



# OPENCL ON INTEL FPGA

Thinks Efficiency !!!

Marc Gaucheron, Intel PSG



## Agenda

### What Is FPGA ?

- Anatomy of FPGA – 15mn
- FPGA development flow – 15mn
- High Level Synthesis Tools – 15mn
- OpenCL for FPGA – 2h00
  - OpenCL SDK & Execution Model
  - Optimizing Kernel For FPGA
- Conclusion – 15mn



# ANATOMY OF FPGA

JDEV2017

## Acronyms

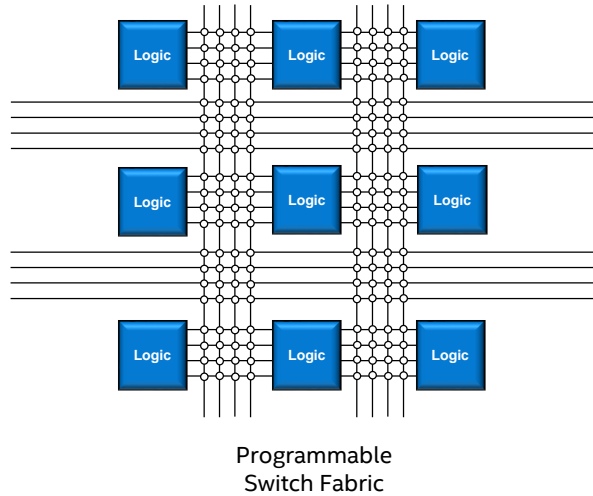
LUT- LookUp Table

LE – Logic Element

ALM – Adaptive Logic Module

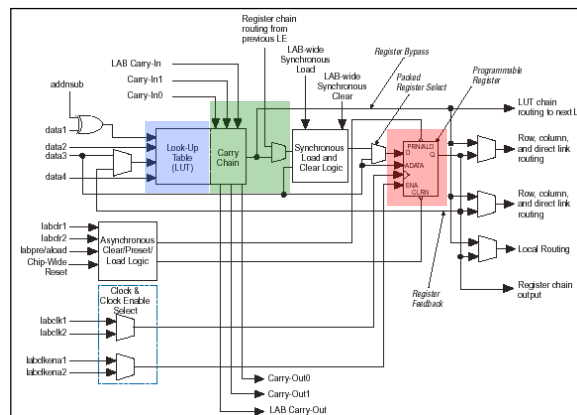
LAB – Logic Array Block

# Field Programmable Gate Array (FPGA)



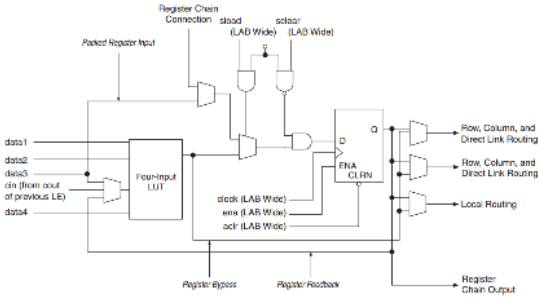
## FPGA Logic blocks

FPGA logic is made up of Logic Elements (LEs) or Adaptive Logic Modules (ALMs)

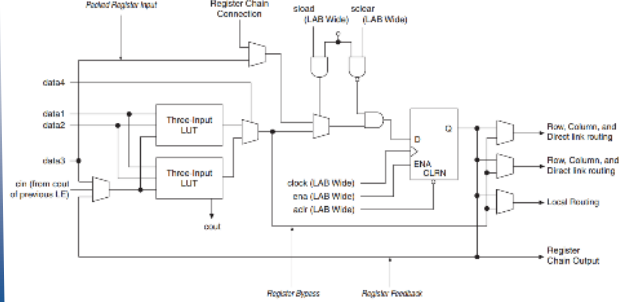


# Logic Element (LE)

**Normal Mode** – Suited for general logic applications and combinational functions



**Arithmetic Mode** – Implementing adders, counters, accumulators, & comparators



The Quartus Prime software automatically chooses the appropriate operation mode for individual functions for optimal performance

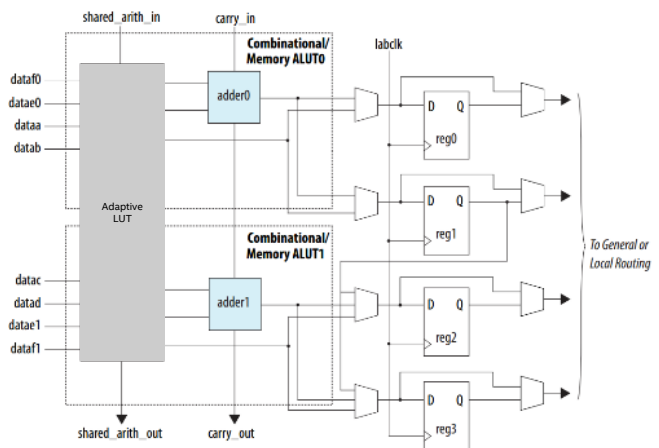
# Adaptive Logic Module (ALM)

One ALM contains 4 programmable registers, each with the following ports:

