

UNIVERSITY OF CAPE TOWN



EEE4120F

HIGH PERFORMANCE EMBEDDED SYSTEMS

---

# Practicals and Projects 2023

---

21 February 2023

## Introduction

Welcome to the practicals for EEE4120F. Please take note of these instructions that are applicable to all practicals in this course.

Practicals, with tutors to assist, generally take place on a Thursday, from 14h00 to 16h00; but please check the schedule and announcements on Vula as the schedule may change. You are responsible for completing these practicals in order to pass this course. Some of these practicals may require access to specialist hardware that will only be available in labs during the practical session - however, for 2022 we have tried to plan all pracs needed for this course to require nothing other than a decent computer, ideally a computer with a nVidia GPU which has CUDA and OpenCL support (most recent nVidia GPUs provide such support). The FPGA pracs have been reworked so that they can be done using simulation, thus you do not need to make use of labs for even those practicals which should enable you to pass the course without needing to use labs on campus if you would prefer not to.

It is assumed you have knowledge regarding tools such as git, L<sup>A</sup>T<sub>E</sub>X and running programs from the Linux command line. If you are wanting to do the pracs on your own computer, and don't have Linux installed, you are suggested to install Ubuntu Linux either as a dual boot option or to run it using virtualization, e.g. [VirtualBox](#).

It is suggested you have a dual boot system. The reason for this is that development tools are generally better on a Linux based system, whereas proprietary tools are often better supported on Windows. We recommend Ubuntu 18.04 LTS or later, though swapping between Ubuntu (or any Linux distro) and Windows is an effective strategy, particularly when using proprietary software.

The source code for all pracs is available on the EE-OCW GitHub Project Page: <https://github.com/UCT-EE-OCW/EEE4120F-Pracs>. You can download the source code for pracs from there.

Please note that all practical's have to be completed in pair's with the exception of practical 1, where you can submit either as a pair or individually.

## Prac Overview

Please note that these dates are subject to change. For due dates, refer to Vula.

Table I: The Prac Overview for EEE4120F, 2022

Prac	Submission	Dates	Tools	Objective
Prac 1	Pairs or Individual	See Vula	JULIA/OCTAVE	Performance measurement
Prac 2.1	Individual	See Vula	OpenCL	Introduction
Prac 2.2	Pairs	See Vula	OpenCL	Matrix Multiplication
Prac 3.1	Individual	See Vula	Vivado	FPGA simulation intro
Prac 3.2	Pairs	See Vula	Vivado	Simulation-based FPGA
Prac 4		See Vula		<i>Choose either Prac4a or Prac4b below (archived*)</i>
Prac 4a		See Vula	Vivado	Prac 4a Module Design and Simulation (archived*)
Prac 4b		See Vula	Vivado	Prac 4b FPGA intro using the Nexys platform (archived*)
Prac 4c		Optional	Vivado	Prac 4c Making a signal generator (archived*)
Project	See Vula		YODA design, implementation & demos	

\* the pracs indicated as 'archived' have been left for potential useful practice or training that may be of use to you in the course project.

# 1 Practical 1 - Julia

## 1.1 Introduction

The focus of this task is on using Julia or OCTAVE (the free sort-of MATLAB program) to do some statistical operations. In later assignments you will make further use of these statistical functions, and perhaps reuse this code, in comparing and discussing results obtained in other practicals and projects (for example, correlation can be used in analysing gold standard results to higher-speed approximation result).

For this practical you will have the option of either submitting the practical in Julia or OCTAVE. Furthermore, this practical can be completed either in pairs or individually. However, pairs are recommended as all subsequent pracs have to be completed in pairs.

For installation and tips and tricks for this course, visit [the EE Wiki](#) (currently only accessible on the UCT Network - though you can use a VPN). Throughout this course you will be utilizing a verity of packages and modules and it is recommended to use a linux distribution or raspbian (if programming on a raspberry pi), due to the setup process often being easier. However, this is not a requirement and you may use any OS.

### 1.1.1 Julia installation tips

If you do not have Julia installed on your device please follow the steps below before attempting the practical. Furthermore, this installation guide will help you setup Julia to run in terminal but you can use an IDE if you would prefer but the tutors are not required to debug IDE issues.

1. Download Julia [here](#)
2. Extract/unzip the downloaded file in the desired directory.
3. Open Julia in terminal to test that the download was successful by running the julia file in the bin directory. In order to do this traverse your devices directories in terminal until you are in the following directory: './julia-1.7.1/bin'. Now run the command 'julia' in terminal to open the julia.
4. Create a symbolic link to the location you have installed julia. This symbolic link will allow you to easily run julia without having to go to the location that julia is saved each time. In order to create the symbolic link run the following cammand:  
'sudo ln -s /julasDirectoryOnYourComputer/bin/julia /usr/local/bin/julia'  
or  
'sudo ln -s ~/julasDirectoryOnYourComputer/bin/julia /usr/local/bin/julia'.
5. Provided you have setup the symbolic link correctly you can run julia script in any directory with the following command: 'julia scriptFileName.jl'. Note julia script has the '.jl' extension.

6. Have fun programming julia scripts in your favourite text editor!

### 1.1.2 Octave installation tips

To install OCTAVE please follow the steps [on the EE wiki](#).

### 1.1.3 Correlation

Correlation is a useful statistical function for comparing two datasets to judge how similar or different they are. The correlation function returns a correlation coefficient,  $r$ , between -1 and 1. A value of 1 for  $r$  implies perfect positive correlation, i.e. the two datasets are the same. Correlation of 0 implies there is no correlation (the two datasets behave totally differently). A correlation of -1 indicates a total opposite - for example if you compare vectors  $x$  to  $-x$  you get a correlation of -1. Generally if  $|r| \geq 0.8$  there is strong correlation, between 0.5 and 0.8 moderate weak, less than 0.5 is weak (towards no) correlation.

Pearson's correlation (which you can read more about on [Wikipedia](#)) is implemented as follows:

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \quad (1)$$

### 1.1.4 Speed Up

Speedup used in this practical and subsequent practicals is defined as follows:

$$Speedup = \frac{T_{p1}}{T_{p2}} \quad (2)$$

Where:

$T_{p1}$  = Run-time of original / non-optimal program

$T_{p2}$  = Run-time of optimised program For obtaining a repeatable timing value, run each version (i.e. initial version and optimised version) of your programs more than once and discard the first measured time. You can, if you want to be complete, indicate what the initial speed up was and then the average speed up.

### 1.1.5 Critical Section

When measuring performance, it's important to ensure that you are only measuring the section of your algorithm that you need to measure. For example, if you are measuring the execution time of an audio filter, you should not be recording the time taken to open and close the files you are applying this filter to.

## 1.2 Julia Requirements

You are required to run the following experiments:

### 1.2.1 Measuring Execution Time of `rand()`

White noise is often generated with Julia's random number generator `rand()`, which generates uniformly distributed random values in the interval  $[0,1)$ . To create a sound wave .wav file of the white noise, the `wavwrite()` function in the WAV package is used. The `wavwrite()` function expects values in  $[-1.0, 1.0)$ , so the `rand()` output must be multiplied by 2 and shifted down by 1. To generate 10 seconds white noise sampled at 48 kHz, the following instruction are called:

```
whiteNoise = (rand(48000)*2).-1
```

And to generate a wave file for this noise, we use the `wavwrite()` function as follows:

```
WAV.wavwrite(whiteNoise, "whiteNoise.wav", Fs=4800) #sample freq is 4800Hz
```

### 1.2.2 White Noise Generator Script

Next, write a function called `createnwhite.m` that implements a function with a **for loop** that generates a white noise signal, one sample at a time, comprising N duration in seconds. You can assume that N will always be positive and a multiple of 10. The white noise must be sampled at either 48 kHz. Name your function `createnwhite(...)`. You need to use the `rand()` function without arguments so that it will generate a single random value, and the main task is figuring out how to scale so that you create a suitable input to `wavwrite(...)` as explained above. The function should include a print statement outputting the number of samples in the white noise before returning the white noise as an array. The returned white noise array should then be saved as an '.wav' file.

```
whiten = createnwhite(1000); #this will create 1000 seconds of white noise  
size(whiten);
```

Outputted print statement:

```
Number of samples: (48000,)
```

should return:

```
an array with the white noise
```

After converting the white noise array into a '.Wav' file, check that the resulting wave/sound file gives the same sound as the white noise sound1.wav generated above by generating the new sound file named white noise sound2.wav and playing it back.

```
WAV.wavwrite(returnedNoise, "white_noise_sound2.wav", Fs=4800)
```

### 1.2.3 Visual Confirmation of Uniform Distribution

Confirm that you've created the sample correctly. Do this using the `Plots.histogram()` function in the `Plots` package.

```
h = Plots.histogram(whiteNoiseArray)
Plots.display(h)
readline() #this will stop the program at this point till you press enter
```

