# 22T2 COMP4601 – Design Project B

LiC: Oliver Diessel

# What "Design Project B" is about

- Many things: capstone course, reconfigurable systems, use of high-level synthesis (HLS), research skills
  - Lectures to discuss FPGA design using HLS
  - Labs to gain experience using an HLS tool and further develop design skills
  - Projects to provide opportunity for group-based open-ended investigation
  - Seminars to practice research and presentation skills
- Resources available: see the course website www.cse.unsw.edu.au/~cs4601

# Assessment

Assessment in this course will be based on the following:

| Criterion | Week due | Contribution |
|---|---|---|
| **Lab component** | | **40%** |
| Lab hand ins | Mondays, Weeks 3, 4, 6 & 8 | 40% i(ndividual) |
| **Seminar component** | | **20%** |
| Seminar presentations | Wednesday lectures, Weeks 7-10 | 10% T(eam) |
| Seminar participation & assessment | Wednesday lectures, Weeks 7-10 | 10% i(individual) |
| **Project component** | | **40%** |
| Individual project contribution | | 10% i(individual) |
| Project plan (presentation & report) | Lab sessions, Week 5 | 10% T(eam) |
| Final project presentation, demonstration & report | Lab sessions, Week 10 | 20% T |

Individual assessments - worth 60% of your assessment in total - will be based on individual efforts, individual contributions to team successes, your learning, and ability to adapt to circumstances.

# Course expectations

1. Read the relevant chapter before each lecture in Weeks 1 – 5

2. Labs (cover quite a bit of ground, so please prepare each week):
   — Complete Vivado tutorial chapters to gain exposure to tool use & familiarity
   — Handins due Mondays at 5pm in Weeks 3, 4, 6 & 8 (each worth 10% of course mark)
   — Start in Week 1 – make sure you have your lab setup working correctly

3. Projects – worth 40%
   — Form a group of 3 local or 3 local +1 remote team members by end of Week 2
     • Decide on spokesperson and let Oliver know who is on your team
   — Agree on a project problem and develop an investigation plan by Week 5

4. Seminars – worth 20%
   — Select the topic you are interested in by end of Week 4
     • Note limit of 3 people per topic; nominations accepted in first-come, first-served order
   — Prepare by reading suggested papers, and beyond, as you see fit

# Student Support - I Need Help With…

| | | |
|---|---|---|
| **Uni and Life in Australia**<br>Stress, Financial, Visas, Accommodation & More | **Student Support** | student.unsw.edu.au/**advisors** |
| **Reporting Sexual Assault/Harassment** | **Equity Diversity and Inclusion (EDI)** | edi.unsw.edu.au/**sexual-misconduct** |
| **Educational Adjustments**<br>To Manage my Studies and Disability / Health Condition | **Equitable Learning Services (ELS)** | student.unsw.edu.au/**els** |
| **Academic and Study Skills** | **Academic Skills** | student.unsw.edu.au/**skills** |
| **Special Consideration**<br>Because Life Impacts our Studies and Exams | **Special Consideration** | student.unsw.edu.au/**special-consideration** |

**My Feelings and Mental Health**
Managing Low Mood, Unusual Feelings & Depression

| | | | |
|---|---|---|---|
| **Mental Health Connect** | student.unsw.edu.au/**counselling**<br>Telehealth | **In Australia Call Afterhours**<br>**UNSW Mental Health Support Line** | 1300 787 026<br>5pm-9am |
| **Mind HUB** | student.unsw.edu.au/**mind-hub**<br>Online Self-Help Resources | **Outside Australia Afterhours**<br>**24-hour Medibank Hotline** | +61 (2) 8905 0307 |

# PP4FPGAs
# Ch 1. Introduction

# Parallel Programming for FPGAs

- Verilog and VHDL are primarily used to specify designs at the register transfer level (RTL)
- High-level synthesis (HLS) enables a designer to focus on larger architectural questions than individual registers and cycle-to-cycle operations
  - A designer captures behaviour in a program that does not include specific registers or cycles and an HLS tool creates the detailed RTL micro-architecture
- The text covers a good deal of the uses of HLS and how to optimize code in order to achieve optimal speedup
- We'll cover the first 5 chapters of the book in lectures and labs
  - Hope that this will inspire you to work through other chapters when time permits
- The text is *not* primarily about HLS compiler algorithms (scheduling, resource allocation, binding, etc.)

# HLS tools

HLS does several things automatically that an RTL designer does manually:

- HLS analyses and exploits the concurrency in an algorithm
- HLS inserts registers as necessary to limit critical paths and achieve a desired clock frequency
- HLS generates control logic that directs the data path
- HLS implements interfaces to connect to the rest of the system
- HLS maps data onto storage elements to balance resource usage and bandwidth
- HLS maps computations onto logic elements performing user specified and automatic optimizations to achieve the most efficient implementation

# Vivado HLS

Mature HLS tools include:
- Xilinx Vivado/Vitis HLS, LegUp and Mentor Catapult HLS

We'll study Vivado HLS as a representative example
- Other tools may use different syntax or semantics, but the principles remain the same

Vivado HLS requires the following inputs:
- A function specified in C/C++ or SystemC
- A testbench that calls the function and verifies its correctness
- A target FPGA device
- The desired clock period
- Directives guiding the implementation process

# HLS tools

❑ HLS can't handle arbitrary code.