- Data
- Clock
- Synchronous & asynchronous clear
- Synchronous load

Backward compatible with 4-input LUT architectures

Quartus Prime automatically configures ALMs for optimized performance

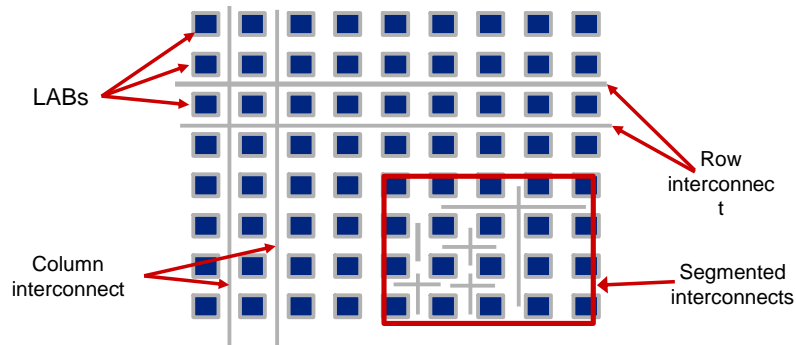


## Logic Array Block (LAB)

LABs are group of LEs or ALMs

Row and column programmable interconnect

Interconnect may span all or part of the array

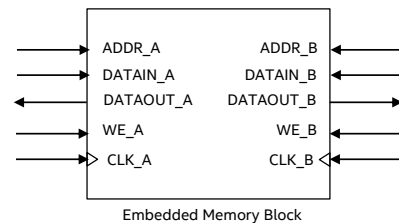


9

## FPGA Embedded Memory

### Memory blocks

- Create on-board memory structures to support design
  - Single/dual-port RAM
  - ROM
  - Shift registers or FIFO buffers
- Initialize RAM or ROM contents on power-on
  - M9K, M10K, M20K
- Memory LABs (MLABs)
- eSRAM



Embedded Memory Block



10

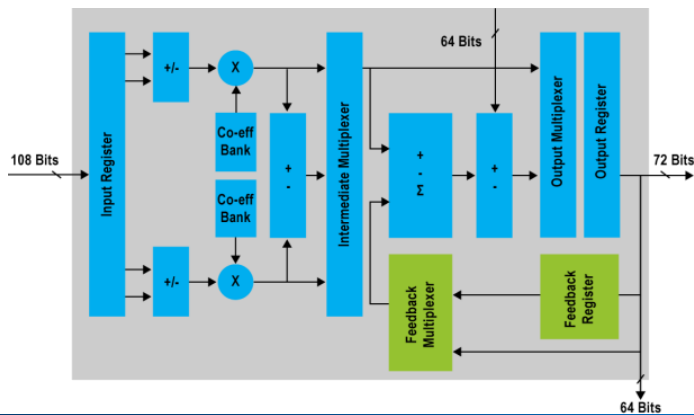
## Arria 10 Embedded Memories

| Feature                                |                           | MLABs   | M20K  |
|--|---------------------------|---------|---|
| Maximum Performance                    |                           | 700 MHz | 730 MHz   |
| Total RAM Bits per Block               |                           | 640     | 20,480  |
| Total M20K Memory Bits (Mb) per Device |                           | 2-13    | 13-54   |
| Port Width Configurations              |                           |         | 16K x 1, 8K x 2<br>4K x 4, 5<br>2K x 8, 10<br>1K x 16, 20<br>512 x 32, 40 |
|  |                           | 32 x 16 |   |
|  |                           | 32 x 18 |   |
|  |                           | 32 x 20 |   |
| Parity                                 |                           | ✓       | ✓   |
| Byte Enable                            |                           | ✓       | ✓   |
| Packed Mode                            |                           | ✓       | ✓   |
| Address Clock Enable                   |                           |         | ✓   |
| Mixed Clock                            |                           | ✓       | ✓   |
| Mixed Width (for Dual Port modes)      |                           | ✓       | ✓   |
| ECC Support                            |                           | ✓       | Hard  |
| Memory Modes                           | Single Port               | ✓       | ✓   |
|  | Simple & True Dual-Port   | ✓       | ✓   |
|  | Shift Register, ROM, FIFO | ✓       | ✓   |

## DSP Block

Useful for DSP functions

High-performance multiply/add/accumulate operations



## Arria 10 DSP block

### Multiplier Modes for Flexibility

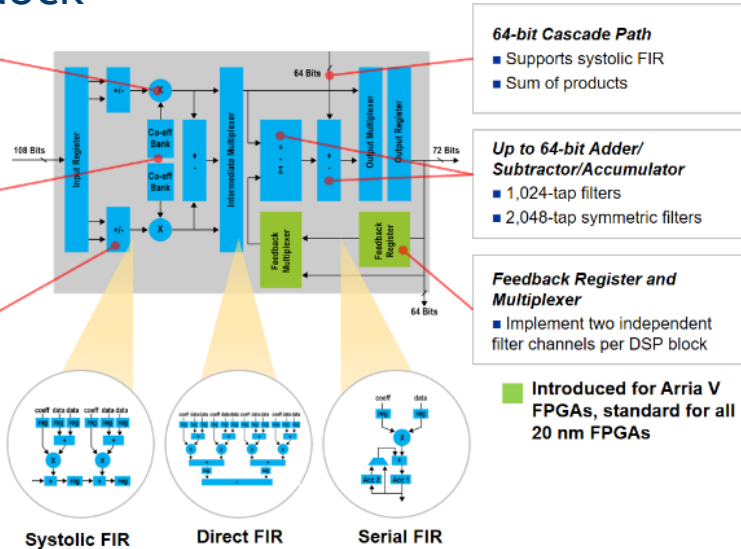
- Two 18x19 multipliers, or
- One 27x27 multiplier per block
- 36x36, 54x54 modes using multiple DSP blocks

### Integrated Coefficient Registers

- Save memory and routing resources
- Built-in timing closure

### Hard Pre-Adders

- Reduce multiplier usage
- Save routing resources



Up to 1.5 TFlops



13

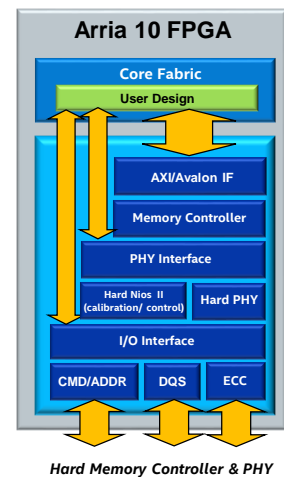
## Arria 10 Hard Memory Controller & PHY