This should give image shown in Figure 1:

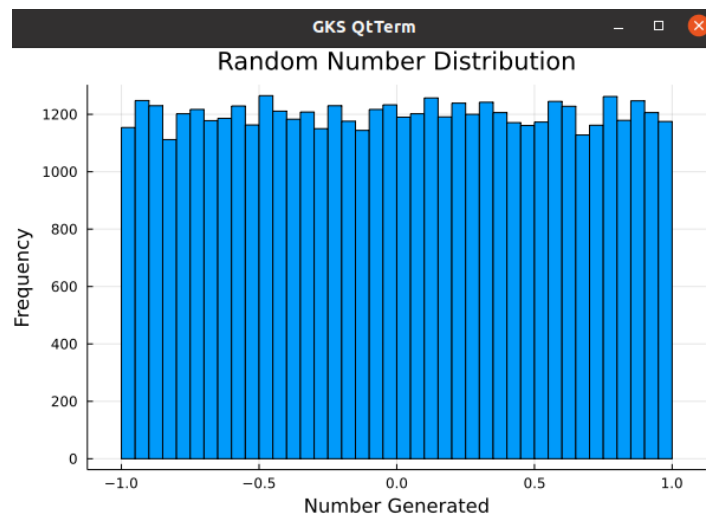


Figure 1: Output of `Plots.histogram(whiteNoiseArray)`

### 1.2.4 Timing Execution

After confirming that your function works correctly use the `TickCount` package to compare the performance of the two approaches to generate white noise. In order to use the `TickCount` package run the command `tick()` to start the timer and the command `tock()` to stop the timer and output the time that has elapsed. Repeat the white noise generation process for varying sample sizes and consider add this to your report.

```
tick()
commands you wish to time
tock()
```

### 1.2.5 Implementing Pearson's Correlation

Implement the Pearson's correlation formula a new function call it `corr()`.

### 1.2.6 Comparing Your Correlation Function to the Statistics Package's Correlation Function

The `cor` function in Statistics package performs a correlation operation. Compare your `corr()` results to the `Statistics.cor` function. Run the following test:

1. Using your output from the `createWhiten(numberOfSeconds)` function compare its correlation against itself using your correlation function and the `Statistics.cor` function.
2. Compare the correlation between the two different approaches to generate white noise, using your correlation function and the `Statistics.cor` function.

Do a table listing sample sizes vs. `corr` speed vs. `Statistics.cor` speed. Indicate the average speed-up, etc. in your report.

### 1.2.7 Correlation of Shifted Signals

For the last experiment, we want to compare signals shifted in time. Generate sin curves of varying frequency and sampling sizes (again sample sizes 100, 1000 and 10000 samples). Compare samples of the same sizes that are shifted in time, consider using the `circshift(sineWave, 10)`, which will cause a shift of 10 samples.

For this step only use `Statistics.cor` to save time. In your report, discuss what you expect the correlation of the identical but shifted signals would be. Run tests to confirm / verify your hypothesis. Provide some screen shots plots of some of the signals you compared. It might be useful to use the `Plots.scatter()` function to generate these plots.

## 1.3 OCTAVE Requirements

You are required to run the following experiments:

### 1.3.1 Measuring Execution Time of `rand()`

White noise is often generated with GNU Octave's random number generator `rand()`, which generates uniformly distributed random values in the interval  $[0,1)$ . To create a sound wave of the white noise, the `wavwrite()` function in Octave is used. The `wavwrite()` function expects values in  $[-1.0, 1.0)$ , so the `rand()` output must be multiplied by 2 and shifted down



by 1. To generate 10 seconds white noise sampled at 48 kHz, the following instruction are called:

```
white = rand(48000*10,1)*2-1;
```

And to generate a wave file for this noise, we use the `wavwrite()` function as follows:

```
wavwrite(white, 48000, 16, 'white_noise_sound.wav');
```

#### NOTE:

`wavwrite` and `wavread` were deprecated in Octave 5.1. If you're using Octave 5.1 or greater, you need to use `audiowrite` and `audioread`.

```
audiowrite('w.wav', white, 48000, 'BitsPerSample', 16);  
[y, fs] = audioread('w.wav', 16);
```

### 1.3.2 White Noise Generator Script

Next, write a function in a new file called `createwhiten.m` that implements a function with a **for loop** that generates a white noise signal, one sample at a time, comprising N duration in seconds. You can assume that N will always be positive and a multiple of 10. The white noise must be sampled at 48 kHz. Name your function `createwhiten(...)`. You need to use the `rand()` function without arguments so that it will generate a single random value, and the main task is figuring out how to scale so that you create a suitable input to `wavwrite(...)` as explained above. Call the function and check output size as follows:

```
whiten = createwhiten(1000);  
size(whiten);
```

should return:

```
ans = 48000000 1
```

Check that the resulting wave/sound file gives the same sound as the white noise sound.wav generated above by generating the new sound file named white noise sound2.wav and playing it back.

```
wavwrite(whiten, 48000, 16, 'white_noise_sound2.wav');
```

### 1.3.3 Visual Confirmation of Uniform Distribution

Confirm that you've created the sample correctly. Since it's a big signal, let's just look at the first 100 samples by plotting using a histogram function as follows:

```
hist(whiten, 100, 1);
```

This should give image shown in Figure 2:

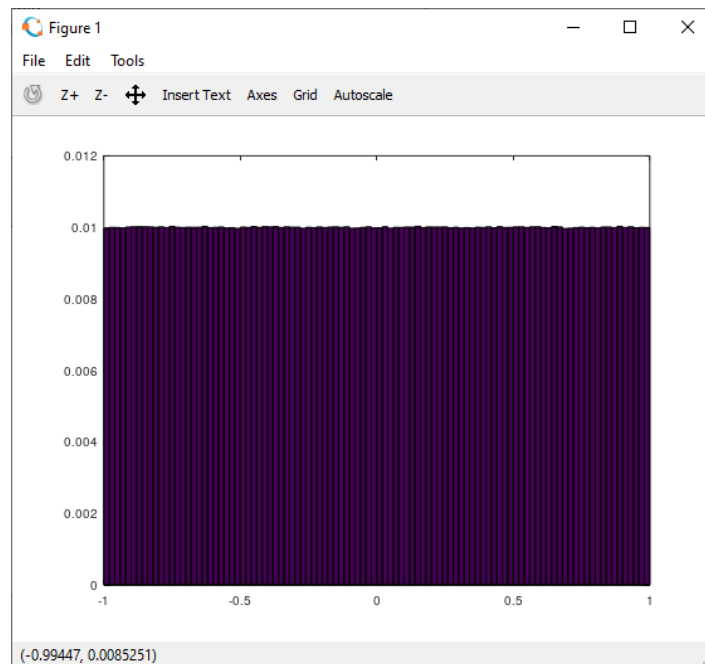


Figure 2: Output of `hist(whiten, 100, 1);`

### 1.3.4 Timing Execution

Now time how long it took to execute the script. Use the `tic()` and `toc()` functions as follows. Note that the call to `tic` is on the same line just before the function - this tends to have less delay between the start of the timer and starting the function. Of course, it is good practice to put it all in a script file.

```
tic; white = rand(48000*1000, 1)*2 - 1; runtime = toc();  
disp(strcat("It took: ", num2str(runtime*1000), " ms to run"));
```

Call the `createwhiten(...)` function you created that does the same thing as the `white = rand(48000*1000,1)*2 - 1;` statement that generate the white noise. Measure the time it took for your `createwhiten(...)` function to run and show the timing difference (in milliseconds) and discuss the speed-up you have achieved (if any).

### 1.3.5 Implementing Pearson's Correlation

Implement the Pearson's correlation formula a new m file called it `mycorr.m`. Provide your code in your report. Marks will be awarded on elegance of your code.

### 1.3.6 Comparing Your Correlation Function to the Built-in Correlation Function

The `corr` function built into OCTAVE also performs a correlation operation. Compare your `mycorr.m` results to the built in correlation function. Run the following test:

1. Using your output from the `createwhiten(numberOfSeconds)` function compare its correlation against itself using your correlation function and the built-in OCTAVE correlation function.
2. Compare the correlation between the two different approaches to generate white noise, using your correlation function and the built-in correlation function.

Do a table listing sample sizes vs. corr speed vs. Built-In speed. Indicate the average speed-up, etc. in your report.

### 1.3.7 Correlation of Shifted Signals

For the last experiment, we want to compare signals shifted in time. Generate sin curves of varying frequency and sampling sizes (again sample sizes 100, 1000 and 10000 samples). Compare samples of the same sizes that are shifted in time.