❑ Many SW concepts are difficult to implement in HW:
  - Dynamic memory allocation, recursion, standard libraries, system calls

❑ On the other hand, HLS can deal with a variety of interfaces:
  - DMA, streaming, on-chip memories

❑ And perform advanced optimizations to create efficient implementations
  - Pipelining, memory partitioning, bitwidth optimization

Thus, HLS is both more restrictive but also (through compiler directives) enhances the input language for the purpose of hardware design

# Assumptions on Vivado HLS input functions

- No dynamic memory allocation (no operators like malloc(), free(), new(), and delete())
- Limited use of pointers-to-pointers (e.g., may not appear at the interface)
- System calls are not supported (e.g., abort(), exit(), printf(), etc.)
  - They can be used in the code, e.g., in the testbench, but they are ignored (removed) during synthesis.
- Limited use of other standard libraries (e.g., common math.h functions are supported, but uncommon ones are not)
- Limited use of function pointers and virtual functions in C++ classes (function calls must be compile-time determined by the compiler).
- No recursive function calls.
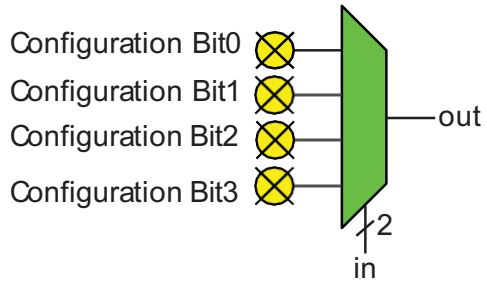- The interface must be precisely defined.

# Vivado HLS outputs

- Synthesizable Verilog and VHDL
- RTL simulations based on the design testbench
- Static analysis of performance and resource usage
- Metadata at the boundaries of a design to facilitate integration into a system

Following RTL output, the design can be synthesized using vendor tools to produce a netlist of FPGA logical elements, which is mapped to the physical FPGA resources during the place and route process. The resulting FPGA configuration of logic elements, wire connections and on-chip memories is captured in a bitstream.

