<float2>((a + complex_mult(w, b)) * mult_factor, pingpong0
, line*sizeof(float2), idx);
    surf2Dwrite<float2>((a - complex_mult(w, b)) * mult_factor, pingpong0
, line*sizeof(float2), idx+shift);
    }

    pingpong = ((pingpong + 1) % 2);

    __syncthreads();
}
}

```

**Listing C.1:** FFT Radix-2 Cooley-Tukey, see [Section 6.3](#)

```

__device__ __forceinline__
float2 complex_mult(float2 v0, float2 v1) {
    return float2{
        v0.x * v1.x - v0.y * v1.y,
        v0.x * v1.y + v0.y * v1.x
    };
}

__device__ __forceinline__
float2 complex_add(float2 v0, float2 v1) {
    return float2{v0.x + v1.x, v0.y + v1.y };
}

__device__ __forceinline__
float2 complex_sub(float2 v0, float2 v1) {
    return float2{v0.x - v1.x, v0.y - v1.y };
}

__device__ __forceinline__
float2 euler(float angle) {
    return float2{cos(angle), sin(angle)};
}

```

```

__global__
void stockham_fft_horizontal(cudaSurfaceObject_t pingpong0, cudaSurfaceObject_t
pingpong1, float fft_dir) {
    int line = threadIdx.x;
    int column = blockIdx.x;
    int pingpong = 0;

    for(int stage = 0; stage < LOG_SIZE; ++stage) {
        int n = 1 << (LOG_SIZE - stage);
        int m = n >> 1;
        int s = 1 << stage;

        int p = line / s;
        int q = line % s;
        float2 wp = euler(fft_dir * 2 * (M_PI / n) * p);
        if(pingpong == 0) {
            float2 a = surf2Dread<float2>(pingpong0, (q + s*(p + 0))*sizeof(
float2), column);
            float2 b = surf2Dread<float2>(pingpong0, (q + s*(p + m))*sizeof(
float2), column);

            surf2Dwrite<float2>(complex_add(a, b), pingpong1, (q
+ s*(2*p + 0))*sizeof(float2), column);
            surf2Dwrite<float2>(complex_mult(wp,complex_sub(a, b)), pingpong1, (q
+ s*(2*p + 1))*sizeof(float2), column);
        }
        else {
            float2 a = surf2Dread<float2>(pingpong1, (q + s*(p + 0))*sizeof(
float2), column);
            float2 b = surf2Dread<float2>(pingpong1, (q + s*(p + m))*sizeof(
float2), column);

            surf2Dwrite<float2>(complex_add(a, b), pingpong0, (q
+ s*(2*p + 0))*sizeof(float2), column);
            surf2Dwrite<float2>(complex_mult(wp,complex_sub(a, b)), pingpong0, (q
+ s*(2*p + 1))*sizeof(float2), column);
        }

        pingpong = ((pingpong + 1) % 2);

        __syncthreads();
    }
}

__global__

```



```

void stockham_fft_vertical(cudaSurfaceObject_t pingpong0, cudaSurfaceObject_t
pingpong1, float fft_dir) {
    int line = blockIdx.x;
    int column = threadIdx.x;
    int pingpong = LOG_SIZE % 2;

    for(int stage = 0; stage < LOG_SIZE; ++stage) {
        int n = 1 << (LOG_SIZE - stage);
        int m = n >> 1;
        int s = 1 << stage;

        float mult_factor = 1.f;
        if ((stage == LOG_SIZE - 1) && fft_dir == 1) {
            mult_factor = 1.f / (FFT_SIZE*FFT_SIZE);
        }

        int p = column / s;
        int q = column % s;
        float2 wp = euler(fft_dir * 2 * (M_PI / n) * p);
        if(pingpong == 0) {
            float2 a = surf2Dread<float2>(pingpong0, line*sizeof(float2), (q + s
*(p + 0)));
            float2 b = surf2Dread<float2>(pingpong0, line*sizeof(float2), (q + s
*(p + m)));

            float2 ra = complex_add(a, b);
            float2 rb = complex_mult(wp, complex_sub(a, b));

            surf2Dwrite<float2>(ra * mult_factor, pingpong1, line*sizeof(float2),
(q + s*(2*p + 0)));
            surf2Dwrite<float2>(rb * mult_factor, pingpong1, line*sizeof(float2),
(q + s*(2*p + 1)));
        }
        else {
            float2 a = surf2Dread<float2>(pingpong0, line*sizeof(float2), (q + s
*(p + 0)));
            float2 b = surf2Dread<float2>(pingpong0, line*sizeof(float2), (q + s
*(p + m)));

            float2 ra = complex_add(a, b);
            float2 rb = complex_mult(wp, complex_sub(a, b));

            surf2Dwrite<float2>(ra * mult_factor, pingpong1, line*sizeof(float2),
(q + s*(2*p + 0)));
            surf2Dwrite<float2>(rb * mult_factor, pingpong1, line*sizeof(float2),
(q + s*(2*p + 1)));
        }
    }
}

```