For this step only use the built-in correlation function to save time. In your report, discuss what you expect the correlation of the identical but shifted signals would be. Run tests to confirm / verify your hypothesis. Provide some screen shots plots of some of the signals you compared. It might be useful to provide a scatter plot of your results.

## 1.4 Submission

Hand in a report about 2-3 pages in length briefly describing your solutions for the tasks above. The format required is the IEEE conference format. Format your report as if it is an article (i.e. don't follow the same chronology as the prac-sheet – format it as “Introduction – Method – Results and Discussion – Conclusion”). In the “Method” section, theorise about what you expect and how you plan on testing said theory. In the “Results” section, confirm that you obtained what you expected (or explain why you obtained something unexpected). Hint: You are trying to answer two questions: 1) “Is my `whitenoise()` function better suited to generated white noise, in comparison to the `rand()` equivalent?” 2) “Is my `corr()` function better suited to run correlation tests, in comparison to the `Statistics.corr` equivalent?” and 2) “Are my theories relating to the correlation of time-shifted sinusoids correct?” Provide key sections of code in your report. Marks will be awarded on elegance of your code and if your results have been portrayed in a logical way (i.e. tables and graphs).

## 1.5 Marking

Table II: Prac 1 Marking Guide

Aspect	Description	Mark Allocation
Report	Intro	2
	Latex/Format	1.5
	Headings etc	1.5
	Captions	1
	Discussion	2
createwhiten	Code	4
	Histogram	2
	Results	3
	Conclusion and Explanation	3
Correlation	Code	4
	Results	2
	Conclusion and Explanation	2
Shifted signals	Code	4
	Plots	3
	Results	2
	Conclusion and Explanation	3
<b>TOTAL</b>		40

## 2 Prac 2.1 - OpenCL

### 2.1 Introduction

This practical focuses on developing an OpenCL kernel, using C++, that you can load, activate and send data back and forth to, from a C / C++ host application. The following quote, from Kronos [1], summarises what OpenCL is:

OpenCL™ (Open Computing Language) is the open, royalty-free standard for cross-platform, parallel programming of diverse processors found in personal computers, servers, mobile devices and embedded platforms. OpenCL greatly improves the speed and responsiveness of a wide spectrum of applications in numerous market categories including gaming and entertainment titles, scientific and medical software, professional creative tools, vision processing, and neural network training and inferencing. [1]

The main purpose of practical 2.1 is to show you how a typical OpenCL program would run by taking a step-by-step approach. Furthermore, this practical has to be completed **individually**.

### 2.2 Requirements

The Blue Lab PCs have all the required packages installed, but if you wish to run this practical on your own machine, you will need to install OpenCL as per the instructions on the Wiki - <http://wiki.ee.uct.ac.za/OpenCL>. Furthermore, you will have to install the appropriate SDK (software development kits) for the device you are planning to run the program on. The following links will direct you to the appropriate websites explaining how to install the SDK for their respective devices.

- Nvidia: <https://developer.nvidia.com/cuda-toolkit>
- Intel: <https://www.intel.com/content/www/us/en/developer/tools/opencl-sdk/choose-download.html>
- AMD: [https://developer.amd.com/wordpress/media/2012/10/AMD\\_APP\\_SDK\\_InstallationNotes.pdf](https://developer.amd.com/wordpress/media/2012/10/AMD_APP_SDK_InstallationNotes.pdf)

Lastly, all code required to start this practical is on the courses github repository <https://github.com/UCT-EE-OCW/EEE4120F>.

### 2.3 The Programming Model

OpenCL uses a programming and memory model similar to OpenGL. The CPU must copy data to the GPU and then tell a kernel (OpenCL worker) to process the data. When the

kernel is finished, the CPU can read back the result. This is an overly simplistic summary, the steps below will explain in more detail how this process takes place.

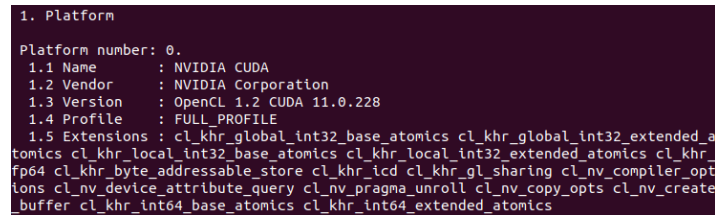
## 2.4 Step-by-Step Explanation

This step-by-step explanation works very closely to the code provided in the github repository, so please clone the repository if you have not already.

### 2.4.1 Preparation

Before, setting up the the main task for this practical you have to check what platforms (SDKs) you have installed on your host computer and determine what device you would like to run the program on. In order to do this the program `platform.cpp` has been given to you. This c++ program will be used to test if you have setup OpenCL correctly and then output the index's for the platforms installed on your host device.

After compiling and running the `platform.cpp` program you should get an output similar to the image below.



```
1. Platform
Platform number: 0.
1.1 Name      : NVIDIA CUDA
1.2 Vendor    : NVIDIA Corporation
1.3 Version    : OpenCL 1.2 CUDA 11.0.228
1.4 Profile    : FULL_PROFILE
1.5 Extensions : cl_khr_global_int32_base_atomics cl_khr_global_int32_extended_a
tomics cl_khr_local_int32_base_atomics cl_khr_local_int32_extended_atomics cl_khr_
fp64 cl_khr_byte_addressable_store cl_khr_icd cl_khr_gl_sharing cl_nv_compiler_opt
ions cl_nv_device_attribute_query cl_nv_pragma_unroll cl_nv_copy_opts cl_nv_create
buffer cl_khr_int64_base_atomics cl_khr_int64_extended_atomics
```

Figure 3: Example Output of `Platform.cpp`

Take note of the the 'platform number' because you will need this later in the practical.

The following steps will require you to open the `main.cpp` file and make the appropriate changes. The `main.cpp` file has sections of code missing which you will have to fill in, these sections are labeled 'TODO code' with the code numbers. Whilst, the manual below will be labeled similarly.

### 2.4.2 Step 1

OpenCL is a heterogeneous programming tool, which means you can program on different devices and differing platforms on the host computer. This in turn requires you to select the correct platform for the device you would like to run your program on. In order to select the platform you would like to use the following command is run.

```
cl_int clGetPlatformIDs(cl_uint num_entries,
                        cl_platform_id *platforms,
```

```
cl_uint *num_platforms)
```

This function saves all available platforms in a list of type `cl_platform_id` with its pointer being set the argument `*platforms`. We can now select which device to use by indexing this list, with the platform number obtained when running the `platform.cpp` program.

```
//TODO: code 1
//Change index depending on the platform you want to use
cl_platform_id platform = platforms[0];
```

### 2.4.3 Step 2

After selecting your platform you need to choose the device that this platform will be utilizing. The following command allows you to do that.

```
cl_int clGetDeviceIDs(cl_platform_id platform,
                     cl_device_type device_type,
                     cl_uint num_entries,
                     cl_device_id *devices,
                     cl_uint *num_devices)
```

Similarly, to the `clGetPlatformIDs()` function, `clGetDeviceIDs()` saves the device id's to the `cl_device_id` variable that you enter as an argument as shown below. Also, note that you will have to indicate the type of device you are looking for and due to most computer only having one CPU and possible one GPU we set the `num_entries` to 1. Resulting in your 'devices' variable containing only one device id, which is your GPU or CPU.

```
\\TODO code 2
err = clGetDeviceIDs(platform, CL_DEVICE_TYPE_GPU, 1, &device, NULL);

OR

err = clGetDeviceIDs(platform, CL_DEVICE_TYPE_GPU, 1, &device, NULL);
```

### 2.4.4 Step 3

Often OpenCL programs have multiple different devices that provide processing for various aspects of the program, possibly in parts of the program there may be different devices running in parallel (for example you might have the main CPU showing displays, a GPU doing serious computations and possible another device as well, such as a configurable sampling card receiving raw data to process). In this practical we are only dealing with a GPU or CPU. However, due to the heterogeneous nature of OpenCL, a context has to be created to house all available devices that the the program is planning on using.