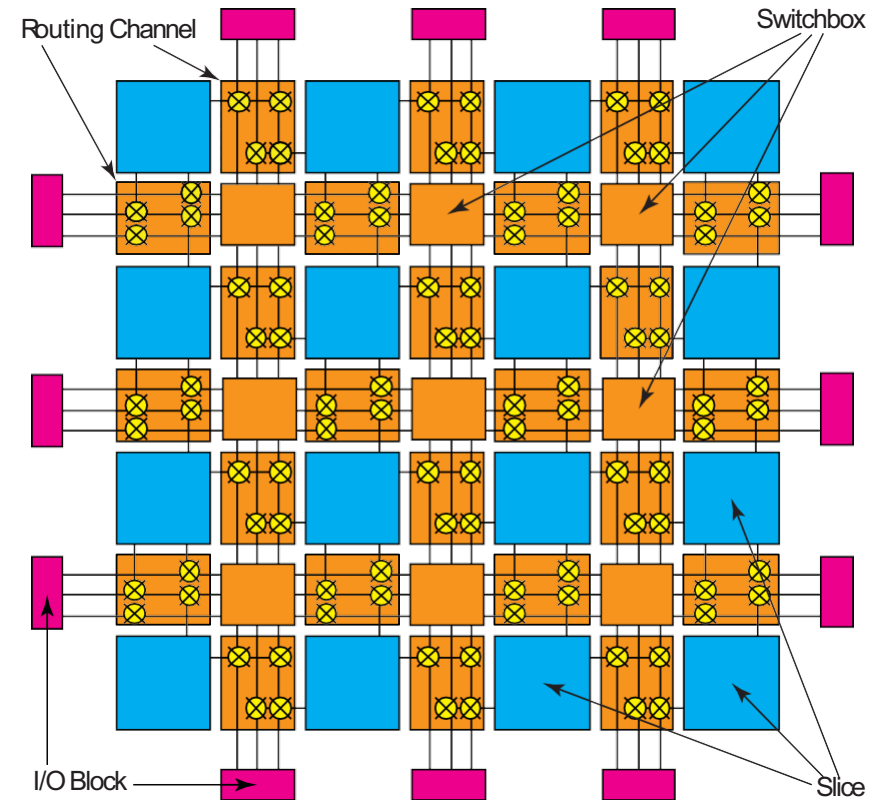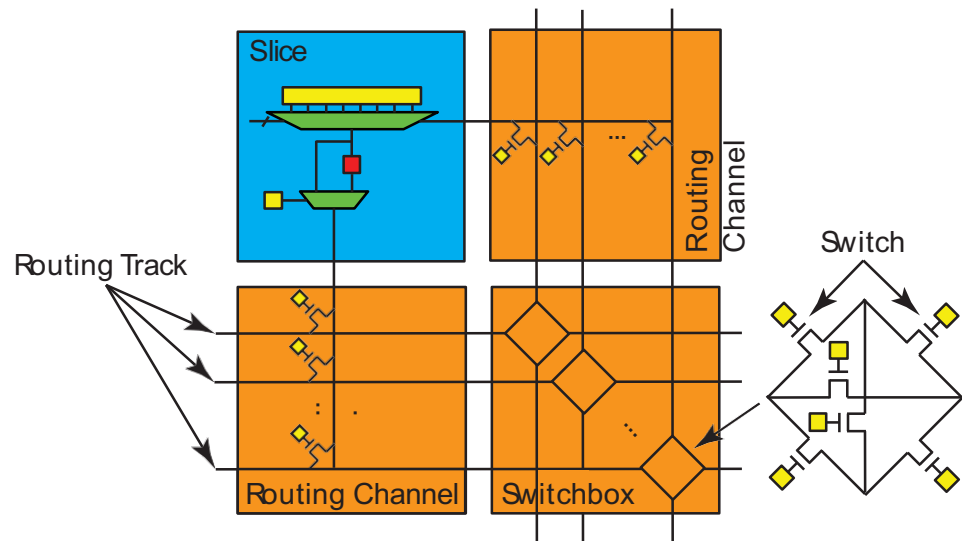# FPGA architecture refresher
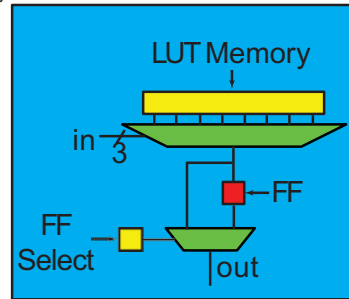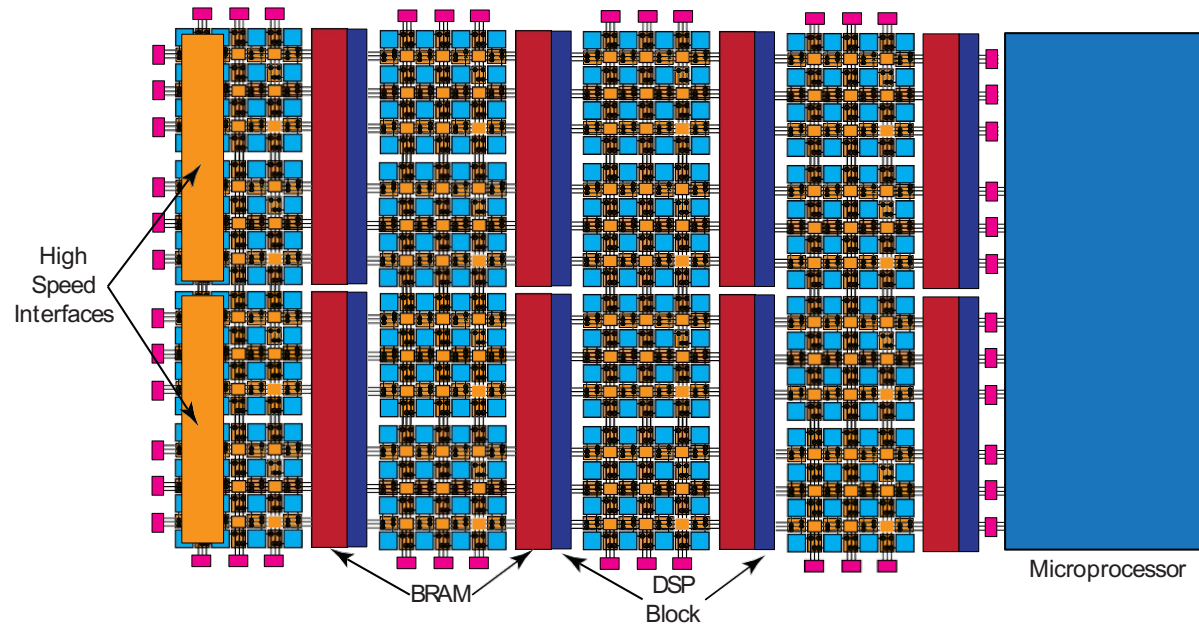


a) Lookup Table (LUT)

