

High Performance Parallel Systems

12 – Computing with accelerators

Kenneth Skovhede & Carl-Johannes Johnsen



What is an accelerator?

- Accelerators are a type of hardware, which is specialized for one particular type of task. In other words, it can accelerate a task.
- They are usually external components, which are connected to a host system. In most cases, the host is a generic CPU/RAM model.
- Due to their specialization, they usually have lower performance than general-purpose hardware for tasks, which do not fit their pattern.



Why use accelerators?

- The use of the Von-Neumann architecture introduces a bottleneck
 - Instructions need to be fetched over same memory channel
 - Only one instruction at a time
- Specialized application
 - Applications that rely heavily on a particular feature, say matrix multiplication, can increase performance if there is special hardware for this operation, instead of using a sequence of CPU instructions
- No need for general-purpose functionality
 - If the application does not need the flexibility of the CPU, those unused elements of the CPU are essentially wasted



Examples of accelerators

- Graphics card
- Bitcoin miner
- TPU
- TCP offload engine
- RAID controller
- Cryptography (AES)
- Video encode/decode (h.264)



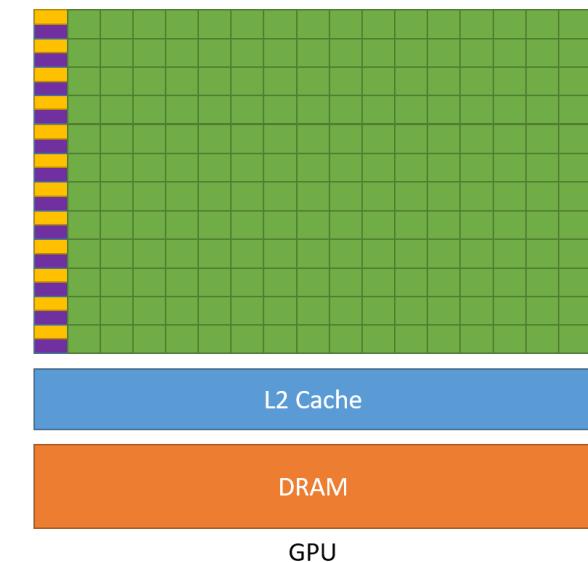
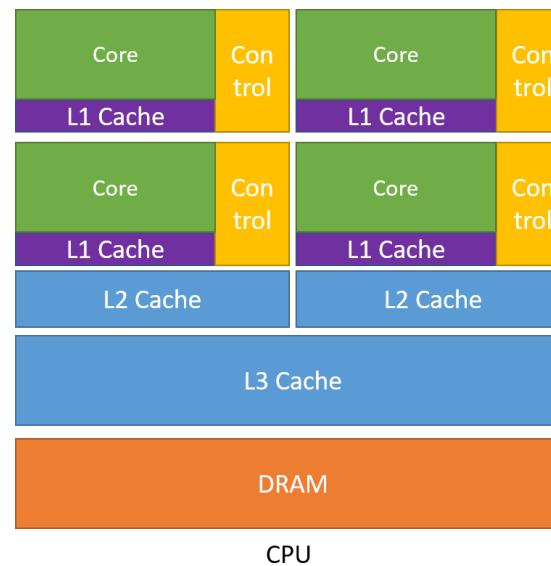
Graphics card

- They were introduced in the 80's, as graphics became more and more popular for computers.
- They supported a wide range of linear algebra methods, as this is what makes up computer graphics.
- They discovered that these computations could be made in parallel and that they didn't need to be flexible in the types of computations.
- Due to the memory access pattern of graphics processing usually being sequential, they can achieve a better memory bandwidth, compared to the CPU RAM.



Graphics card – The good

- As they grew more and more popular, graphics cards became faster and cheaper.
- Their performance increased by leveraging the parallelism through a lot of computation cores, or shaders, and by increasing memory bandwidth.
- All these hundreds, and even thousands, of cores run in lock-step. This means that every core runs the exact same program at the same time.
- They can be more energy efficient, compared to CPUs, due to more of the circuitry being dedicated for compute, rather than branch prediction, random access, etc.



Graphics card – The bad

- They are not good at single/lightly threaded programs.
- They are not good at random access memory patterns.
- They are not good at small problems, as there is an overhead when using external components.
- They are not good at branching.
- They are not good at running multiple different programs.
- They are not as easy as CPUs for getting “quick” results.



