Processors

Ehsan Yousefzadeh-Asl-Miandoab ehyo@itu.dk April-2022

• Computing

- Processors
- How do electrons work for us?!
- Tradeoff of processors
- CPU
- GPU
- FPGA
- Accelerator



What Computers Do Computing 3 ********** Processor (CPU) device(s) Output device input Process Produce Out Information USING COMPUTER GAINING TRANSFORMING **TECHNOLOGY TO** THE INPUT DATA KNOWLEDGE COMPLETE A INTO OUTPUT (INSIGHT) **GOAL-ORIENTED** DATA. TASK. -

Richard Hamming (1915 - 1998)



- Computing
- Processors
- How do electrons work for us?!
- Tradeoff of processors
- CPU
- GPU
- FPGA
- Accelerator

Processors

- Making Computing possible
- Making electrons to work for us to do the computation



- Computing
- Processors
- How do electrons work for us?!
- Tradeoff of processors
- CPU
- GPU
- FPGA
- Accelerator





Examples:

- 1. Calculating the First n Numbers of Fibonacci Series
- 2. Train a Machine Learning Model





finclude <iostream>
using namespace std;

int main() {
 int n, t1 = 0, t2 = 1, nextTerm = 0;

cout << "Enter the number of terms: "; cin >> n;

```
cout << "Fibonacci Series: ";</pre>
```

```
for (int i = 1; i <= n; ++i) {
    // Prints the first two terms.
    if(i == 1) {
        cout << t1 << ", ";
        continue;
    }
    if(i == 2) {
        cout << t2 << ", ";
        continue;
    }
    nextTerm = t1 + t2;
    t1 = t2;
    t2 = nextTerm;
    }
}</pre>
```

cout << nextTerm << ", ";
}
return 0;</pre>

```
# Program to display the Fibonacci sequence up to n-th term
nterms = int(input("How many terms? "))
# first two terms
n1, n2 = 0, 1
count = 0
# check if the number of terms is valid
if nterms <= 0:
  print("Please enter a positive integer")
# if there is only one term, return n1
elif nterms == 1:
  print("Fibonacci sequence upto", nterms, ":")
  print(n1)
# generate fibonacci sequence
  print("Fibonacci sequence:")
  while count < nterms:</pre>
      print(n1)
      # update values
      n1 = n2
      n2 = nth
       count += 1
```



















d1	d2	sum	Carry
0	0	0	0
0	1	1	0
1	0	1	0
1	1	0	1











- Computing
- Processors
- How do electrons work for us?!
- Tradeoff of processors
- CPU
- GPU
- FPGA
- Accelerator

Processors

- Making Computing possible
- Making electrons to work for us to do the computation



Tradeoff

Processor	Programmability	Goal Program	Flexibility (Different Programs)	Promised Performance
Central Processing Unit (CPU)	Easy	Latency Oriented	Super High	Fair
Graphical Processing Unit (GPU)	Medium	Throughput Oriented	High	Medium
Field Programmable Gate Array	Hard	Both	Low	High
Accelerator	Super Hard	Both	Super Low	Super High

- Computing
- Processors
- How do electrons work for us?!
- Tradeoff of processors
- CPU
- GPU
- FPGA
- Accelerator





Modern CPUs



- They fetch and execute more than one instruction (a windows of instruction)
 - Higher throughput
- Advanced Hardware Execution Mechanisms to execute faster
- Employ Cache Hierarchy to fill the Memory-Processor performance gap
 - Temporal/ Spatial Locality
- They have several cores (parallel computing)





Cache Hierarchy

Less Access latency More Data Locality

Less Storage Capacity More Expensive per bit





Overall Time of Executing three Instruction: 3 * (1ns + 1ns + 2ns + 1ns + 2ns) = 3 * 7ns = 21ns

Pipelining

• Pipelined Basic Processor





Inst1

Inst2

Inst3

= (2ns + 2ns + 2ns + 2ns + 2ns) + 2ns + 2ns= 10ns + 2ns + 2ns = 14ns



- Computing
- Processors
- How do electrons work for us?!
- Tradeoff of processors
- CPU
- GPU
- FPGA
- Accelerator



GPU

- Parallel Processor
- Throughput
- Low Working Frequency























- Computing
- Processors
- How do electrons work for us?!
- Tradeoff of processors
- CPU
- GPU
- FPGA
- Accelerator

Field Programmable Gate Array (FPGA)

First, used for Prototyping by Electronic people

High-Level Language

Very Low-Level Language (requires Digital Electronics knowledge)

Field Programmable Gate Array (FPGA)

First, used for Prototyping by Electronic people

High-Level Language

Very Low-Level Language (requires Digital Electronics knowledge)

FPGAs in the Age of Al

- Supply the practitioners with the highest parallelism
- Built a specific Neural Network on Hardware
- Change it whenever you decide

FPGAs in the Age of Al

- Supply the practitioners with the highest parallelism
- Built a specific Neural Network on Hardware
- Change it whenever you decide

FPGAs in the Age of Al

- Supply the practitioners with the highest parallelism
- Built a specific Neural Network on Hardware
- Change it whenever you decide

FPGAs in the Age of AI

- Supply the practitioners with the highest parallelism
- Built a specific Neural Network on Hardware
- Change it whenever you decide

- Computing
- Processors
- How do electrons work for us?!
- Tradeoff of processors
- CPU
- GPU
- FPGA
- Accelerator

Accelerator

• Application Specific Integrated Circuit (ASIC)

- Everything is specified by the designers
- No flexibility is expected but depends on the designers!

Google

Tensor Processing Unit

Cerebras Wafer Scale Engine

Cerebras WSE 1.2 Trillion Transistors 46,225 mm² Silicon

Largest GPU 21.1 Billion Transistors 815 mm² Silicon

Conclusion

- Computation
- Processors make the computation possible
- Tradeoff of different processors
- CPU
- GPU
- FPGA
- Accelerator

Questions?

Thanks for your attention!