Configuration Bit0
Configuration Bit1
Configuration Bit2
Configuration Bit3

out

2

in

b)

| in[1] | in[0] | out |
|-------|-------|-----|
| 0 | 0 | 0 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |
| 1 | 1 | 1 |

out = in[1] & in[0]

c) Slice

LUT Memory

in 3

FF

FF Select

out

Routing Channel

Switchbox

Slice

Routing Channel

Routing Track

Switch

Routing Channel    Switchbox

I/O Block

Slice

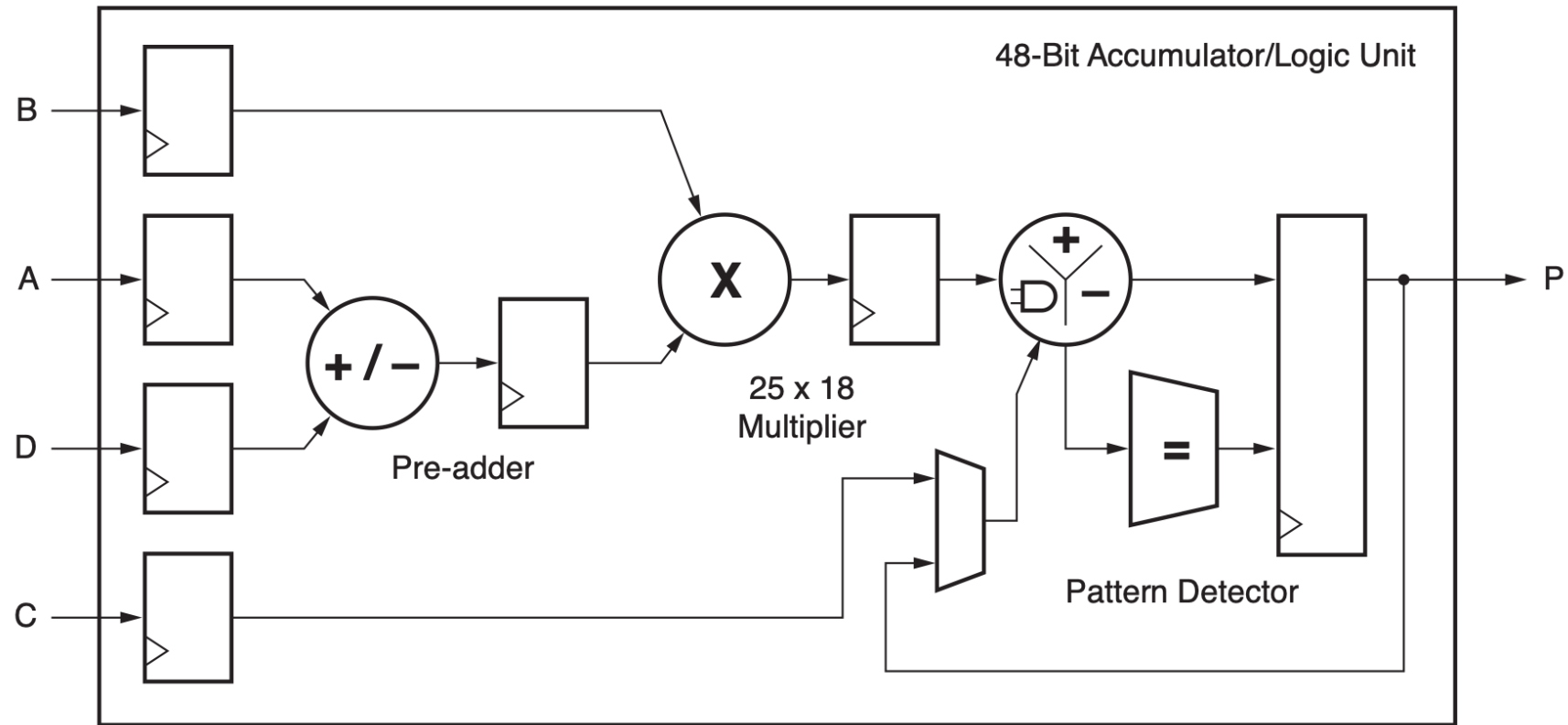# Modern FPGAs



High Speed Interfaces

BRAM

DSP Block

Microprocessor

- Numerous hardened resources for improved performance
  - DSP slices, block RAM, microprocessors, AXI interfaces, high-speed transceivers, clock managers etc.
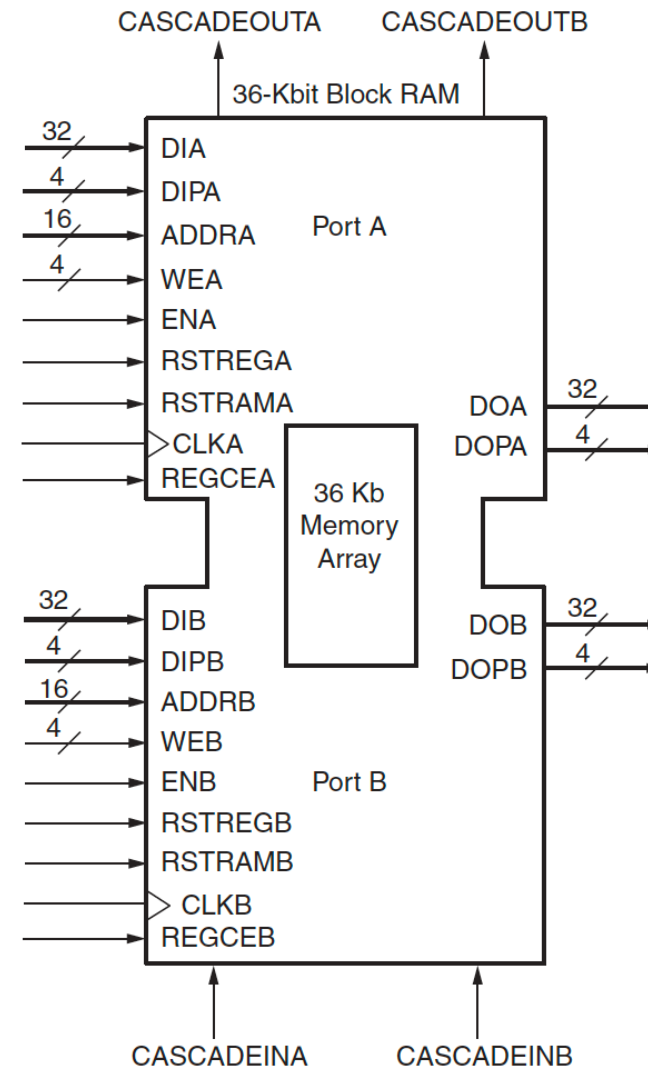
# DSP slices



UG479_c1_21_032111

# BRAM

Configurable RAM modules that support different memory layouts and interfaces

- E.g., Xilinx 36Kb BRAMs can be organized as 64K/32Kx1-, 16Kx2-, 8Kx4-, 4Kx9-, 2Kx18-, 1Kx36- or 512x72-bit storage and connected to local, on-chip buses or processor buses

- Configuration done via vendor tools; a major advantage of Vivado HLS is that the designer does not need to worry about these low-level details