### Hard memory controller

- Saves logic and memory resources
  - 5K LEs and 29 M20K blocks per x72 DDR3 IF
- Up to x144 support
- Up to 4x72 DDR3 interfaces in a single device

### Hard memory controller supports

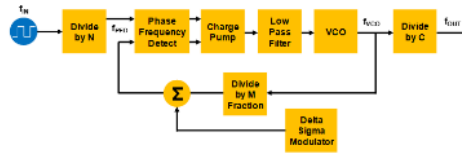
- DDR4, DDR3, LPDDR3
- DDR4 up to 2400 Mbps



14

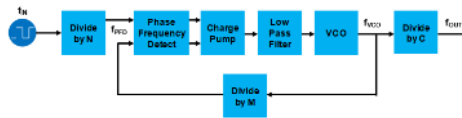
# Core PLLs

| Fractional PLLs           |   |
|---------------------------|---|
| Feature                   | Description   |
| # Available               | 1 for every 3 transceivers on device  |
| Location on die           | In the core, adjacent to transceivers   |
| Operating modes           | Fractional-synthesis Integer (M/N where M,N=Integer)  |
| FPGA clock network access | Transceiver reference clock<br>GCLK (global clock)<br>PCLK (periphery clock)<br>RCLK (regional clock) |



*fPLLs Reduce Cost, Power and Circuit Board Space*

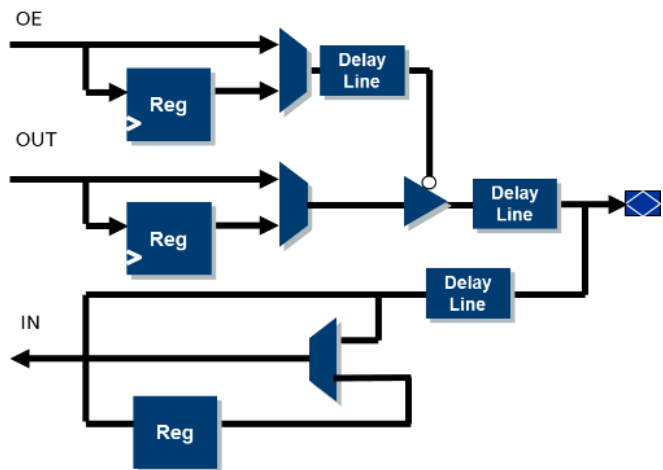
| IO PLLs                   |  |
|---------------------------|--|
| Feature                   | Description  |
| # Available               | 1 for every IO bank (48 GPIOs)   |
| Location on die           | In the core, adjacent to IO banks  |
| Operating modes           | Integer (M/N where M,N=Integer)  |
| FPGA clock network access | External Memory Interface<br>LVDS SerDes Interface<br>GCLK (global clock)<br>PCLK (periphery clock)<br>RCLK (regional clock) |



*IO PLLs Enable High-Bandwidth IO Interfaces*

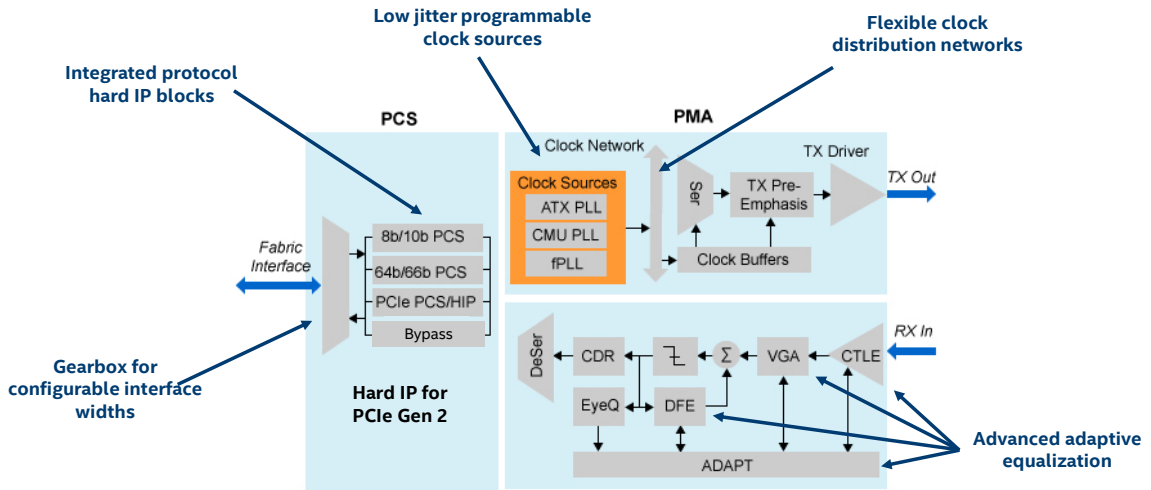
# FPGA IO Elements

- Input/output/bidirectional
- Multiple I/O standards
- Differential signaling
- Current drive strength
- Slew rate
- On-chip termination/pull-ups
- Open drain/tri-state