The following command creates a context with all the devices you are planning on using in the program, which you obtained in step 2.

```
cl_context clCreateContext(cl_context_properties *properties,
                          cl_uint num_devices,
                          const cl_device_id *devices,
                          void *pfn_notify(const char *errinfo,
                          const void *private_info, size_t cb,
                          void *user_data),
                          void *user_data, cl_int *errcode_ret)
```

For this practical we are only using the one device obtained in step two and the following command has to be added to main.cpp.

```
\\TODO code 3
cl_context context;
    context = clCreateContext(NULL, 1, &device, NULL, NULL, NULL);
```

#### 2.4.5 Step 4

Now that the platforms and devices have been added to the context, the program we want to run on these devices need to be loaded into a buffer, as a string. So that we can isolate what we want to run on the various devices. The file used in this practical is, 'OpenCL/Kernel.cl'.

```
//TODO code 4
program_handle = fopen("OpenCL/Kernel.cl", "r");
```

#### 2.4.6 Step 5

A program variable has to then be created from the source code obtained in step 4, essentially adding the program to the context. Note, you could have multiple programs each with their own program ID.

```
cl_program clCreateProgramWithSource(cl_context context,
                                    cl_uint count,
                                    const char **strings,
                                    const size_t *lengths,
                                    cl_int *errcode_ret)
```

```
//TODO code 5
cl_program program = clCreateProgramWithSource(context, 1,
                                              (const char**)&program_buffer, &program_size, NULL);
```



### 2.4.7 Step 6

In order for our program to run on the device we want, the program has to be compiled for that specific device. The function `clBuildProgram()` compiles your selected program for a specific device, using the device IDs obtained in step 2.

```
cl_int clBuildProgram(  
    cl_program program,  
    cl_uint num_devices,  
    const cl_device_id* device_list,  
    const char* options,  
    void (CL_CALLBACK* pfn_notify)(cl_program program, void* user_data),  
    void* user_data);
```

```
\\TODO code 6  
cl_int err3= clBuildProgram(program, 0, NULL, NULL, NULL, NULL);
```

### 2.4.8 Step 7

Now that the program has been compiled for the device you want to run it on. You need to select what section of this program that you have compiled, you want to run on the selected device. These sections of code are referred to as kernels, and are represented as functions in the 'cl\_program' setup in step 6. If you look in the file `OpenCL/Kernel.cl` you will see a function called `HelloWorld`, this is the kernel for this practical.

```
cl_kernel clCreateKernel(cl_program program,  
                        const char* kernel_name,  
                        cl_int* errcode_ret);
```

```
\\TODO code 7  
cl_kernel kernel = clCreateKernel(program, "HelloWorld", &err);
```

### 2.4.9 Step 8

Currently, we now have all the platforms, devices and kernels setup for the practical. However, the kernels needs a way to get from the host to your selected device (ie GPU or CPU). A command queue is required to facilitate this process.

```
cl_command_queue clCreateCommandQueueWithProperties(cl_context context,  
                                                  cl_device_id device,  
                                                  const cl_queue_properties* properties,  
                                                  cl_int* errcode_ret);
```

For this practical the queue that you have to setup is between the host and the device you have selected. This queue is created with the following command.

```
\\TODO code 8
cl_command_queue queue = clCreateCommandQueueWithProperties(context, device,
                                                         0, NULL);
```

#### 2.4.10 Step 9

Note the queue only allows for kernels to be sent, which is contains only the program that must be run for each work item. So we need to setup data buffers to allow for communication between devices. Which is our case is communication between the host and target device.

This can be done with the `clCreateBuffer()` function which sets up shared memory between devices that both the host and target device can access. This has to be done as your GPU by default uses its VRAM, while the host will use the computers RAM, and these memory units are completely separate, so the GPU can't access the RAM and the host CPU can't access the VRAM. Therefore, the memory blocks we are setting up now, are block's of memory where both devices can access them.

```
cl_mem clCreateBuffer(cl_context context,
                     cl_mem_flags flags,
                     size_t size,
                     void* host_ptr,
                     cl_int* errcode_ret);
```

For this practical you need to setup 3 buffers, the first 2 buffers are only one integer in size and will allow for the target device to access these arguments. Note in the code below that your target device may read these buffers but may not alter them.

The last buffer will be your output buffer used to store the outputs from your kernel, so that the host can retrieve what the kernels calculated.

```
\\TODO code 9.2
argument1_buffer = clCreateBuffer(context,
                                  CL_MEM_READ_ONLY | CL_MEM_COPY_HOST_PTR,
                                  sizeof(int), &argument1, &err);

argument2_buffer = clCreateBuffer(context,
                                  CL_MEM_READ_ONLY | CL_MEM_COPY_HOST_PTR,
                                  sizeof(int), &argument2, &err);

output_buffer = clCreateBuffer(context,
                               CL_MEM_READ_WRITE | CL_MEM_COPY_HOST_PTR,
                               global_size*local_size*sizeof(int), output, &err);
```

Furthermore, if you look at the output buffer, we set the size to `global_size*sizeof(int)`, this is due to the fact that each work item needs to have its own output of one integer in size.

Therefore, before setting the memory buffers we need to determine how many work items and work groups we want for the program. So we set the following variables that contain the required information, which we will use at a later stage.

```
\\TODO step 9.1
size_t global_size = 16; //total number of work items
size_t local_size = 4; //Size of each work group
cl_int num_groups = global_size/local_size; //number of work groups needed
```

#### 2.4.11 Step 10

In order for our target device to utilize the memory blocks created in step 9, we have to add arguments to the kernel, giving the kernel the pointers to the blocks of memory. The `clSetKernelArg()` function links these memory blocks to the kernel. Note, when looking at the `OpenCL/Kernel.cl` file you will see the same arguments needed to run the kernel `HelloWorld`.

```
cl_int clSetKernelArg (cl_kernel kernel,
                      cl_uint arg_index,
                      size_t arg_size,
                      const void *arg_value)
```

In this practical we have setup 3 memory blocks, hence we need 3 kernel arguments as shown below.

```
\\TODO step 10
clSetKernelArg(kernel, 0, sizeof(cl_mem), &argument1_buffer);
clSetKernelArg(kernel, 1, sizeof(cl_mem), &argument2_buffer);
clSetKernelArg(kernel, 2, sizeof(cl_mem), &output_buffer);
```

#### 2.4.12 Step 11

Finally, we can now deploy the kernels to the target device using the `clEnqueueNDRangeKernel()` function and set the number of work items per work groups. The work items each run the `kernel.cl` file and can only access memory from other work items in the same group.

```
cl_int clEnqueueNDRangeKernel (cl_command_queue command_queue,
                               cl_kernel kernel,
                               cl_uint work_dim,
                               const size_t *global_work_offset,
                               const size_t *global_work_size,
                               const size_t *local_work_size,
```

```

        cl_uint num_events_in_wait_list,
        const cl_event *event_wait_list,
        cl_event *event)

```

For this practical we have already setup the number of work items (`global_size`) and work groups, which we set in 9.1, so we just deploy the kernels as follows.

```

//TODO code 11
cl_int err4 = clEnqueueNDRangeKernel(queue, kernel, 1, NULL,
                                     &global_size, &local_size, 0, NULL, NULL);

```

### 2.4.13 Step 12

After the kernels have run, we now want the host to read from the output memory buffer, using the following command.

```

//read from a buffer object to host memory
cl_int clEnqueueReadBuffer(
    cl_command_queue command_queue,
    cl_mem buffer,
    cl_bool blocking_read,
    size_t offset,
    size_t size,
    void* ptr,
    cl_uint num_events_in_wait_list,
    const cl_event* event_wait_list,
    cl_event* event);

//write to a buffer object from host memory
cl_int clEnqueueWriteBuffer(
    cl_command_queue command_queue,
    cl_mem buffer,
    cl_bool blocking_write,
    size_t offset,
    size_t size,
    const void* ptr,
    cl_uint num_events_in_wait_list,
    const cl_event* event_wait_list,
    cl_event* event);

```

In this practical we only need the host to read from the output buffer object, as shown below.

```

\\TODO Step 12
err = clEnqueueReadBuffer(queue, output_buffer, CL_TRUE, 0,
                          sizeof(output), output, 0, NULL, NULL);

```

#### 2.4.14 Step 13

The last step is not actually part of the OpenCL process but rather just checking that the output received by the host is the same as the outputs calculated by the kernels.

```
printf("\nOutput in the output_buffer \n");
for(int j=0; j<global_size; j++) {
    printf("element number:%i \t Output:%i \n",j ,output[j]);
}
```

### 2.5 Kernel Creation

Provided that you completed the Step by Step Explanation correctly when running the program you should get the following output, where each work item prints 'Hello World'. Notice that the output value from step 13, are random numbers because the kernel has not updated the output buffer in any way. Therefore, the host's output will be whatever values were in those memory addresses before the program ran.

```
Name       : NVIDIA CUDA
Vendor     : NVIDIA Corporation
Version    : OpenCL 1.2 CUDA 11.0.228
Profile    : FULL_PROFILE
Device ID  = 0
program ID = 0
Hello World
Hello World
Hello World
Hello World
Hello World
Hello World
Hello World
Hello World
Hello World
Hello World
Hello World
Hello World
Hello World
Hello World
Hello World
Hello World
Hello World
Hello World
Hello World
Hello World
Kernel check: 0

Output in the output_buffer
element number:0      Output:-1131117520
element number:1      Output:21940
element number:2      Output:-1156875328
element number:3      Output:21940
element number:4      Output:-1208024240
element number:5      Output:32767
element number:6      Output:-1208024480
element number:7      Output:32767
element number:8      Output:-1208024000
element number:9      Output:32767
element number:10     Output:0
element number:11     Output:0
element number:12     Output:0
element number:13     Output:0
element number:14     Output:-1156876259
element number:15     Output:21940
```

Figure 4: Hello World Example Output

The final task for this practical is to edit the OpenCL/Kernel.cl file to complete the following tasks.