- Used to transfer data between on-chip resources (e.g. fabric and processor) and store large datasets on chip

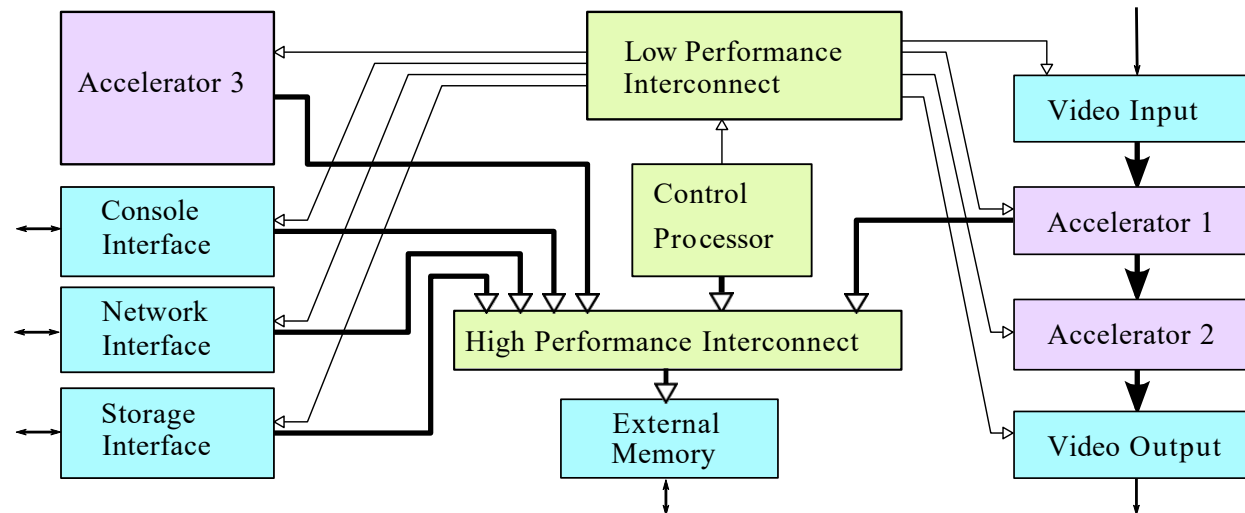- Only support one or two ports i.e. accesses per cycle

# Memory choices

The decision about where to place your application's data is crucial and one that we will revisit a number of times throughout this course

|  | External Memory | BRAM | FFs |
|---|---|---|---|
| count | 1-4 | thousands | millions |
| size | GBytes | KBytes | Bits |
| total size | GBytes | MBytes | 100s of KBytes |
| width | 8-64 | 1-16 | 1 |
| total bandwidth | GBytes/sec | TBytes/sec | 100s of TBytes/sec |

Vivado HLS provides options to allow the designer to specify exactly where and how to store data

# Example FPGA design for video processing



A block diagram showing a hypothetical embedded FPGA design, consisting of I/O interface cores (shown in blue), standard cores (shown in green), and application specific accelerator cores (shown in purple). Note that accelerator cores might have streaming interfaces (Accelerator 2), memory-mapped interfaces (Accelerator 3), or both (Accelerator 1).