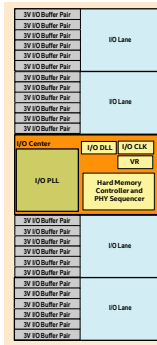


# High Speed IO Transceivers

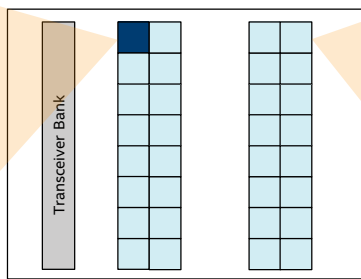
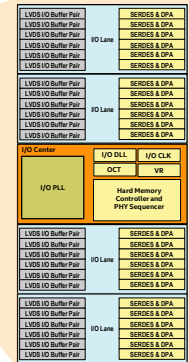


# Cyclone® 10 GX: GPIO Blocks

## 3V I/O Block Structure\*



## LVDS I/O Block Structure

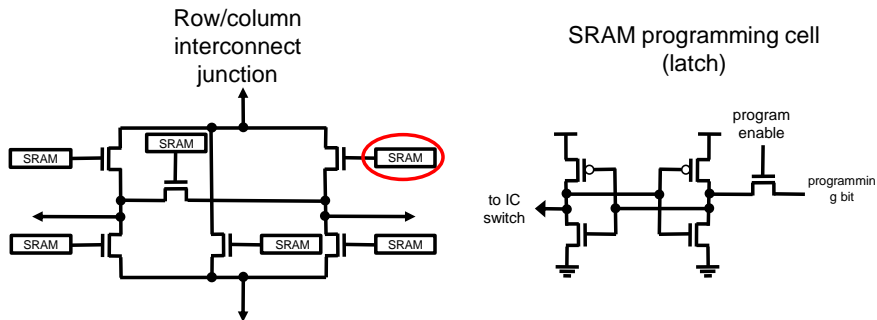


\* Only a single 3V I/O Block in Cyclone 10 GX

## FPGA Programming

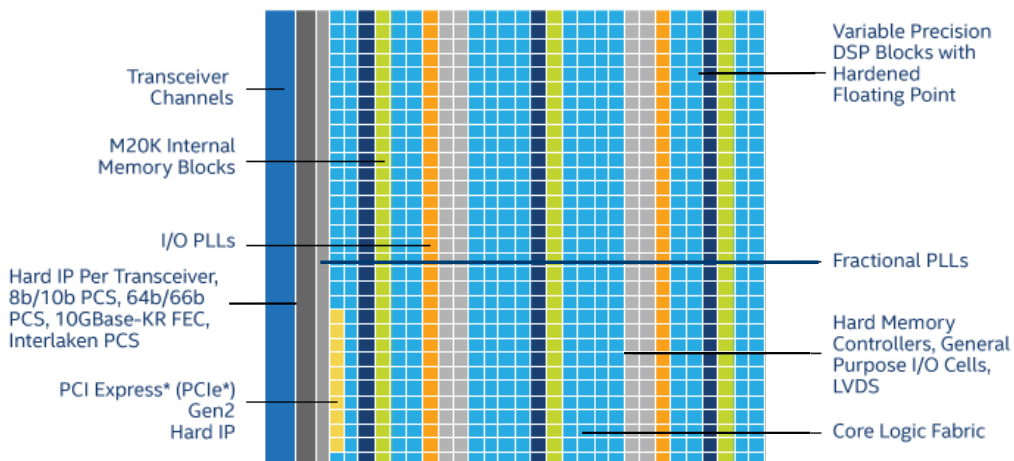
FPGAs use SRAM cell technology to program interconnect and LUT function levels

