

# OPENCL ON THE FPGA

SO FOUNDERS CAN DRINK BEER NOW

Huh?



Huh?



# OPENCL ON THE FPGA

SO FOUNDERS CAN DRINK BEER NOW



Huh?

# HI! I'M VINCENT HINDRIKSEN

- FOUNDER OF A COMPANY CALLED STREAMCOMPUTING
- MY COLLEAGUES KNOW HOW TO MAKE SOFTWARE FAST, REALLY FAST!
- WE USE EVERYTHING, FROM ALGORITHM-OPTIMISATIONS TO DUCT-TAPE.
- ONE OF THE TARGET DEVICES WE USE ARE FPGAS

# SEE OUR HOMEPAGE FOR BRAGGING POSTS ON SPEEDUPS WE GOT

The screenshot shows the Stream Computing website header with the logo and navigation menu. The main article title is "Porting Manchester's UNIFAC to OpenCL@XeonPhi: 160x speedup". Below the title, it says "Written by Vincent Hindriksen" and "November 04, 2015". The article text begins with "As we cannot use the performance results for most of our commercial projects because they contain sensitive data, we were happy that Dr. David Topping from the University of Manchester was so kind to allow us to share the data for the UNIFAC project. The goal for this project was simple: port the UNIFAC algorithm to the Intel XeonPhi using". A line graph is partially visible, showing performance metrics for different configurations.

← From 33 seconds to 0.068 seconds

# OPENCL – THE COMPUTE LANGUAGE

- YOUNG LANGUAGE, STARTED ONLY IN 2009
- ENABLES COMPUTE PARALLELISM AND EXPLICIT MEMORY MANAGEMENT IN ONE C-LIKE LANGUAGE.
- OPEN STANDARD SUPPORTED BY MANY BIG COMPANIES, LIKE AMD, NVIDIA, INTEL, ALTERA, XILINX, ARM, VIVANTE, IMAGINATION, TEXAS INSTRUMENTS, QUALCOMM, ETC.
- NOW AT VERSION 2.1, BUT NO DRIVERS AVAILABLE YET.
- OPENCL 1.2 IS THE MOST USED VERSION, WHERE 2.0 IS GETTING TRACTION QUICKLY.

# OPENCL – THE COMPUTE LANGUAGE

- WORKS ON CPUs – YOU KNOW
- WORKS ON GPUS – GRAPHICS UNITS THAT NOW ALSO PROCESS MATRICES
- WORKS ON DSPs – DIGITAL SIGNAL PROCESSORS – LIKE YOUR SOUND CARD
- WORKS ON FPGAs – OUR SUBJECT
- WORKS ON STEROIDS – BECAUSE IT'S FAST

# OPENCL PARALLELISM IN 3 SIMPLE STEPS

- TAKE A SMALL FUNCTION (KERNEL) THAT OPERATES ON ONE DATA ELEMENT
- DEFINE THE DATA TO OPERATE ON
- COMBINE THE ABOVE TWO AND TADA!
- (THEN OPTIMISE THE HELL OUT OF IT, WHICH TAKES MOST OF THE EFFORT)

APPLY  
KERNEL-FUNCTION  
TO EACH  
ELEMENT

2





# WHAT IS AN FPGA?

- IT'S A PROCESSOR WHERE THE PROGRAM DOESN'T RUN ON IT, BUT THE PROCESSOR IS THE SOFTWARE ITSELF. SO... IT IS ACTUALLY HARDWARE-ONLY PROGRAMS.
- IT RUNS PROGRAMS WITH VERY LOW LATENCY – THE TIME TO PROCESS SIMPLE INPUT CAN BE DONE UNDER 0.1 MICROSECOND = 0.0001 MILLISECOND = 0.0000001 SECOND.
- IT'S MUCH COOLER THAN A GPU. MOSTLY BECAUSE IT USES 10X LESS POWER.
- MANY BRANDS, TWO BIGGEST ARE ALTERA AND XILINX

# WHY PROGRAM ALSO FPGAS WITH OPENCL?

- STANDARD LANGUAGES VHDL AND VERILOG ARE FOR DESIGNING IP (FUNCTIONALITY)
- OPENCL IS FOR DESIGNING FULL SYSTEMS (CONNECT IP TO RESOURCES LIKE MEMORY)
- OPENCL IS GOOD IN DEFINING PARALLELISM