### 2.5.1 Task 1

Perform a basic numeric calculation that utilizes all the arguments for the kernel. The required calculation is given below.

$$output[pos] = workItemNumber \times argument1 + argument2$$

### 2.5.2 Task 2

Print the work item, work group, arguments and output values, as follows.

```
Hi from work item: 12    work group:3    Arg1: 10    Arg2: 20    Output: 140
Hi from work item: 13    work group:3    Arg1: 10    Arg2: 20    Output: 150
Hi from work item: 14    work group:3    Arg1: 10    Arg2: 20    Output: 160
Hi from work item: 15    work group:3    Arg1: 10    Arg2: 20    Output: 170
Hi from work item: 4     work group:1    Arg1: 10    Arg2: 20    Output: 60
Hi from work item: 5     work group:1    Arg1: 10    Arg2: 20    Output: 70
Hi from work item: 6     work group:1    Arg1: 10    Arg2: 20    Output: 80
Hi from work item: 7     work group:1    Arg1: 10    Arg2: 20    Output: 90
Hi from work item: 0     work group:0    Arg1: 10    Arg2: 20    Output: 20
Hi from work item: 1     work group:0    Arg1: 10    Arg2: 20    Output: 30
Hi from work item: 2     work group:0    Arg1: 10    Arg2: 20    Output: 40
Hi from work item: 3     work group:0    Arg1: 10    Arg2: 20    Output: 50
Hi from work item: 8     work group:2    Arg1: 10    Arg2: 20    Output: 100
Hi from work item: 9     work group:2    Arg1: 10    Arg2: 20    Output: 110
Hi from work item: 10    work group:2    Arg1: 10    Arg2: 20    Output: 120
Hi from work item: 11    work group:2    Arg1: 10    Arg2: 20    Output: 130
```

Figure 5: Expected Output

### 2.5.3 Task 3

For the final task you have to add the output values of all the work items in each work group.

In order to solve this issue consider isolating one work item, with an if statement, from each work group. This work item can then iterate through all the work items in the work group and add up each of their output values. Furthermore, make use of the barrier function that stops work items from continuing until all the work item in the group reach this barrier point.

```
groupValue: 300          Work item:4      Work group: 1
groupValue: 620          Work item:12     Work group: 3
groupValue: 140          Work item:0      Work group: 0
groupValue: 460          Work item:8      Work group: 2
```

Figure 6: Expected Output

## 2.6 Submission

For the final submission, submit your code in a zipped file labeled as follows ‘Prac2.1\_STUDENTNUMBER’. The zipped file should contain the following files:

1. main.cpp
2. platforms.cpp
3. OpenCL/Kernel.cl
4. Any other compiled files

Table III: Prac 2.1 Marking Guide

Aspect	Description	Mark Allocation
Main.cpp		
	Steps	13
Kernel.cl		
	Initial Calculation	1
	Print Output	2
	Work Group Total	4
TOTAL		20

## 3 Prac 2.2 - OpenCL Matrix Multiplication

### 3.1 Introduction

This practical focuses on developing an OpenCL kernel, using C++, that you can will load, activate and send data back and forth to, from a C / C++ host application. The focus of this practical is investigating the performance metrics of an OpenCL matrix multiplication implementation and how depending on the type of hardware and algorithms being running performances can vary. Please note that this practical has to be completed in groups of two.

### 3.2 Requirements

Practical 2.2 is the extension of practical 2.1. In the previous practical we covered all the steps involved in creating your own OpenCL program. In practical 2.2 you will create your own kernels and run the appropriate tests.

The Blue Lab PCs have all the required packages installed, but if you wish to run this practical on your own machine, you will need to install OpenCL as per the instructions on the Wiki - <http://wiki.ee.uct.ac.za/OpenCL>.

### 3.3 Programming

#### 3.3.1 Initial Setup

For this practical, you have to create your own kernels and the appropriate buffers needed to implement these kernels. Skeleton programs `multiplication.cpp`, `kernel.cl`, and `multiplicationGoldenStandard.cpp` have been provided in the github repository. In order for these programs to run you have to fill in all the missing sections. The missing sections are indicated with a TODO comment and you should use practical 2.1 as a reference.

#### 3.3.2 Skeleton Files

The following files can be found in the github repository:

1. `multiplication.cpp`: This C++ file includes all the setup and running of the openCL program.
2. `kernel.cl`: Contains the kernel needed to complete this practical.
3. `multiplicationGoldenStandard.cpp`: C++ file where you should implement the golden standard.



### 3.4 Desired Programs

Using the given skeleton code you have to implement two matrix multiplications. The multiplication.cpp file already includes two functions createKnownSquareMatrix and createRandomSquareMatrix that automatically create square matrices in single 1D arrays. Please note that the dimensions of these matrices are controlled by the variable Size.

#### 3.4.1 Matrix Multiplication

For the OpenCL kernel you are to implement matrix multiplication as shown below:

$$\begin{array}{ccccccc} 1 & 2 & 3 & & 2 & 4 & 6 & & 1 & 2 & 3 & & 72 & 144 & 216 \\ 1 & 2 & 3 & \times & 2 & 4 & 6 & \times & 1 & 2 & 3 & = & 72 & 144 & 216 \\ 1 & 2 & 3 & & 2 & 4 & 6 & & 1 & 2 & 3 & & 72 & 144 & 216 \end{array}$$

Watch out for potential memory errors and that the sequencing of operations is done correctly to get the right answer.

### 3.5 Report

After completing the code you will have to run tests to do a report on the practical. In this report you should evaluate the openCL program by running tests that you think apply. If you are unable to get the program running you should still complete the report where you indicate what you think the outcomes for the tests would have been given that the code was working. Please note that marks are given for test results, etc. and if your code does not run you will not receive these marks but any explanations will be marked.

#### 3.5.1 Data Transfer Overhead and Speed-up

Speed-up is a concept that was explained in practical 2.1, you should consider looking into the speed-up of the OpenCL GPU implementation and include graphs and tables where applicable.

The asynchronous nature of the OpenCL interface makes it difficult to obtain accurate timing information of the various steps of the process. Try to come up with a way to measure the data transfer overhead and processing time separately.

Use this new information to comment on the sources of OpenCL processing delay. Also comment on the speed-up factor achieved when transfer overhead is not taken into account. Does this relate well to the number of threads that are running on the GPU? If not, provide an argument for why this is the case. Do this for large N (pick a value that takes long enough to dominate the transfer overhead, but does not let you finish a whole coffee between runs).

## 3.6 Helpful Information

This section includes some tips that may be helpful. You probably want to have a quick look over these tips before diving straight into your OpenCL coding.

### 3.6.1 Creating Local Memory on the GPU

In practical 2.1 we setup memory blocks globally resulting in no local memory exclusively accessible within work groups. Meaning we did not use any VRAM (Video Random Access Memory) on the GPU and only used the host computers RAM.

If you would like to setup VRAM for each work group, you do not need to setup any memory buffers as this memory is found on the target device and it is not shared with the host device. Therefore, you just need to indicate in the argument section that such memory block exists, as shown below.

```
//NULL is used for the buffer pointer as there
//is no memory block on the host computer
clSetKernelArg(kernel, ArgumentNumber, MemoryBlockSize, NULL);
```

Then in you kernel you will have to indicate that this argument is a local memory block and not a global memory block, as shown below.

```
__kernel void kernel(__global int* globalMemoryBlockPointer,
                    __local int* localMemoryBlockPointer){
    //Kernel code
}
```

### 3.6.2 Timing

If you would like to get accurate run times consider using the clock\_t objects, as shown below.

```
start = clock(); //start running clock
\\Some code
end = clock();
printf ("Run Time: %0.8f sec \n", ((float) end - start)/CLOCKS_PER_SEC);
```

Furthermore, note that the clFinish command acts as a barrier, stopping the program at this point until everything in the queue has been run.

### 3.7 Submission

Compile your experiments and findings into an IEEE-style conference paper. Make sure you include ALL hardware details, including GPU and CPU clock rates, if available. Of particular importance is the local work group size and number of compute units.

The page limit is 3 pages. Submit your report and code in a zip file to the Vula Assignment for this practical, only one member in the group has to submit. The zipped file should contain the following files:

1. multiplication.cpp
2. multiplicationGoldenStandard.cpp
3. OpenCL/Kernel.cl
4. Report ('Prac2\_2.STUDENTNUMBER'.pdf)

### 3.8 Marks

Note that 33 marks are available, but you will still cannot score a mark above 100%.