*Volatile! Must be programmed at power-on!*



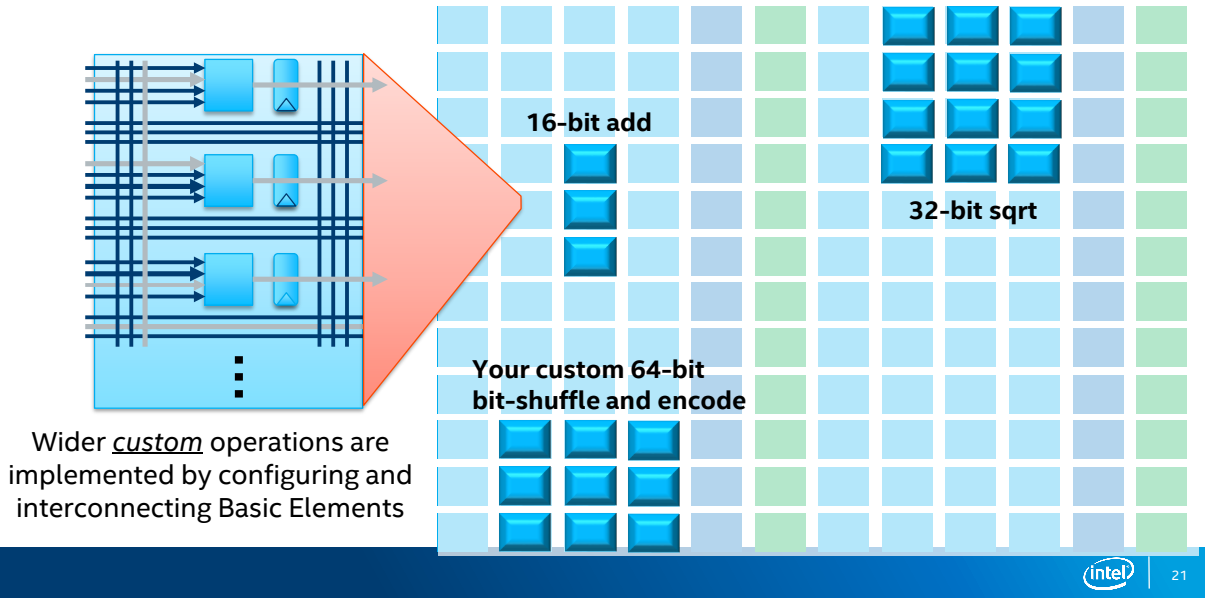
19

## FPGA Floorplan

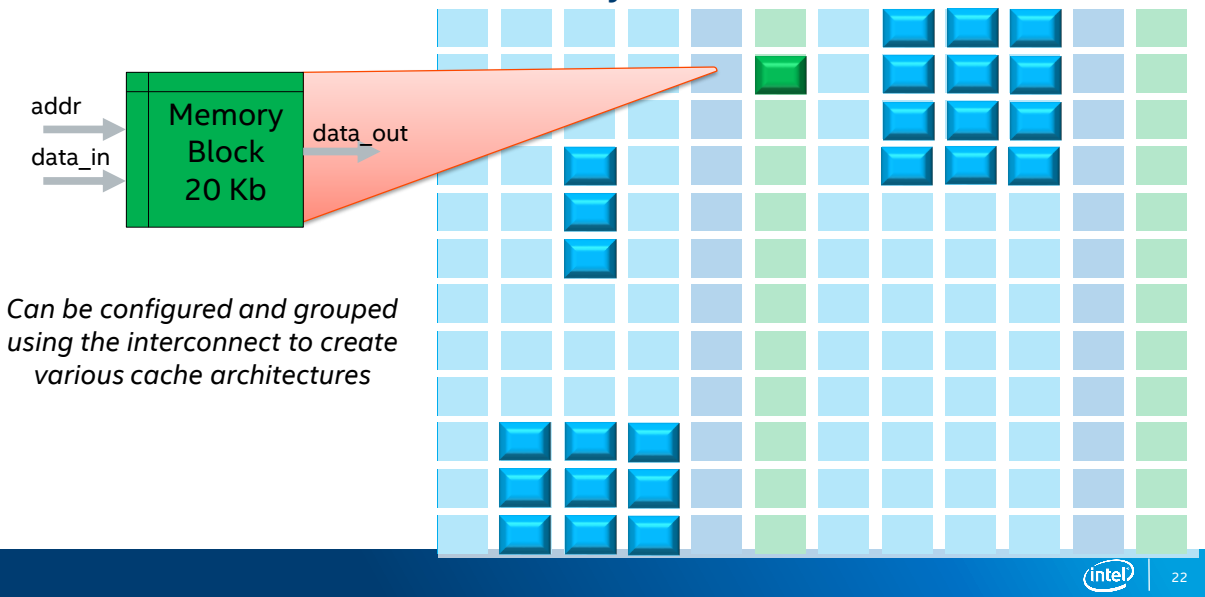


20

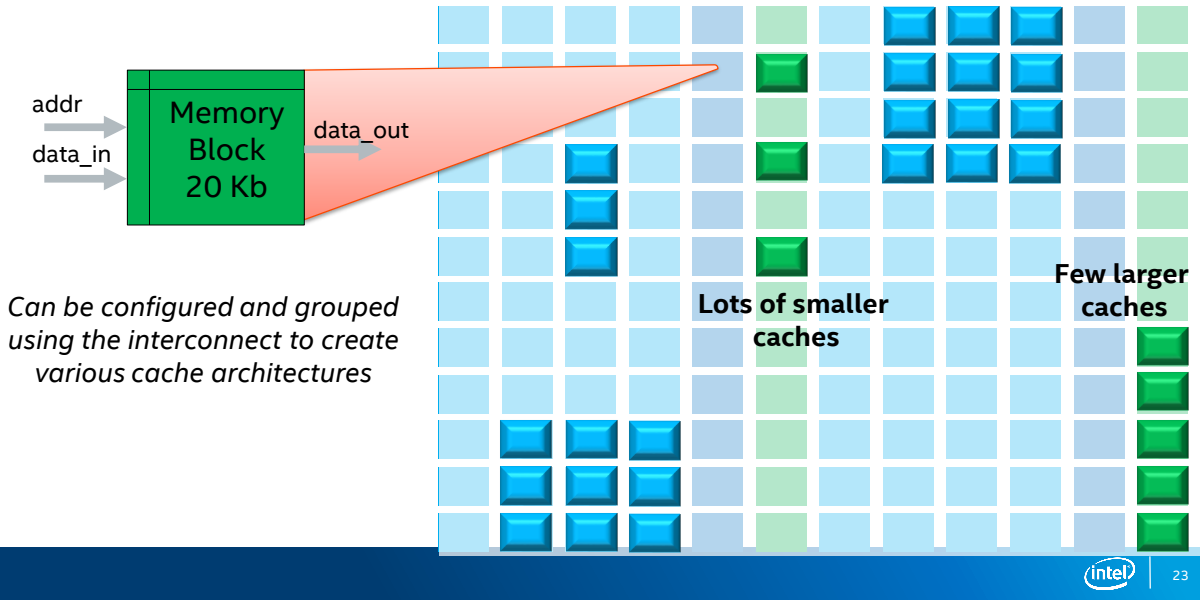
## FPGA Architecture: Custom Operations Using Basic Elements