OpenCL / CUDA

- As graphics card evolved, they started to have quite powerful computing capabilities.
- Traditionally when using a graphics card, one had to program them using graphical operations. Generally, this was done through shaders, where you specify how to draw and color graphics.
- This wasn't a great fit for many general-purpose programs, as this required a substantial amount of work.
- To solve this, two frameworks were introduced: OpenCL by the Khronos group and CUDA by NVIDIA.



OpenCL / CUDA

- Both frameworks uses the host/device model, where the programmer must write a host program and a device program.
- The device program is made up of one or more compute kernels.
- Each kernel is responsible for moving memory throughout the device memory hierarchy.
- The host program is responsible for moving memory to/from the host and device memory, and for invoking the device kernels.
- The host program can be written in C or C++ and as such also any language which can utilize C bindings, such as Python.
- Kernels are written in a C derivation in OpenCL and extended C++ in CUDA.



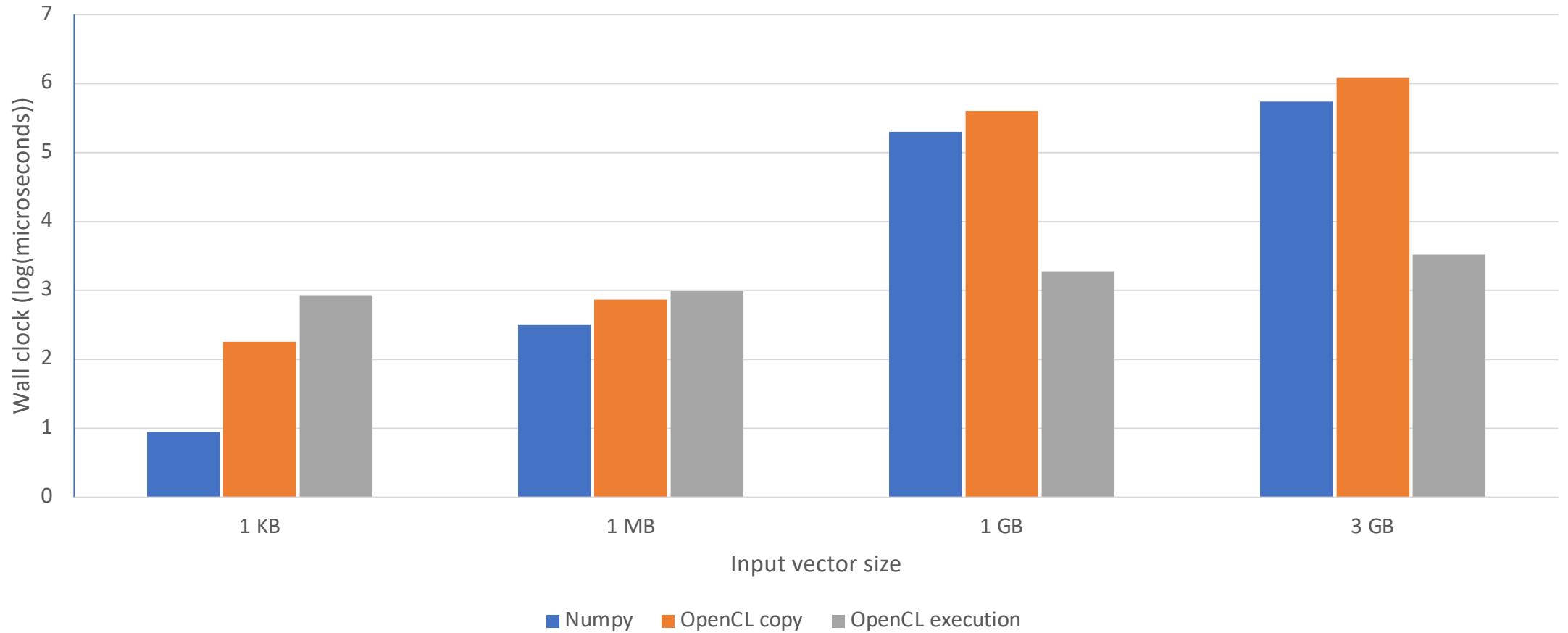
OpenCL example – vector add

```
1 import numpy as np
2 import pyopencl as cl
3
4 # Generate input data
5 n = np.uint32(1024) # 4 KB
6 a = np.random.rand(n).astype(np.float32)
7 b = np.random.rand(n).astype(np.float32)
8 c = np.empty_like(a)
9
10 # Compute the expected result using CPU
11 c_expected = a + b
12
13 # Create the OpenCL context
14 context = cl.create_some_context()
15 queue = cl.CommandQueue(context)
16
17 # Allocate device memory
18 mf = cl.mem_flags
19 a_device = cl.Buffer(context, mf.READ_ONLY, a.nbytes)
20 b_device = cl.Buffer(context, mf.READ_ONLY, b.nbytes)
21 c_device = cl.Buffer(context, mf.WRITE_ONLY, c_expected.nbytes)
22
23 # Load and compile the kernel
24 with open('vector_add.cl', 'r') as f:
25     kernel_source = f.read()
26 program = cl.Program(context, kernel_source).build()
27
28 # Copy input to device
29 cl.enqueue_copy(queue, a_device, a)
30 cl.enqueue_copy(queue, b_device, b)
31
32 # Execute kernel
33 program.vector_add(queue, a.shape, None, a_device, b_device, c_device, n)
34
35 # Copy result from device
36 cl.enqueue_copy(queue, c, c_device)
37
38 # Check results
39 print ('Difference:', np.linalg.norm(c - c_expected))
```

```
1 __kernel void vector_add(
2             __global float *a,
3             __global float *b,
4             __global float *c,
5             const unsigned int n) {
6
7     // Get the thread id
8     int id = get_global_id(0);
9
10    // Don't go out of bounds
11    if (id < n)
12        c[id] = a[id] + b[id];
13 }
```



