



# UPPSALA UNIVERSITET

Accelerator-Based Programming

Assignment 3  
A TENSORFLOW EXAMPLE

Jinglin Gao

September 27, 2023

## 1 Compute time of GPU node on UPPMAX

In this task, we use a GPU node on Uppmax to compute the implementation of Conway's Game of Life. First we use 1000 iterations and verify the code's correctness. The result is shown below:

```
(tf) [jinglin@s152 ass-3l$ python3 assignment3_template.py
2023-09-27 23:18:13.531800: W tensorflow/compiler/tf2tensorrt/utils/py_utils
.cc:38] TF-TRT Warning: Could not find TensorRT
2023-09-27 23:18:16.846856: I tensorflow/core/common_runtime/gpu/gpu_device.
cc:1639] Created device /job:localhost/replica:0/task:0/device:GPU:0 with 13
775 MB memory: -> device: 0, name: Tesla T4, pci bus id: 0000:08:00.0, comp
ute capability: 7.5
2023-09-27 23:18:46.742106: I tensorflow/compiler/xla/stream_executor/cuda/c
uda_dnn.cc:432] Loaded cuDNN version 8600
Compute time: 32.43176509696059
Cells alive at start: 603
Cells alive at end: 2658
2658
```

Figure 1: Compute time on GPU node

From the result, we can see that the compute time is 32.43 seconds. The result after 1000 iterations is 2658, which indicates my implementation is correct.

## 2 Compute time of CPU node on UPPMAX

Next, we disable the GPU by adding two lines to the start of the node and compute the time needed to finish 1000 iterations again. The result is shown below:

```
2023-09-27 23:26:02.763471: I tensorflow/core/common_runtime/eager/execute.c
c:1678] Executing op __inference_runlife_27006 in device /job:localhost/repl
ica:0/task:0/device:CPU:0
Compute time: 69.00890642288141
Cells alive at start: 603
Cells alive at end: 2658
2658
```

Figure 2: Compute time on CPU node

So when we disable the GPU node we can notice that there is a significant performance drop. Now we need 69.01 seconds to finish 1000 iterations on the CPU node, which is around  $2.15\times$  slower than the GPU node.

## 3 Switch Data Type and Compare Run Time

Data Type	Run Time	
	GPU node (sec)	CPU node (sec)
float16	32.43	69.01
float32	34.56	77.89
int32	48.69	85.16
bfloat16	34.42	61.38

Generally, we can expect that GPU will provide better performance than CPU, and our experiments align with this.

GPUs are optimized for certain data types and operations. Here we use Tesla T4, which is optimized for FP16 and FP64, so here in the result, we can see that we get better performance for **float** datatype than **int** datatype on GPU. If we change the iteration times from 1000 to 3000, we will see a more significant difference between the compute time of float16 and float32, which is 101.63 and 108.77, respectively.

In conclusion, GPU has better performance than CPU in general. But for different data types, GPU performs better for some specific data types since it is optimized for some specific computing.