```

        pingpong = (pingpong + 1) % 2;

        __syncthreads();
    }
}

```

**Listing C.2:** FFT Radix-2 Stockham, see [Section 6.3](#)

```

__device__ __forceinline__
float2 complex_mult(float2 v0, float2 v1) {
    return float2{
        v0.x * v1.x - v0.y * v1.y,
        v0.x * v1.y + v0.y * v1.x
    };
}

__device__ __forceinline__
float2 complex_add(float2 v0, float2 v1) {
    return float2{v0.x + v1.x, v0.y + v1.y};
}

__device__ __forceinline__
float2 complex_sub(float2 v0, float2 v1) {
    return float2{v0.x - v1.x, v0.y - v1.y};
}

__device__ __forceinline__
float2 euler(float angle) {
    return float2{cos(angle), sin(angle)};
}

__global__
void stockham_fft_horizontal(cudaSurfaceObject_t pingpong0, cudaSurfaceObject_t
    pingpong1, float fft_dir) {
    int line = threadIdx.x;
    int column = blockIdx.x;
    int pingpong = 0;

    for(int stage = 0; stage < HALF_LOG_SIZE; ++stage) {
        int n = 1 << ((HALF_LOG_SIZE - stage)*2);
        int s = 1 << (stage*2);

        int n0 = 0;
        int n1 = n/4;
        int n2 = n/2;
        int n3 = n1 + n2; // 3N/4
    }
}

```

```

    int p = line / s;
    int q = line % s;
    float2 w1p = euler(-2*(M_PI / n) * p);
    float2 w2p = complex_mult(w1p,w1p);
    float2 w3p = complex_mult(w1p,w2p);
    if(pingpong == 0) {
        float2 a = surf2Dread<float2>(pingpong0, (q + s*(p + n0))*sizeof(
float2), column);
        float2 b = surf2Dread<float2>(pingpong0, (q + s*(p + n1))*sizeof(
float2), column);
        float2 c = surf2Dread<float2>(pingpong0, (q + s*(p + n2))*sizeof(
float2), column);
        float2 d = surf2Dread<float2>(pingpong0, (q + s*(p + n3))*sizeof(
float2), column);

        float2 apc = complex_add(a,c);
        float2 amc = complex_sub(a,c);
        float2 bpd = complex_add(b,d);
        float2 jbmd = complex_mult(float2{0,1}, complex_sub(b,d));

        surf2Dwrite<float2>(complex_add(apc, bpd)
,
pingpong1, (q + s*(4*p + 0))*sizeof(float2), column);
        surf2Dwrite<float2>(complex_mult(w1p, complex_sub(amc, jbmd)),
pingpong1, (q + s*(4*p + 1))*sizeof(float2), column);
        surf2Dwrite<float2>(complex_mult(w2p, complex_sub(apc, bpd)) ,
pingpong1, (q + s*(4*p + 2))*sizeof(float2), column);
        surf2Dwrite<float2>(complex_mult(w3p, complex_add(amc, jbmd)),
pingpong1, (q + s*(4*p + 3))*sizeof(float2), column);
    }
    else {
        float2 a = surf2Dread<float2>(pingpong1, (q + s*(p + n0))*sizeof(
float2), column);
        float2 b = surf2Dread<float2>(pingpong1, (q + s*(p + n1))*sizeof(
float2), column);
        float2 c = surf2Dread<float2>(pingpong1, (q + s*(p + n2))*sizeof(
float2), column);
        float2 d = surf2Dread<float2>(pingpong1, (q + s*(p + n3))*sizeof(
float2), column);

        float2 apc = complex_add(a,c);
        float2 amc = complex_sub(a,c);
        float2 bpd = complex_add(b,d);
        float2 jbmd = complex_mult(float2{0,1}, complex_sub(b,d));

        surf2Dwrite<float2>(complex_add(apc, bpd)
,
pingpong0, (q + s*(4*p + 0))*sizeof(float2), column);

```