OpenCL example – performance



Application Specific Integrated Circuit (ASIC)

- All the examples given are ASICs, where the circuit is final once printed.
- This will always give the best performance, compared to general-purpose hardware, as one could make said hardware in the worst case.
- However, the starting price is 500 000 \$, where one has to order 1000s of samples.
- If there is an error in the design the samples are useless, as they cannot be modified once printed.



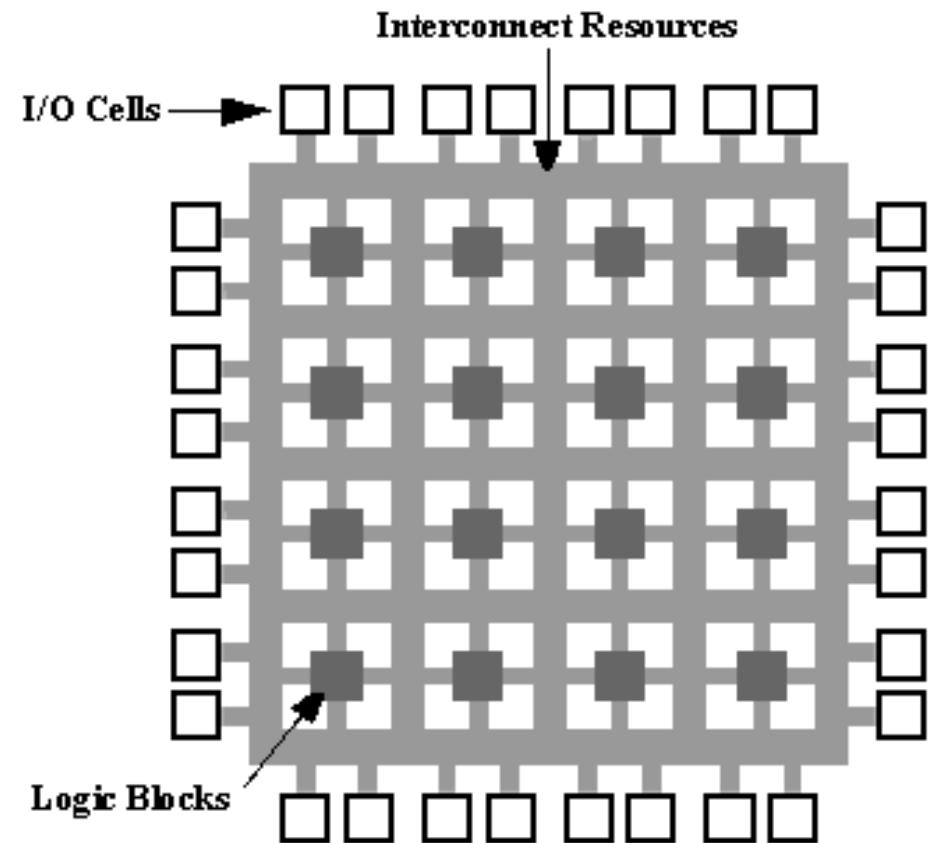
Reconfigurable hardware

- Reconfigurable hardware is the middleground between general-purpose hardware and ASICs.
- Rather than supporting an instruction set, the programmer has to describe the circuit performing the computations. This enables better utilization, as only the circuit needed is implemented, but becomes harder to implement complex designs.
- Their design can be reconfigured multiple times, which make them more flexible than ASICs. This also lowers the performance, compared to ASICs.