# OPENCL ON FPGAS

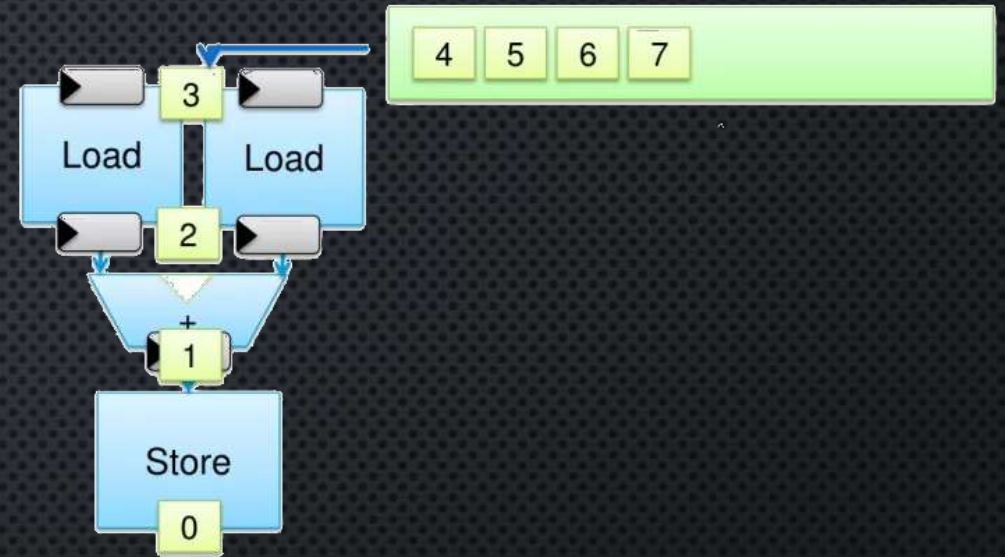
- ONLY ON ALTERA AND XILINX
- ALTERA IS SLOWLY GETTING TO VERSION 1.2
- XILINX IS STILL PRETTY MUCH BETA-STATE, BUT IT WORKS.

# STANDARD “SIMD” PARALLELISM

- CORE 1 PROCESSES TASK 1
- CORE 2 PROCESSES TASK 2
- CORE 3 PROCESSES TASK 3
- CORE 4 PROCESSES TASK 4
- ETC

THING IS: FPGA DOESN'T HAVE  
“CORES”. SO WHAT NOW?

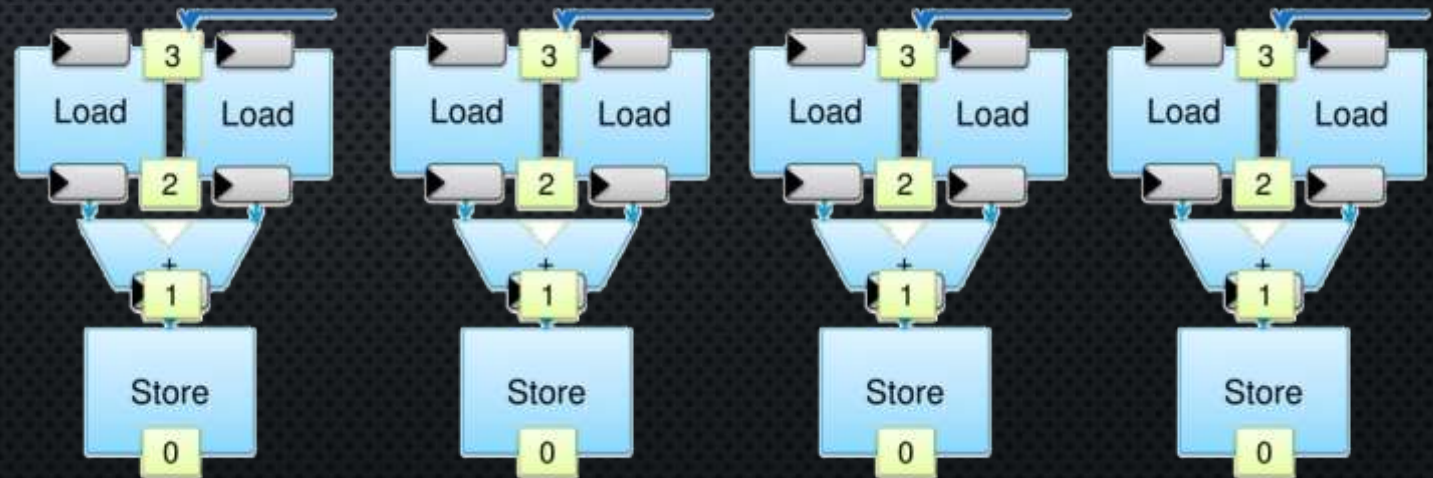
# PARALLELISM ON FPGAS 1/2



- FPGA USES PIPELINE
- EACH TASK IS PIPED THROUGH. WHEN A PHASE IS FINISHED, ALL TASKS HOP FORWARD ONE POSITION.
- IF THERE IS SPACE, MULTIPLE IDENTICAL PIPELINES CAN BE CREATED.