Table IV: Prac 2 Marking Guide

Aspect	Description	Mark Allocation
Code		
	Buffer Setup	3
	Kernel	5
	Timers	2
	Golden Standard	3
Report		
	Introduction	3
	Layout/Captions etc	2
	PC details	2
	Graphs/Tables *	6
	Results Comparison	3
	3 Metrics reported on	3
	Overhead discussion	4
	Discussion	4
TOTAL		40

## 4 Prac 3

This practical is split up into two sections, practical 3.1 FPGA simulation introduction and practical 3.2 FPGA Verilog. The simulation introduction covers installing and the basics of Vivado. Whilst, the Verilog section delves into debugging and creating your own module.

### 4.1 Prac 3.1 - FPGA Simulation Introduction

#### 4.1.1 Introduction

The aim of this practical is to familiarise you with Verilog and Xilinx Vivado. Vivado is a software suite produced by Xilinx, that we will be using to perform simulation and HDL programming suited to FPGA-based development.

For this practical you will be doing a tutorial produced by Xilinx to provide you with knowledge of all the main tools and development procedures that are used in the programming and testing of FPGAs. This knowledge is needed for you to complete Prac 3.2 and the YODA project.

Please note that this practical serves as a prelude to practical 3.2; and will count for 10% of your total practical 3 mark and should be completed individually. The reason for doing this individually is so that you can proceed through the tasks at your own pace, and understand them thoroughly rather than relying on a teammate to be responsible for the programming and testing activities.

#### 4.1.2 Required Resources

In order to run the tutorial you will need to have an installation of Vivado please refer the the following site [Xilinx Vivado](#), in order to do so. If you need remote access to a server with Vivado already installed, you can request a login for this.

Once you have a working installation of Vivado, you can download the tutorial here [Vivado Tutorial](#).

The practical files can then be downloaded from this link, [practical resources](#). Please note you will have to login with your Xilinx account to download this file.

Currently, Vivado is compatible with numerous OS's including Ubuntu and Windows 10. The tutorial does not explicitly include issues relating to specifications of the software and hardware requirements; but the practical should run fine on either of these mentioned operating systems.

### 4.1.3 Practical Instructions

For this practical you have to complete Chapters 1 to 3, which is the first 46 pages of the Xilinx Vivado Design Suite Tutorial. Moreover, you will have to make 2 changes from the tutorial so that it falls more in line with this course.

Changes:

- In Chapter 1, step 1.3: save the project as 'practical\_3', not 'project\_xsim'.
- In Chapter 1, step 1.10: install and select Nexys4 DDR, as your board.

The purpose of these changes is to set up the project so that you run the simulation for the FPGA boards used in this course.

### 4.1.4 Submission

There is no report required for this submission. However, the following files have to be submitted in a single compressed file, named STDNUM\_prac3 (eg HLPCHZ001\_prac3.zip).

Located in the downloaded practical file:

- ./ug937-design-files/sources/sinegen.vhd

Files generated whilst completing the practical:

- ./practical\_3/practical\_3.cache
- ./practical\_3/practical\_3.hw
- ./practical\_3/practical\_3.ip\_user\_files
- ./practical\_3/practical\_3.runs
- ./practical\_3/practical\_3.sim
- ./practical\_3/practical\_3.srscs
- ./practical\_3/practical\_3.xpr
- ./practical\_3/tutorial\_1.wcfg

Please do not include practical\_3/practical\_3.gen.

#### 4.1.5 Mark Allocations

Table V: Prac 3.1 mark allocation

<u>Part</u>	<u>Includes</u>	<u>Marks</u>
<b>Files</b>		
	Submission labelled correctly	1
	Correct files submitted	1
<b>Chapter 2</b>		
	IP cores setup correctly	3
<b>Chapter 3</b>		
	Object names correct	1
	Waveforms	2
	Groups	1
	Dividers	1
	Markers	1
	Debugged Waveform	4
	<b>TOTAL</b>	15

The marks schema table is separated into the files that you provide, essentially a mark for having produced the needed files. Chapter 2 involves just setting up the needed IP core and marks for that. Chapter 3 involved a sequence of tasks which are allocated marks according to the effort concerned for completing these.

## 4.2 Prac 3.2 - FPGA Verilog

For this practical the aim is to debug and create your own modules for an FPGA board using verilog (although you can just do everything in simulation, a physical board is not needed). The practical should be completed in pairs, with only one group member submitting a report. The module you will be creating is a amplitude and phase shifting module, that is given an input as a sine wave and imposes an amplitude and phase shift on it. Moreover, you will be editing a sine wave module so that the values are called from BRAM that you setup.

Source files are available on the EEE4120F OCW GitHub:

<https://github.com/UCT-EE-OCW/EEE4120F-Pracs>.

### 4.2.1 Given Modules

1. Sine Generation:

The module called 'sine\_gen.v', reads sine values from a the memory file 'sine.LUT\_values.mem' and outputs one sine values per clock cycle.

2. sine LUT values:

The 1024 sine values are stored in a memory file called, 'sine.LUT\_values.mem'.

3. Test Bench:

A test bench module is uses the sine gen module to generate a sine wave at a specific frequency by controlling the clock speed.

### 4.2.2 Requirements

For this practical you are required to complete the following three tasks and write up one report covering the three tasks.

1. Task 1: Edit the frequency

Currently, the frequency of the sinusoid is 2MHz, when outputted on the test bench. Edit the test bench so that the sine wave frequency is 100MHz. Please include your calculations, code and include a screen shot of the test bench's output.

2. Task 2: BRAM (optional)

The given code calls the sinudoids values from a memory file. For this task you have to import the sine values into BRAM and call the data from BRAM instead of from the memory file. A short tutorial is given below on how to create the BRAM IP and initialize it.

3. Task 3: Amplitude and Shifting Module

Lastly, create a module called 'shifter', that's inputs include the clock, shift and sinusoid from the 'sine\_gen.v' file. Add this module to the test bench with an amplitude increase of 2 and a phase shift of 90°. Then display the shifted sinusoid below, the non-shifted

sinusoid. Note, that you may have to edit some of the inputs and outputs from the given code.

### 4.2.3 Report Requirements

For your final report include a section for each task covering all relevant information, including code, explanations and screenshots of the test bench's output. You should also provide an explanation on what BRAM is and provide an in depth critical analysis of the final system.

### 4.2.4 BRAM setup

1. Select "IP Catalogue" on the left hand side of the IDE
2. Search "BRAM".
3. Under "RAMs & ROMs & BRAM", select the "Block Memory Generator"
4. Select "Port A Settings"
5. Note the write and read width. These relate to the bitwidth of the data in the memory. Change this to 16 for both read and write. We use 16 because the sinusoidal data is in 16 bit hexadecimal values.
6. Browse to the `Prac 3_2_code\srcs\sources_1\new`, and load "sine\_LUT\_values.coe"  
ata is a list of 1024, 16 bit data points.
7. The write and read depth relate to how many samples can be stored. Change this to 1024, as our example sine table has 1024 samples, and this we need 1024 addresses.
8. Select "Other Options"
9. Select "Load Init File"
10. Browse to the `Prac 3_2_code\srcs\sources_1\new`, and load "sine\_LUT\_values.coe"
11. Click "Ok" at the bottom of the dialog
12. Click "Generate" on the next dialog. Press "Ok" when the information dialog shows up.
13. In sources view, there will be a template available to instantiate the IP.
14. Congrats! You've added a 16-bit 1024 sample Block RAM IP!



#### 4.2.5 Rubric

Aspect	Description	Mark Allocation
General	Introduction	3
	Layout	2
	Hardware Details	4
Task 1	Explanation and Calculations	4
	Clock Speed	1
	Frequency	1
	Code	2
	Screen Shot	2
Task 2 (optional)	Explanation on BRAM	4
	Code	3
	Screen Shot	2
Task 3	Module explanation	4
	Correct Phase Shift	3
	Correct Amplitude Increase	3
	Screen Shot	2
	Code	4
	Critical Analysis	6
<b>Total</b>		<b>50</b>

## 5 Prac 4A - FPGA Introduction (Archived)

### 5.1 Introduction

For this practical, it is important to complete the tutorial first to become acquainted with the tools required to complete the practical. It is recommended that you work in groups of three, and that all three members work through the tutorial individually.

It is also imperative that you read through the [wiki.ee](#). Operation of Vivado can be intimidating at first use. The Wiki has been carefully written to include all you need for the practicals, see the page [here](#).