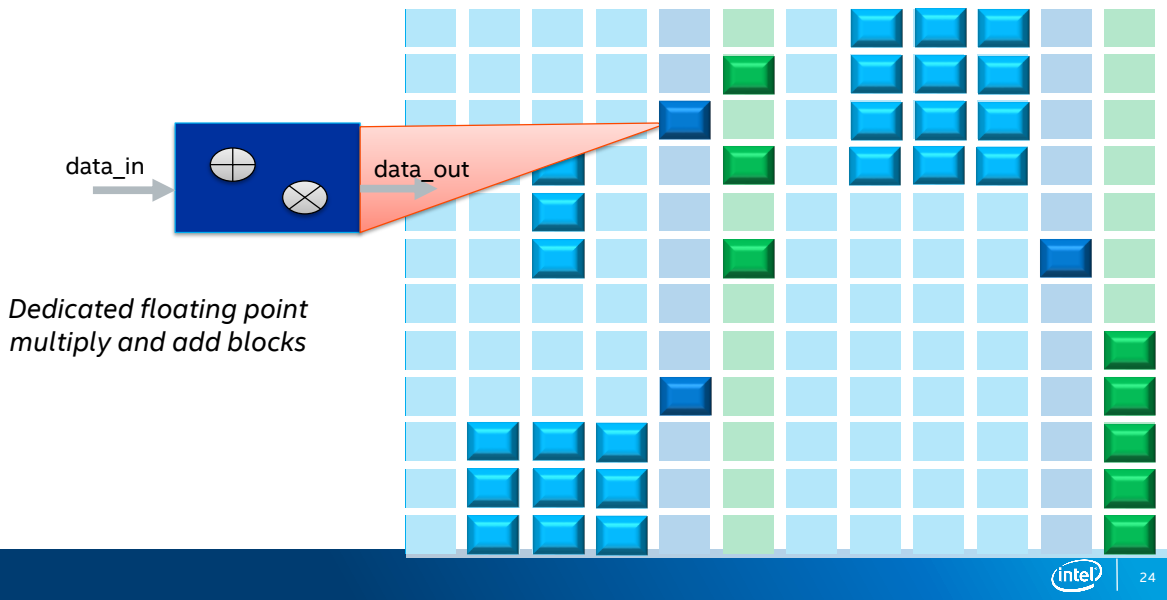
## FPGA Architecture: Memory Blocks



## FPGA Architecture: Memory Blocks

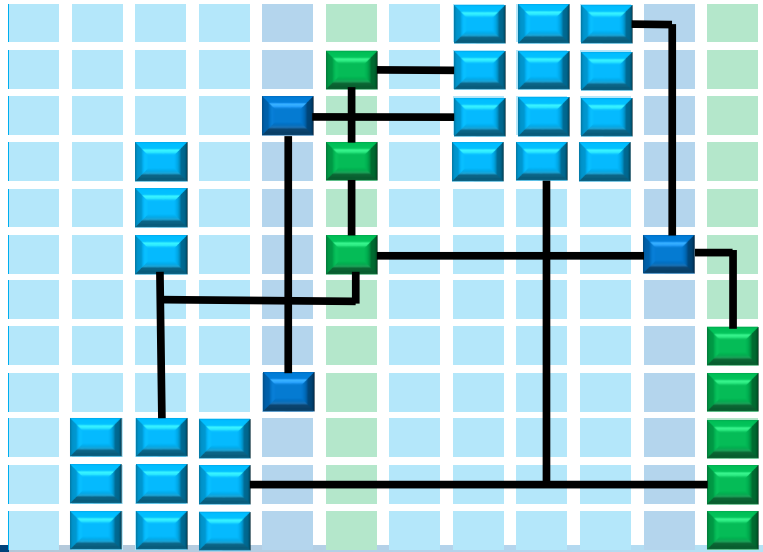


## FPGA Architecture: Floating Point Multiplier/Adder Blocks



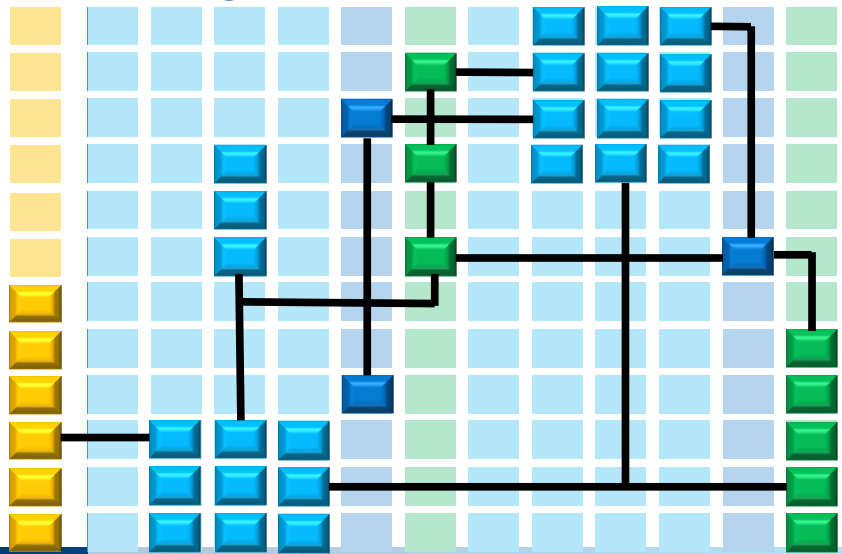
## FPGA Architecture: Configurable Routing

Blocks are connected into a **custom data-path** that matches your application.