# PARALLELISM ON FPGAS 2/2

- MULTIPLE PIPELINES



# EXAMPLE! DEVICE-SIDE

```
__KERNEL VOID VECTOR_ADD(__GLOBAL CONST FLOAT *X,  
                        __GLOBAL CONST FLOAT *Y,  
                        __GLOBAL FLOAT *RESTRICT Z)  
{  
    // GET INDEX OF THE WORK ITEM  
    INT INDEX = GET_GLOBAL_ID(0);  
    // ADD THE VECTOR ELEMENTS  
    Z[INDEX] = X[INDEX] + Y[INDEX];  
}
```

# HOST-SIDE BOILERPLATE (STANDARD)

1. GET A LIST OF AVAILABLE PLATFORMS
2. SELECT DEVICE
3. CREATE CONTEXT
4. CREATE COMMAND QUEUE
5. CREATE MEMORY OBJECTS
6. READ KERNEL FILE
7. CREATE PROGRAM OBJECT
8. COMPILE KERNEL
9. CREATE KERNEL OBJECT
10. SET KERNEL ARGUMENTS
11. EXECUTE KERNEL (ENQUEUE TASK)
12. READ MEMORY OBJECT
13. FREE OBJECTS



# HOST-SIDE BOILERPLATE (WHEN FPGA-ONLY)

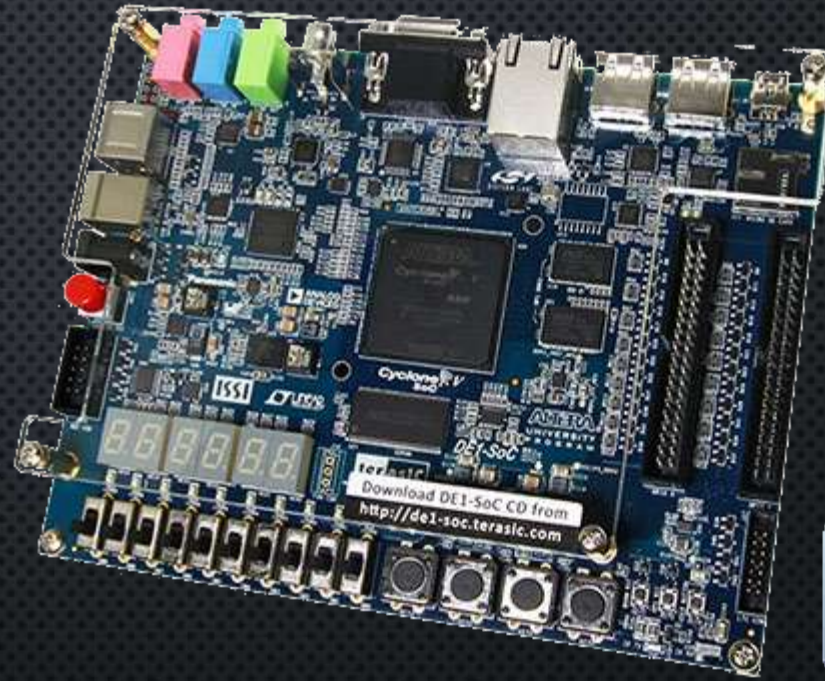
1. SELECT FPGA-DEVICE
2. CREATE CONTEXT AND COMMAND QUEUE
3. CREATE MEMORY OBJECTS
4. READ PRE-COMPILED KERNEL FILE
5. CREATE KERNEL OBJECT OF PRE-COMPILED KERNEL
6. SET KERNEL ARGUMENTS
7. EXECUTE KERNEL (ENQUEUE TASK)

8. READ MEMORY OBJECT
9. FREE OBJECTS

OR:

1. COPY HOST-CODE FROM PREVIOUS PROJECT
2. ALTER TILL IT WORKS

# TRY YOURSELF FPGA



- CHEAP BOARD - \$250
- GIVES ONE YEAR LICENCE (NOT SURE IF STILL APPLICABLE...)
- IT RUNS "HELLO WORLD" AND MORE!

# WANT TO KNOW MORE?

- UHM, UH... GOOGLE?
- FOLLOW @OPENCLONFPGAS ON TWITTER
- CHECK WEBSITES OF ALTERA AND XILINX
- ASK US