### 5.2 Tutorial

The tutorial is not for marks but it is suggested you complete it, as it will give you the basics of Vivado, and allow you to upload a simple program to the FPGA. In the tutorial (as in [wiki.ee](#)) we use the Digilent Nexys A7 as an example board, but the process should be the same for the Nexys 4 DDR and Nexys 4. The only thing that will change between the three is the board you select when creating the project, and the constraints file you use. Note the constraints file does determine naming of some of the I/O, so double check that.

Most of the details on how to complete these steps are on the wiki under [.](#)

1. Install Vivado

2. Install boards

The only boards you might work with in this course are the Nexys 4, Nexys 4 DDR, and the Nexys A7. So you only need to worry about installing those.

3. Download and save the constraint files. That, or copy the simple example given on the wiki if you're using the A7.

4. Create a new project.

- You can call it "Tutorial".
- Set it as an RTL project and select "Do not add sources at this time". For the prac later, you can add the given source files here.
- Select your board appropriately.

5. Add constraint source

- Right click on constraints, select "Add sources". In the dialog box, make sure "Add or create constraints" is selected. Hit next.
- Select "Create file" if you are copy pasting the constraints in, or "add files" if you have the constraints file downloaded locally. Press "Finish".

- If you were intending on copy-pasting constraints in, do so now. Ensure your constraints are correct for the board, and have the clock, two switches and two LEDs enabled.
  - Right click on the constraints file, and select "Set as Target Constraint File"
6. Add the Verilog source
    - Right click on "Design Sources", select "Add sources". In the dialog box, make sure "add or create design sources" is selected. Hit next.
    - Select "Create file". Call it the same name as your project (good practice for the top level module). Press "Finish".
    - A dialog will open, with ports. You can just press "okay", as we'll define ports in the Verilog code.
    - Right click on the Verilog file, and select "Set as Top" (if it is not already set as top - indicated by being in bold).
    - In it, paste code to, each clock cycle, write switch[0] to LED[0] and the inverse of switch[1] to LED[1]. Example code can be found on the wiki.
  7. Select "Run Synthesis"
  8. Select "Run Implementation"
  9. Select "Generate Bitstream"
  10. Upload your bitstream to your target board (speak to a tutor to get a board)
    - (a) Plug in the board
    - (b) Select "Open Hardware Manager"
    - (c) Select "Open Target" and then "Auto Connect"
    - (d) Vivado should find the board. Select "Program Device" and select the board you've plugged in (The A7 shows as "xc7a100t\_0"). Press the "Program" button.
  11. Hooray! You've created your first FPGA circuit. Toggle the switches to see if it operates as expected. Now let's move on to the fun stuff.

## 5.3 Practical

### 5.3.1 Introduction

In this practical, you will create a digital clock on an FPGA. There is no report for the practical, but you will need to submit screenshots of your testbenches and demonstrate your implementation to a tutor.

Source files are available on the EEE4120F OCW GitHub: <https://github.com/UCT-EE-OCW/EEE4120F-Pracs>.

### 5.3.2 Given Modules

1. TLM  
The top level module, called "Clock.v" in the source files on GitHub, contains the primary logic for your wall clock and allows you to implement I/O and other modules
2. Delay\_Reset  
It's also useful as many components require a set up time. So by using a delayed reset signal, we can cater for reset times of peripherals.
3. Seven-Segment Driver  
This module takes 4 BCD values and displays them on the seven segment display.
4. Decoder  
Used by the Seven-Segment Driver to decode decimal to the appropriate cathode pins.
5. Debounce  
A debounce module you'll need to implement in order to debounce button presses.
6. PWM  
A module you'll need to implement in order to give the seven segment displays changing brightness. This can be tricky, it's suggested you leave it for last.

### 5.3.3 Requirements

The following outcomes are required to pass the demonstration:

1. Implement a simple state-machine to display the real time (hours and minutes) on the 7-segments display. You can start your clock at 00:00 upon reset. Make your clock faster in order to test that the time overflows correctly. Use a 24-hour time format.  
  
The easiest way to do this is with deeply nested if statements. Run a counter that overflows every second and increment the seconds counter on every overflow. Every time this seconds counter equals 59 (i.e. it will overflow on this clock-cycle), increment a minutes units counter. Every time this minutes units counter is about to overflow, increment a minutes tens counter, etc.  
  
Only use non-blocking assignments ( $<=$ ). Blocking assignments ( $=$ ) inside clocked structures are much more difficult to debug. Remember that the entire always block is evaluated at once: the statements are not evaluated sequentially.
2. Display the seconds on the LEDs, in binary format. This is done with a simple assignment outside the always block.
3. Use one of the buttons, properly debounced, to set the minutes. It must increment time by one minute every time it is pressed. Make sure that your time overflows correctly. You do not need to increment the hours when changing the minutes.  
  
It is recommended to write a Debounce module for this. On every clock cycle of the system clock (the fast one), check the state of the button. If it is not the same as the

current module output, change the output and start a dead time counter. While this counter is counting, do not change the output of the module, no matter what the input is doing. Use a deadtime of between 20 ms and 40 ms. To prevent unstable states, register the button before use.

In the clock state machine, you can use a register to store the button's 'previous' state. If the current state is high, and the previous state is low, the button signal went through a rising-edge. Do not use always @(posedge Button) – keep everything in the same clock domain.

4. Use another one of the buttons, properly debounced, to set the hours. It must increment time by one hour every time it is pressed. Make sure that your time overflows correctly.
5. Use the slide-switches to represent a binary "brightness" word. Make use of pulse-width modulation (PWM) to dim the brightness of the LED display.

Ensure that the phasing between the driver signals and the PWM signals is correct. The easiest way to do this is to select a PWM frequency such that the PWM signal goes through exactly one period between driver signal state changes. You can implement this within the SS Driver module.

## 5.4 Mark Allocations

Table VI: Prac 4 mark allocation

		<b>Marks</b>
<b>Testbench</b>		12
	Minute seconds reaching 60 and increasing minutes	
	Minutes reaching 59 and increasing hours	
	Time reaching 23:59 and wrapping back to 00:00	
	(3 each +3 for neatness)	
<b>Demo</b>		
	Time flows correctly (minutes overflow to an increase in hours, and 23:59 flows to 00:00) [6, subtracting 2 for each missed objective]	6
	Time can be scaled through a variable (i.e. the count to increase seconds isn't fixed)	2
	Minute button increases minutes and doesn't increase hours	2
	Hour button increases hours	1
	Debounce module implemented	2
	PWM module	5
	<b>TOTAL</b>	30

**For the 5 marks on PWM:** 1 mark for attempted, 2 marks for reading switches and adjusting, 3 marks for "flashing" implementation, 4 marks for some mix between flashing and decent PWM, 5 marks for correctly implemented.

## 6 Prac 4B - Vivado IP and Resource Usage (archived)

FPGAs are often used in DSP applications due to their highly parallelizable nature and ability to process data at high speeds. Usually FPGAs are used to sample signals and process them, but in this application we're going to generate waveforms and output them to a speaker!

If you aren't all that well acquainted with the physics of music, [this video by Vihart](#)<sup>1</sup> is a really great introduction! She speaks quite quickly, so you may need to watch it twice.

These videos don't have much to do with this prac, but are fun to learn about. If you'd like to learn more about modular synthesis (which is super cool!) check out [this video by Andrew Huang](#)!<sup>2</sup> LookMumNoComputer also does a great bunch of custom modular synthesis projects - just check out his [Furby Organ](#)!<sup>3</sup> He also has guides on how to [build your own modular synth](#).<sup>4</sup>

We'll be implementing a sine wave - which you don't really hear much in music. But as [Composerly shows us](#)<sup>5</sup>, you can create some great tunes with nothing but sine waves and enough processing.

### 6.1 Tutorial

The tutorial starts with guidance on getting familiar with the Xilinx Vivado IP facilities, which leads in to the requirements for the prac.

#### 6.1.1 Resource usage

Please read the [Implement section](#) in the Xilinx Vivado page on the UCT EE Wiki.

#### 6.1.2 Vivado IP

For general information on the Vivado IP, please read the "Xilinx Vivado IP" section of the [Xilinx Vivado wiki page](#).

Of course, when using an IP core, you need to have some knowledge of the technology. You should know how to use the IP cores - and, since you are university students and not necessarily going to just be users - you should also have some sense of what is happening behind the scenes.<sup>6</sup> In the tutorial, you will be guided through the relevant settings and

---

<sup>1</sup><https://youtu.be/i.0DXxNeaQ0>

<sup>2</sup><https://youtu.be/cWslSTTKiFU>

<sup>3</sup><https://youtu.be/GYLBjScgb7o>

<sup>4</sup><https://youtu.be/4Kz8YopLTCQ>

<sup>5</sup><https://www.youtube.com/watch?v=xLtTMkMr2Wg>

<sup>6</sup>You'll get a better sense of the theory behind these concepts from the lectures building on the basics of the PLDs and FPGAs covered in ES1 and ES2.

what they mean, but when using other IP, you should do your research as to what you're adjusting.