Field-Programmable Gate Array (FPGA)

- FPGAs are an example of reconfigurable hardware.
- It consists of a big interconnect and a grid of components, such as logic gates, Look-Up Tables (LUT), multiplexors and Flip-Flop (FF) registers.
- Hardware is designed through how these components are connected through the interconnect. The resulting implementation is equivalent to a printed circuit of the same design, although at a performance hit.
- FPGAs are generally more efficient at tasks, compared to general-purpose hardware, due to the same benefits as GPUs for compute.



Working with hardware models

- While parallel programming becomes harder, as multiple events occur at the same time, programming hardware models become even harder as there is a temporal dimension.
- Events occurs concurrently, and signals take time to propagate on the wires.
- Most designs handles this signal propagation, by guarding all the registers with a clock signal.
- This clock signal is then used to synchronize events, so the programmer can better control the order of events.



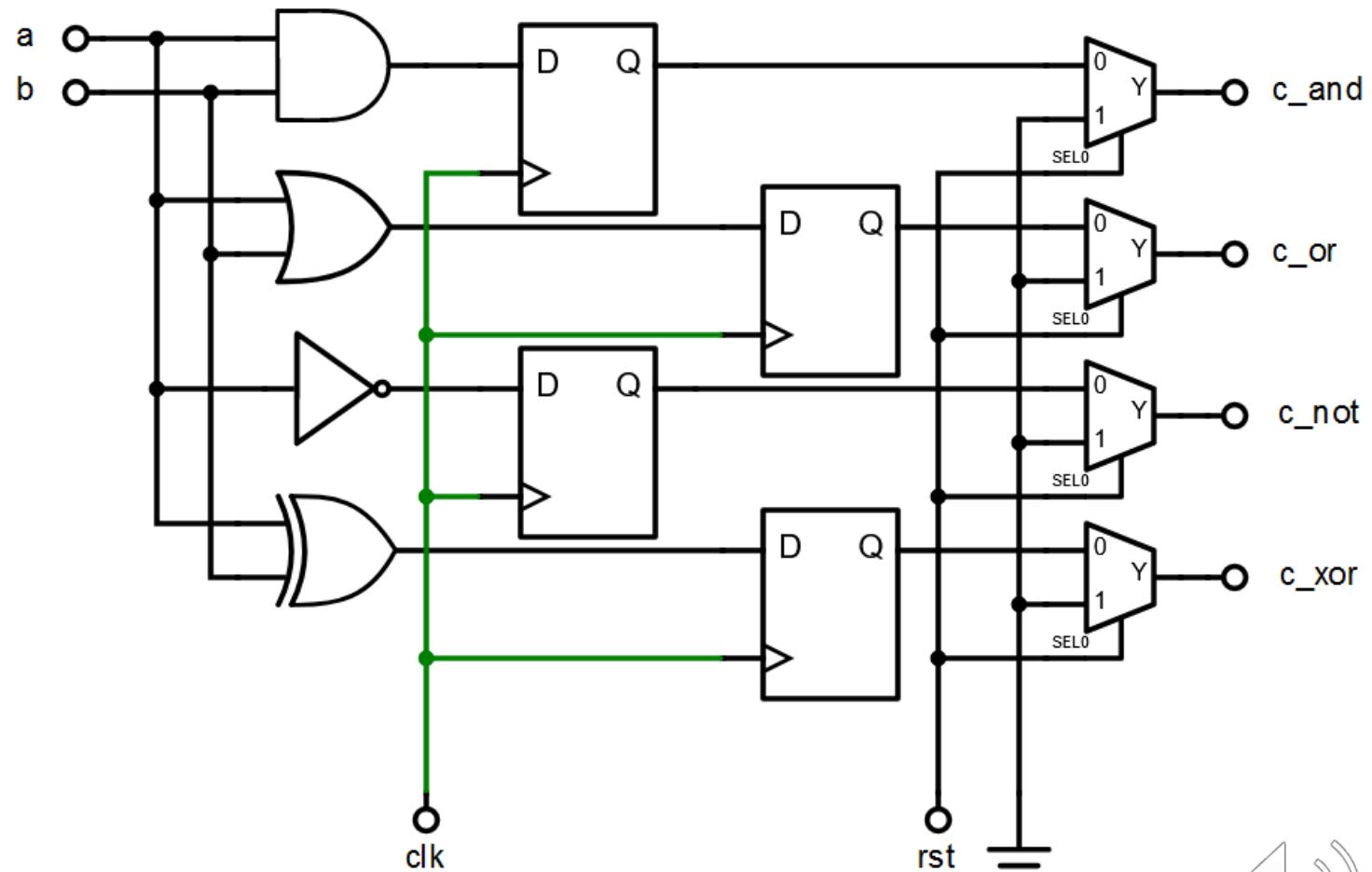
Hardware Description Language (HDL)