## FPGA Architecture: Configurable IO

The **Custom data-path** can be connected directly to **custom or standard IO interfaces** for inline data processing





# FPGA DEVELOPMENT FLOW



## VHDL, VERILOG vs C

### VHDL

```

1 library ieee;
2 use ieee.std_logic_1164.all;
3
4 entity parity_generator is
5     generic( m           : integer);
6     port( input_stream : in std_logic_vector
7           (m-1 downto 0);
8           clk           : in std_logic;
9           parity        : out bit);
10 end parity_generator1;
11
12 architecture odd of parity_generator is
13 begin
14     P1: process
15         variable odd : bit;
16     begin
17         wait until clk'event and clk = '1';
18         odd := '0';
19         for i in 0 to m-1 loop
20             odd := odd xor input_stream(i);
21         end loop;
22         parity <= odd;
23     end process;
24 end odd;

```

### Verilog

```

1 module parity_generator (
2     input_stream,
3     parity,
4     clk
5 );
6
7     parameter M = 8;
8     input logic [M-1:0] input_stream;
9     input logic clk;
10    output logic parity;
11
12    always@(posedge clk) begin
13        parity <= ^input_stream;
14    end
15
16 endmodule;

```

### C++

```

1 #include "HLS/ac_int.h"
2 #define M=8
3
4 bool parity(
5     ac_int<M> input_stream
6 ) {
7     bool parity = true;
8     for(int i=0; i<M; ++i)
9         parity ^= ap_int[i];
10    return parity;
11 }

```

# VHDL, VERILOG vs C

VHDL

```

1 entity ODD_PARITY_TB is
2 end;
3
4 library ieee;
5 use ieee_std_logic_1164.all;
6 use WORK_anu.all;
7 use textio.all;
8
9 architecture OP_TB_ARCH of ODD_PARITY_TB is
10 component Parity_Generator1
11 port (
12     input_stream : in input;
13     clk          : in std_logic;
14     parity       : out bit );
15 end component;
16
17 signal input_stream : input;
18 signal clk          : std_logic;
19 signal parity       : bit;
20 variable dbg       : line;
21 begin
22     U1: Parity_Generator1
23     port map (
24         input_stream => input_stream,
25         clk          => clk,
26         parity       => parity
27     );
28     clk_process : process (clk)
29     begin
30         if clk <= '0' then
31             clk <= '0' after 1 ns;
32         else
33             clk <= not clk after 1 ns;
34         end if;
35     end process;
36 end OP_TB_ARCH;
    
```

Verilog

```

1 module parity_tb;
2     reg clk, input_stream, parity;
3
4     parity_generator dut (
5         .clk(clk),
6         .input_stream(input_stream),
7         .parity(parity)
8     );
9
10    initial
11    clk = 0;
12    always
13    #5 clk = !clk;
14
15    initial begin
16        wait @(posedge clk);
17        input_stream = 8'b10100110;
18        wait @(posedge clk);
19        $display("%b => %b",
20            input_stream, parity);
21        input_stream = 8'b01111100;
22        wait @(posedge clk);
23        $display("%b => %b",
24            input_stream, parity);
25    end
26 endmodule
    
```

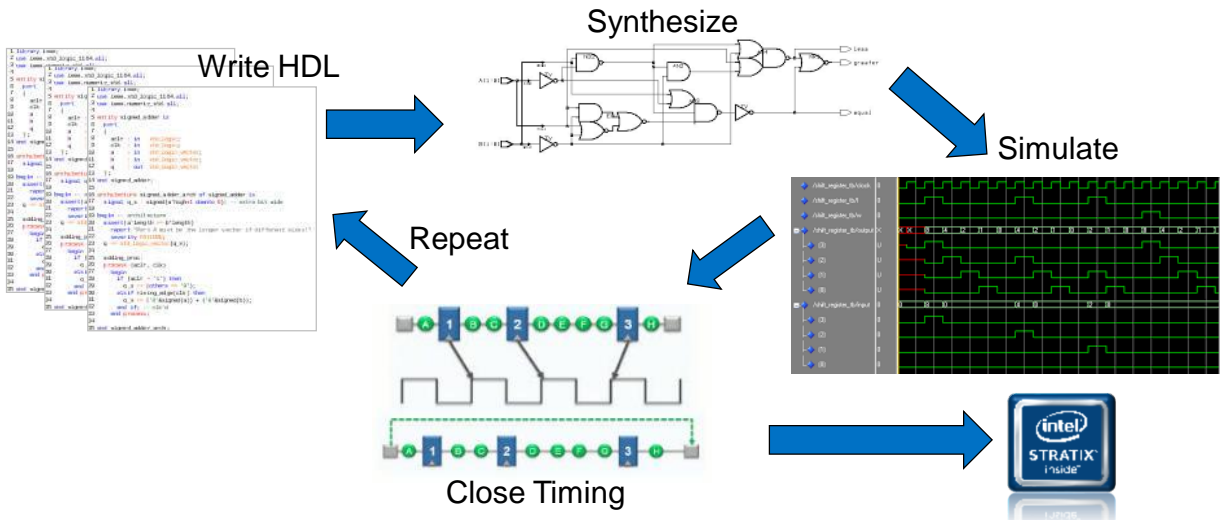
C++

```

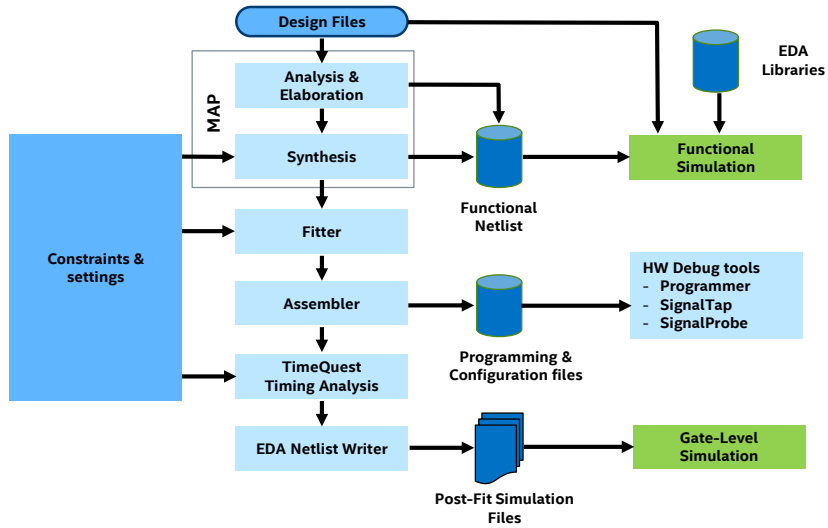
1 #include "stdio.h"
2
3 int main() {
4     int input;
5     input = 0b10100110;
6     printf("%h => %d\n", input,
7         parity_generator(input));
8     input = 0b10100110;
9     printf("%h => %d\n", input,
10        parity_generator(input));
11    return 0;
12 }
    
```