Once you've created your project for your board, you can do the following to instantiate memory to hold the lookup table. You can repeat this process for different tables or waves (for example sawtooth, triangle, or perhaps another periodic waveform - maybe the waveform for a specific instrument).

1. Select "IP Catalogue" on the left hand side of the IDE
2. Search "BRAM".
3. Under "RAMs & ROMs & BRAM", select the "Block Memory Generator"
4. Select "Port A Settings"
5. Note the write and read width. These relate to the bitwidth of the data in the memory. Change this to 11 for both read and write. We use 11 because the audio module expects a signal 11 bits big.
6. The write and read depth relate to how many samples can be stored. Change this to 256, as our example sine table has 256 samples, and this we need 256 addresses.
7. Select "Other Options"
8. Select "Load Init File"
9. Browse to the Prac 4 sources, and load LUT\_sinefull.coe
10. Click "Ok" at the bottom of the dialog
11. Click "Generate" on the next dialog. Press "Ok" when the information dialog shows up.
12. In sources view, there will be a template available to instantiate the IP.
13. Congrats! You've added a 11-bit 256 sample Block RAM IP!

It would be a good idea to set up a test bench to ensure you're reading the correct values from the BRAM as you expect. Here's some hints which might make setting up the test bench easier:

- You don't need to do any writing to the BRAM (write enable can be left low)
- You can leave the read enable signal high
- You can just increase the address by 1 each time - it will auto wrap around to 0

## 6.2 Practical

We're going to start by making a simple waveform generator to output middle C at  $f=261.625565\text{Hz}$  (or as close as you can get to it!). Then we're going to try free up some resources in our first implementation. From there, we're going to create an arpeggiator. If you're familiar with synthwave - it's usually an arpeggiator (or *arp*) that creates the repetitive notes in the background.

### 6.2.1 Provided files

You are provided with:

1. `top.v`  
The top level module for the project. A template to get you started.
2. `pwm.v`  
A PWM module to convert the BRAM samples to a PWM signal for the audio jack.
3. `LUT_sinefull.co`  
A full sinewave table, 11 bits wide with 256 samples

### 6.2.2 Creating a waveform generator

In this section you need to create a simple output waveform generator at as close to  $261.625565\text{Hz}$  as you can. The hardware operates at  $100\text{MHz}$ , so  $100\text{MHz}/(261.62556\text{Hz}\cdot 256)$  gives us  $382225.643736$ . But! We have 256 samples in our look up table. So we need to operate 256 times faster to complete one wave in the expected time.  $382225.643736$  divided by 256 is about 1493. So count to that value to produce a tone around middle C.

1. Start by creating a new project, and adding the full sine wave table to BRAM as per the tutorial above. A reminder to use the correct board settings when creating a new project and calling it `fullsine`.
2. Create a new test bench, showing loading a sample from the BRAM
3. Record resource usage for the full sine table (we're looking for total power, LUT, FF, BRAM)
4. Output this wave to the audio port, and record video of it playing. To do so, you need to tie in the `AUD_SD` and `AUD_PWM` signals from your constraints file. `AUD_SD` can be written high (as an enable). `AUD_PWM` needs to be a PWM signal. To generate this, pass the data from the BRAM to an instantiation of `pwm.v`, and pass the output of that PWM module to `AUD_PWM`.



Create a new project called **quartersine** and implement a quarter sine wave table. Record resource usage for this project too, as you will need to compare the implementations later. This resource will be helpful in implementation: <https://zipcpu.com/dsp/2017/08/26/quarterwave.html>. Make sure to remember that you only need  $256/4 = 64$  samples when you instantiate the BRAM on this project.

### 6.3 Create a simple major arpeggiator

A major scale in music sounds “happy”. Minor progressions sound “sad”. We’re going to create a major arpeggio generator. A major arpeggio consists of a base note, a major third and a fifth<sup>7</sup>. As we saw in the Vihart video, this is just maths. If we have a frequency  $f$ , then the major third is just  $f*1.25$ , and the fifth is just  $f*1.5$ . We will finish off the arpeggio by completing the octave, which will be  $f*2$ .

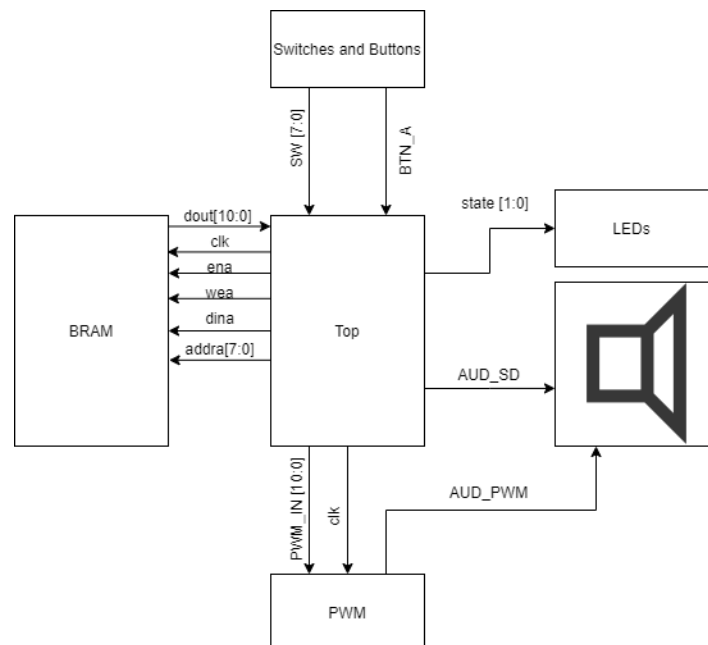


Figure 7: The overview block diagram

So let’s think about what we need to do to implement this:

- We need a counter to change the note every 500ms
- We will need a case statement, which, depending on the note, will change the rate at which we switch through the BRAM address
- We will need to read in switches to add to the base note we’ve chosen (middle C)

<sup>7</sup>[Super simple arpeggio explanation for guitarists](#)

- We will need a state to hold whether we are sending out an arpeggio, or just out base note.

With these items in mind, we know our always block might look something like this. In the code I wrote, I made f\_base the highest frequency (lowest count). The example is very simple, and doesn't cater for switching between f\_base output and arpeggio output. It is just meant to serve as a guide.

Listing 1: Example of “Top” for switching between notes

```
always @(posedge CLK100MHZ) begin
    PWM <= douta; // tie memory to the PWM out

    f_base[8:0] = 746 + SW[7:0]; // get our base frequency

    note_switch = note_switch + 1; // keep track of when to change notes
    if (note_switch == 50000000) begin
        note = note + 1;
        note_switch = 0;
    end

    // Output divider to control frequency
    clkdiv <= clkdiv + 1;

    case(note)
        0: begin // base note
            if (clkdiv >= f_base*2) begin
                clkdiv[12:0] <= 0;
                addra <= addra + 1;
            end
        end
        1: begin //1.5 times faster
            if (clkdiv >= f_base*3/2) begin
                clkdiv[12:0] <= 0;
                addra <= addra + 1;
            end
        end
        2: begin // 1.25 faster
            if (clkdiv >= f_base*5/4) begin
                clkdiv[12:0] <= 0;
                addra <= addra + 1;
            end
        end
        3: begin //2 times faster
            if (clkdiv >= f_base) begin
                clkdiv[12:0] <= 0;
                addra <= addra + 1;
            end
        end
    end
```

```

        end
        default: begin // Don't know what's happening, just output middle C
            if (clkdiv >= 1493) begin
                clkdiv[12:0] <= 0;
                addra <= addra +1;
            end
        end
    end
endcase;
end

```

## 6.4 Hand In

Submit the following as a single PDF:

1. Show evidence that the two implementations produce the same results. To do this:
  - Show screen captures of the test benches at  $t=0$ ,  $t= 1/2 \pi$  and  $t = \pi$  (we're looking to see the changes in values in the output sine value at these points.) [6 marks]
  - To further elaborate on the above - we want to see that the value passed to the PWM module on either side of those values of  $t$  follow the same progression.
2. List the resource and power usage of the both implementations [5 marks]
3. Write a paragraph or two on resource consumption of the FPGA. Talk about resource availability, resources used, effort in terms of implementation [5 marks]

Submit a video on YouTube showing the arpeggiator working. Make sure to have the volume loud enough. Please describe what you are showing in the video. [10 marks] The video should show the following:

- In base.f mode, toggle the switches to show that the frequency changes.
- Enable arpeggio mode and show that works by recording the output from a speaker
- Change the output frequency while in arpeggiator mode

## References

- [1] “Opencl - the open standard for parallel programming of heterogeneous systems,” Jul 2013. [Online]. Available: <https://www.khronos.org/opencl/>