```

        surf2Dwrite<float2>(complex_mult(w1p, complex_sub(amc, jbmd)),
pingpong0, (q + s*(4*p + 1))*sizeof(float2), column);
        surf2Dwrite<float2>(complex_mult(w2p, complex_sub(apc, bpd)) ,
pingpong0, (q + s*(4*p + 2))*sizeof(float2), column);
        surf2Dwrite<float2>(complex_mult(w3p, complex_add(amc, jbmd)),
pingpong0, (q + s*(4*p + 3))*sizeof(float2), column);
    }

    pingpong = ((pingpong + 1) % 2);

    __syncthreads();
}
}

__global__
void stockham_fft_vertical(cudaSurfaceObject_t pingpong0, cudaSurfaceObject_t
pingpong1, float fft_dir) {
    int line = blockIdx.x;
    int column = threadIdx.x;
    int pingpong = HALF_LOG_SIZE % 2;

    for(int stage = 0; stage < HALF_LOG_SIZE; ++stage) {
        int n = 1 << ((HALF_LOG_SIZE - stage)*2);
        int s = 1 << (stage*2);

        int n0 = 0;
        int n1 = n/4;
        int n2 = n/2;
        int n3 = n1 + n2;

        int p = column / s;
        int q = column % s;
        float2 w1p = euler(-2*(M_PI / n) * p);
        float2 w2p = complex_mult(w1p,w1p);
        float2 w3p = complex_mult(w1p,w2p);
        if(pingpong == 0) {
            float2 a = surf2Dread<float2>(pingpong0, line*sizeof(float2), q + s*(
p + n0));
            float2 b = surf2Dread<float2>(pingpong0, line*sizeof(float2), q + s*(
p + n1));
            float2 c = surf2Dread<float2>(pingpong0, line*sizeof(float2), q + s*(
p + n2));
            float2 d = surf2Dread<float2>(pingpong0, line*sizeof(float2), q + s*(
p + n3));

            float2 apc = complex_add(a,c);
            float2 amc = complex_sub(a,c);

```

```

        float2 bpd = complex_add(b,d);
        float2 jbmd = complex_mult(float2{0,1}, complex_sub(b,d));

        surf2Dwrite(complex_add(apc, bpd) , pingpong1,
line*sizeof(float2), q + s*(4*p + 0));
        surf2Dwrite(complex_mult(w1p, complex_sub(amc, jbmd)), pingpong1,
line*sizeof(float2), q + s*(4*p + 1));
        surf2Dwrite(complex_mult(w2p, complex_sub(apc, bpd)) , pingpong1,
line*sizeof(float2), q + s*(4*p + 2));
        surf2Dwrite(complex_mult(w3p, complex_add(amc, jbmd)), pingpong1,
line*sizeof(float2), q + s*(4*p + 3));
    }
    else {
        float2 a = surf2Dread<float2>(pingpong1, line*sizeof(float2), q + s*(
p + n0));
        float2 b = surf2Dread<float2>(pingpong1, line*sizeof(float2), q + s*(
p + n1));
        float2 c = surf2Dread<float2>(pingpong1, line*sizeof(float2), q + s*(
p + n2));
        float2 d = surf2Dread<float2>(pingpong1, line*sizeof(float2), q + s*(
p + n3));

        float2 apc = complex_add(a,c);
        float2 amc = complex_sub(a,c);
        float2 bpd = complex_add(b,d);
        float2 jbmd = complex_mult(float2{0,1}, complex_sub(b,d));

        surf2Dwrite(complex_add(apc, bpd) , pingpong0,
line*sizeof(float2), q + s*(4*p + 0));
        surf2Dwrite(complex_mult(w1p, complex_sub(amc, jbmd)), pingpong0,
line*sizeof(float2), q + s*(4*p + 1));
        surf2Dwrite(complex_mult(w2p, complex_sub(apc, bpd)) , pingpong0,
line*sizeof(float2), q + s*(4*p + 2));
        surf2Dwrite(complex_mult(w3p, complex_add(amc, jbmd)), pingpong0,
line*sizeof(float2), q + s*(4*p + 3));
    }

    pingpong = (pingpong + 1) % 2;

    __syncthreads();
}
}

```

**Listing C.3:** FFT Radix-4 Stockham, see [Section 6.3](#)