# FPGA design process

FPGA designs are often composed of components or IP cores, structured like on the previous slide

- *I/O interface cores* at the periphery implement timing critical I/O functions or protocols e.g. memory controller, video interface, A/D converter
  - I/O interfaces are often highly customized for a particular FPGA architecture and hence are typically provided by the FPGA vendor as reference designs or off-the-shelf components
- Synchronous *standard cores*, such as processors, on-chip memories, and generic, fixed-function processing components, such as filters and codecs, are also often provided by an FPGA vendor, as these do not differentiate a product
- FPGA designs also typically contain customized, application-specific, synchronous *accelerator cores*, that are usually created by the system designer, since they contribute the "secret sauce" that differentiates a design from others

# Accelerators

- Ideally a designer can quickly and easily generate high-performance custom cores, perform a design space exploration of feasible designs, and integrate these into their system in a short timeframe.

- This course focuses on using HLS to design and implement high-performance custom cores quickly and efficiently.

- A *core-based design methodology* composes the custom cores with other provided cores using design integration tools.

- Alternatively, a *platform-based design methodology* makes use of standard design templates that combine a stable, verified composition of standard and I/O cores targeting a particular board. This enables a high-level programmer to integrate different algorithms or *roles* within the interface provided by a platform or *shell*. Accelerators designed this way can be ported between platforms whose shells have similar interfaces.

# Performance characterization

Best to use <u>run time in seconds</u> to compare performance

- A target clock period is provided as a constraint to Vivado HLS, and the tool reports the number of clock cycles the generated architecture needs to complete processing
    - The specified clock period impacts the feasible architectures for the design vis-à-vis which operations (and how many) can be performed within one clock cycle

Vivado HLS counts cycles by determining the maximum number of registers between any input and output
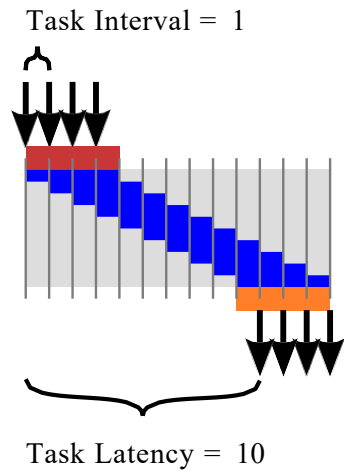
# Task interval & latency

The term *task* is used to mean a fundamental unit of behaviour, which corresponds to a C/C++ function invocation in Vivado HLS

- The *task latency* is the time from when a task starts until it finishes
- The *task interval* is the time between one task starting and the next starting – it is the difference between the start times of two consecutive tasks
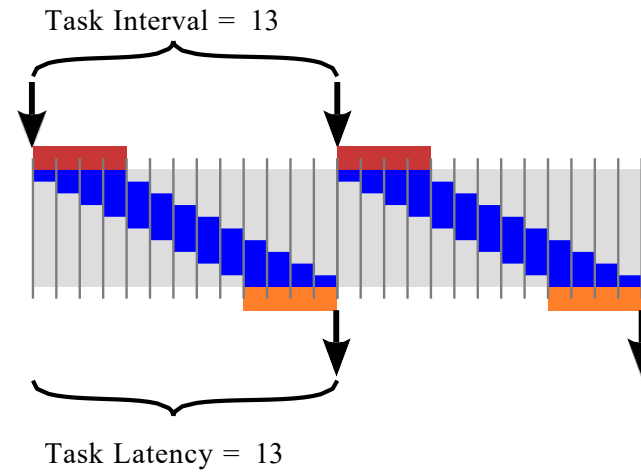
In many designs, data rate is a key design goal, and depends on the task latency and the size of the arguments to the function

While all task I/O and computation is bounded by the task latency, the time at which I/O occurs does not necessarily correspond with the start or end time of the task
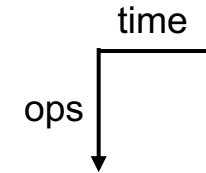
# Task interval & latency

Task Interval = 1

Task Interval = 13

time

ops

Task Latency = 10

Task Latency = 13

Pipelined design

Non-pipelined design

Similar to instruction
pipeline in a
microprocessor, but
customized

# A typical HLS function

Let's take a look at the implementation of a finite impulse response (FIR) filter

- An FIR filter performs a convolution on an input sequence with a fixed set of coefficients
  e.g. a moving average filter

An HLS tool will analyze the code and produce a functionally equivalent RTL circuit