- FPGAs are programmed using HDLs, such as VHDL and Verilog.
- Both languages are parallel languages for describing circuits.
- They are very tedious to work with, as they are very verbose and very low-level.
- Even though the languages have evolved, FPGA vendors are slow to adopt these changes / additions to the languages.



VHDL example – Basic logic gates

```
1 library IEEE;
2 use IEEE.STD_LOGIC_1164.ALL;
3 use IEEE.NUMERIC_STD.ALL;
4
5 entity logic_gates is
6   port(
7     a: in std_logic;
8     b: in std_logic;
9
10    c_and: out std_logic;
11    c_or: out std_logic;
12    c_xor: out std_logic;
13    c_not: out std_logic;
14
15    rst: in std_logic;
16    clk: in std_logic
17  );
18 end logic_gates;
19
20 architecture RTL of logic_gates is
21
22  process(clk, rst)
23  begin
24    if RST = '1' then
25      c_and <= '0';
26      c_or <= '0';
27      c_xor <= '0';
28      c_not <= '0';
29    elsif rising_edge(CLK) then
30      c_and <= a and b;
31      c_or <= a or b;
32      c_xor <= a xor b;
33      c_not <= not a;
34    end if;
35  end process;
36
37 end RTL;
```



Test benches

- In order to test an implementation written in HDLs, test benches are needed.
- Generally, they are an additional program, which drives the input wires and verifies whether the output wires behave as expected.
- They have to be written in HDLs as well, although they are no longer limited to a synthesizable subset of the language.



VHDL testbench example – Basic logic gates

```
1 library IEEE;
2 use IEEE.STD_LOGIC_1164.ALL;
3 USE IEEE.NUMERIC_STD.ALL;
4
5 entity logic_gates_tb is
6 end logic_gates_tb;
7
8 architecture tb of logic_gates_tb is
9     signal a, b : std_logic;
10    signal c_and, c_or, c_not, c_xor : std_logic;
11
12    signal clk, rst : std_logic;
13    signal stop_clock : boolean := false;
14 begin
15     inst_logic_gates : entity work.logic_gates
16     port map(
17         a => a,
18         b => b,
19         c_and => c_and,
20         c_or => c_or,
21         c_not => c_not,
22         c_xor => c_xor,
23         clk => clk,
24         rst => rst
25 );
26
27 clk_proc: process
28 begin
29     while not stop_clock loop
30         clk <= '1';
31         wait for 5 ns;
32         clk <= '0';
33         wait for 5 ns;
34     end loop;
35     wait;
36 end process;
37
38 tester: process
39 begin
40     rst <= '1';
41     wait for 5 ns;
42     rst <= '0';
43
44     wait until rising_edge(clk);
45     a <= '0';
46
47     b <= '0';
48     wait until rising_edge(clk);
49     a <= '0';
50     b <= '1';
51
52     wait until falling_edge(clk);
53     assert(c_and = '0');
54     assert(c_or = '0');
55     assert(c_not = '1');
56     assert(c_xor = '0');
57
58     wait until rising_edge(clk);
59     a <= '1';
60     b <= '0';
61
62     wait until falling_edge(clk);
63     assert(c_and = '0');
64     assert(c_or = '1');
65     assert(c_not = '1');
66     assert(c_xor = '1');
67
68     wait until rising_edge(clk);
69     a <= '1';
70     b <= '1';
71
72     wait until falling_edge(clk);
73     assert(c_and = '0');
74     assert(c_or = '1');
75     assert(c_not = '0');
76     assert(c_xor = '1');
77
78     wait until rising_edge(clk);
79
80     wait until falling_edge(clk);
81     assert(c_and = '1');
82     assert(c_or = '1');
83     assert(c_not = '0');
84     assert(c_xor = '0');
85
86     stop_clock <= true;
87     wait;
88 end process;
89 end architecture tb;
```



VHDL testbench example – Basic logic gates