- This is a complex process, which we won't get into the detail of, but think of it as a compiler, like *gcc*, that outputs an RTL description
- It's not necessary to understand exactly how the compiler works in order for it to allow the programmer/designer to work at a higher level of abstraction
- Yet knowledge of how the compiler works can enable you to write more efficient code. This is particularly true for HLS because of the synthesis options e.g. memory layout, pipelining etc. that are not typically obvious to someone that only understands the software flow

```
#define NUM_TAPS 4

void fir(int input, int *output, int taps[NUM_TAPS]) {

  static int delay_line[NUM_TAPS] = {};

  int result = 0;
  for (int i = NUM_TAPS-1; i > 0; i--) {
   delay_line[i] = delay_line[i-1];
  }
  delay_line[0] = input;

  for (int i = 0; i < NUM_TAPS; i++) {
   result += delay_line[i] * taps[i];
  }

  *output = result;
}
```
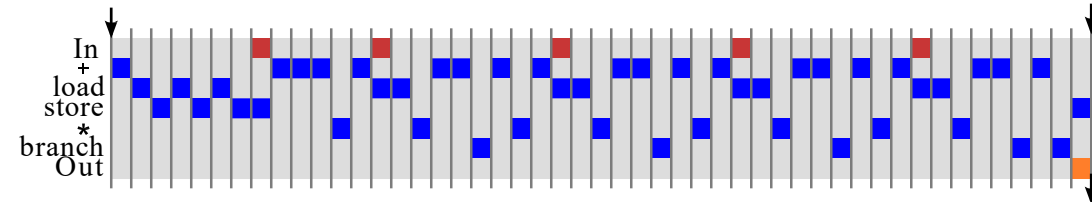
# Area/throughput tradeoffs

There are numerous possible circuits generated by an HLS tool

One possible circuit executes the code sequentially, like a simple RISC processor (see next slide)

# Compiled FIR function

```
fir:
.frame r1,0,r15 #  vars= 0, regs= 0, args= 0
.mask 0x00000000
addik r3,r0,delay line.1450
lwi r4,r3,8 #  Unrolled loop to shift the delay line
swi r4,r3,12
lwi  r4,r3,4
swi r4,r3,8
lwi  r4,r3,0
swi r4,r3,4
swi r5,r3,0 #  Store the new input sample into the delay line
addik r5,r0,4 #  Initialize the loop counter
addk r8,r0,r0 #   Initialize accumulator to zero
addk r4,r8,r0 #  Initialize index expression to zero
$L2:
muli r3,r4,4 #  Compute a byte offset into the delay line array
addik r9,r3,delay line.1450
lw r3,r3,r7 #  Load filter tap
lwi r9,r9,0 #  Load value from delay line
mul r3,r3,r9 #  Filter Multiply
addk r8,r8,r3 #  Filter Accumulate
addik r5,r5,−1 #  update the loop counter
bneid r5,$L2
addik r4,r4,1 #  branch delay slot, update index expression

rtsd r15, 8
swi r8,r6,0 #  branch delay slot, store the output
.end fir
```



Task takes 49 CC to compute one output sample if one instruction/cycle is executed sequentially

# Area/throughput tradeoffs

There are numerous possible circuits generated by an HLS tool

One possible circuit executes the code sequentially, like a simple RISC processor (see previous slide)

One characteristic of HLS is that architectural tradeoffs can be made without needing to fit to the constraints of an instruction set architecture

- HLS designs may generate architectures that issue hundreds or thousands of RISC-equivalent instructions per clock with pipelines that are hundreds of cycles deep

# Sequential architecture

By default, Vivado HLS will generate an optimized, but largely sequential architecture

- Loops and branches are transformed into control logic that enables the registers, functional units and the rest of the datapath
- Conceptually similar to the execution of a RISC processor except that the program to be executed is converted to an FSM in the generated RTL rather than being fetched from program memory
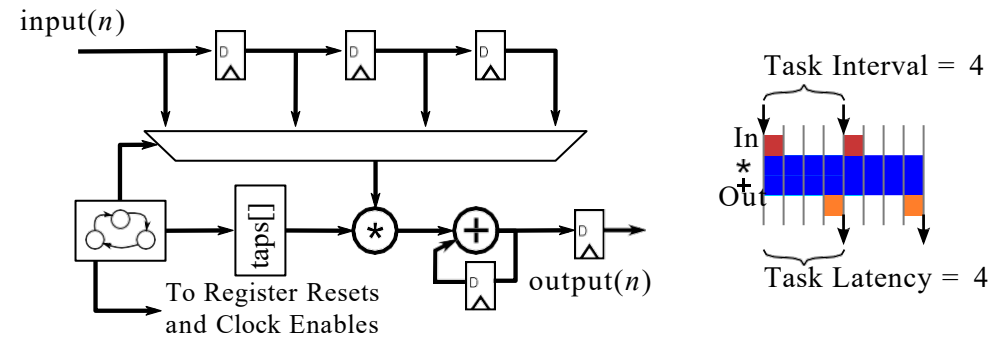
A sequential architecture tends to limit the number of functional units in a design with a focus on resource sharing over massive parallelism

- Complexity of the control logic can hamper analysis; the behaviour of the control logic may also be data dependent

# One tap per clock FIR function

The Vivado HLS tool can be directed to generate a pipeline by placing a `#pragma HLS pipeline` directive into the body of a function
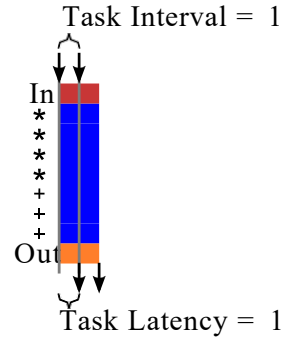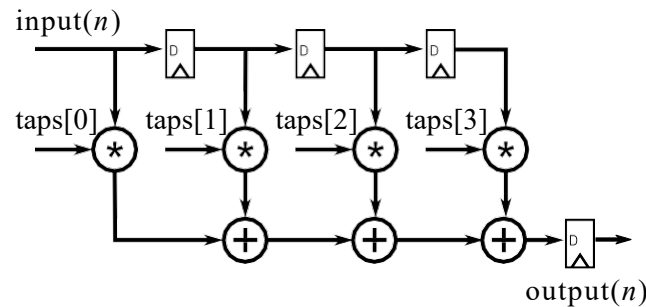
- Includes a parameter to specify the *initiation interval* of the pipeline



A "one tap per clock" architecture

- One mult, one add – results in a task latency and task interval of 4

# One sample per clock FIR function



input($n$)

taps[0]   taps[1]   taps[2]   taps[3]

output($n$)

Task Interval = 1

In

Out

Task Latency = 1

A "one sample per clock" architecture

- 4 mults, 3 adds with task latency and interval of 1

Other implementations, such as "two taps per clock" or "two samples per clock", which may be needed to meet higher throughput requirements, are also possible

# Controlling the output generated by HLS

In practice, complex designs often include complicated tradeoffs between sequential architectures and parallel architectures, in order to achieve the best overall design

- In Vivado HLS, these tradeoffs are largely controlled by the user, through various tool options and code annotations, such as `#pragma` directives

# Restrictions on processing rate

The task interval (and hence throughput) is fundamentally limited by *recurrences* (feedback loops) and *resource limits*

A recurrence is any case where a computation by a component depends upon a previous computation by the same component

- Examples include the static variables in the FIR code and the accumulator in the one tap per clock design
- Recurrences fundamentally limit throughput even when pipelining
- Analyzing recurrences and generating hardware that is guaranteed correct are key functions of HLS tools
- Similarly, understanding algorithms and selecting those without tight recurrences is an important onus on the designer using HLS

# Resource limitations

Another key factor limiting the processing rate is resource limitations

One form of resource limitation is associated with the wires at the boundary of a design, since a synchronous circuit can only capture or transmit one bit per wire per clock cycle

Another form of resource limitation arises from memories since the number of accesses per cycle is usually limited

The designer might also create a limitation by restricting the number of operators that can be instantiated during synthesis

Finally, the coding style used e.g. choice of `#pragma` directives can limit the range of architectures that can be generated from the code

# The importance of code restructuring

Using HLS is not just a matter of adding `#pragma` directives to your software code

Kastner stresses the need to have a good understanding of the application at hand so that optimizations that require rewriting of the code (code restructuring) can be taken advantage of

- Standard, off-the-shelf code typically yields very poor quality of results that are <span style="color:red">orders of magnitude slower</span> than CPU designs, thus it is also important to know how to write code that the HLS tool will synthesize in an optimal manner
- Restructuring code is an essential step to generate an efficient FPGA design
- Writing restructured code requires significant hardware design expertise and domain-specific knowledge
- Many of the examples in the text show how to restructure the code for more efficient hardware design

# Chapter outline

| Chapter | FIR 2 | CORDIC 3 | DFT 4 | FFT 5 | SPMV 6 | MatMul 7 | Histogram 8 | Video 9 | Sorting 10 | Huffman 11 |
|---|---|---|---|---|---|---|---|---|---|---|
| Loop Unrolling | X | | X | X | X | | X | | X | |
| Loop Pipelining | X | | X | X | X | | X | X | X | X |
| Bitwidth Optimization | X | X | | | | | | | | X |
| Function Inlining | X | | | | | | | | | X |
| Hierarchy | X | | | X | | | X | X | X | X |
| Array Optimizations | | | X | X | X | X | X | X | X | X |
| Task Pipelining | | | | X | | | X | X | X | X |
| Testbench | | | | | X | X | | | X | X |
| Co-simulation | | | | | X | | | | | |
| Streaming | | | | | | X | | X | X | |
| Interfacing | | | | | | | | X | | |

We'll be studying the first 5 chapters of the text in class. You will be introduced to the parallelization techniques listed above the red dashed line.