## Traditional FPGA Design Process



# Tools Flow Overview



# QSYS Prime

**System Contents**

| Name                | Description                     | Export           | Clock    | Base        |
|---------------------|---------------------------------|------------------|----------|-------------|
| clk                 | Clock Source                    | clk_clk_in       | exported |             |
| clk_in              | Clock Input                     | clk_clk_in_reset | clk      |             |
| clk_reset           | Reset Input                     |                  | clk      |             |
| clk_reset           | Clock Output                    |                  | clk      |             |
| clk_reset           | Reset Output                    |                  | clk      |             |
| pll                 | Altera PLL                      |                  | clk      |             |
| reset               | Reset Input                     |                  | clk      |             |
| out0                | Clock Output                    |                  | clk      |             |
| out1                | Clock Output                    |                  | clk      |             |
| locked              | Clockout                        |                  | clk      |             |
| external_connection | Conduit                         |                  |          |             |
| push_button_reg     | Avakon Push Button R/W Register |                  | sys_clk  | 0x0000_0000 |
| clock               | Clock Input                     |                  | sys_clk  |             |
| reset               | Reset Input                     |                  | sys_clk  |             |
| buttonreg           | Avakon Memory Mapped Slave      |                  | sys_clk  | 0x0000_0000 |
| avm_mem_master      | Avakon Memory Mapped Master     |                  | sys_clk  |             |
| avm_mem_slave       | Avakon Memory Mapped Slave      |                  | sys_clk  |             |
| clock               | Clock Input                     |                  | sys_clk  |             |
| reset               | Reset Input                     |                  | sys_clk  |             |
| led_driver          | LED Driver                      |                  | sys_clk  |             |
| clock               | Clock Input                     |                  | sys_clk  |             |
| reset               | Reset Input                     |                  | sys_clk  |             |
| streaming_sink      | Avakon Streaming Sink           |                  | sys_clk  |             |
| debug               | Avakon Memory Mapped Slave      |                  | sys_clk  | 0x0000_0000 |

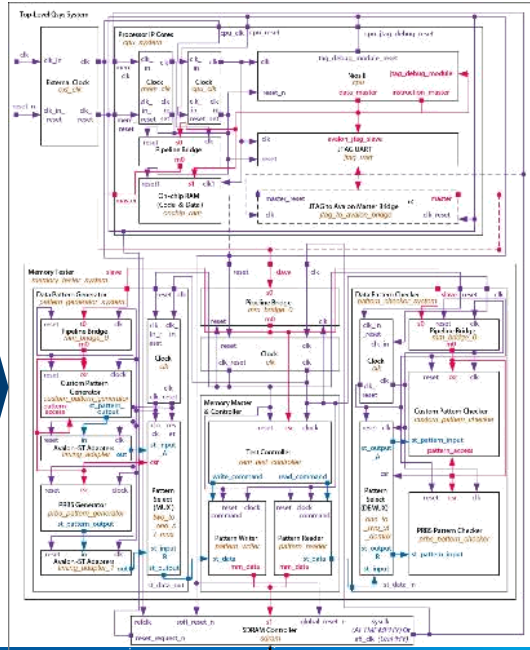
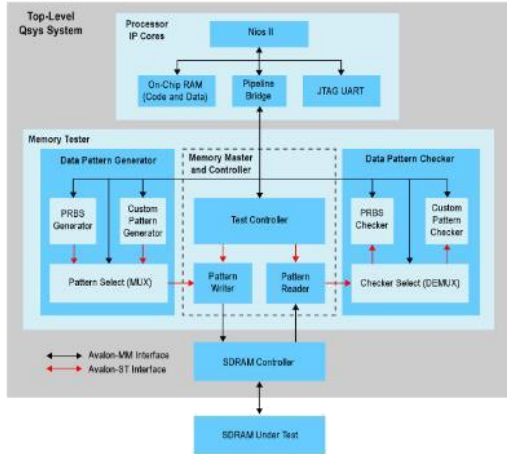
**Messages**

| Type    | Path                                  | Message   |
|---------|---------------------------------------|---|
| Warning | nm_transfer_system_controller         | Properties (Memory Device) have been set on interface <code>avm</code> - in composed mode these are ignored |
| Warning | sm_transfer_system_pll                | Unable to implement PLL - Actual settings differ from Requested settings                                    |
| Warning | nm_transfer_system_pll                | plllocked must be exported, or connected to a matching conduit.   |
| Warning | sm_transfer_system_dma_source_to_sram | Interrupt sender <code>dma_source_to_sram_irq</code> is not connected to an interrupt receiver              |

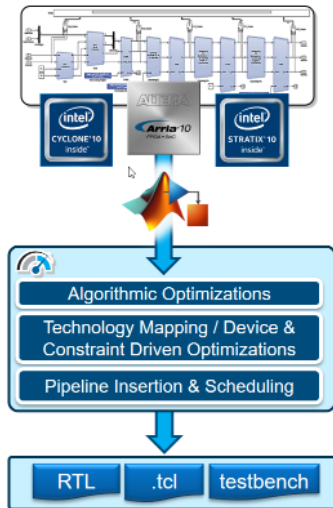


# QSYS Example Design

## Memory Tester

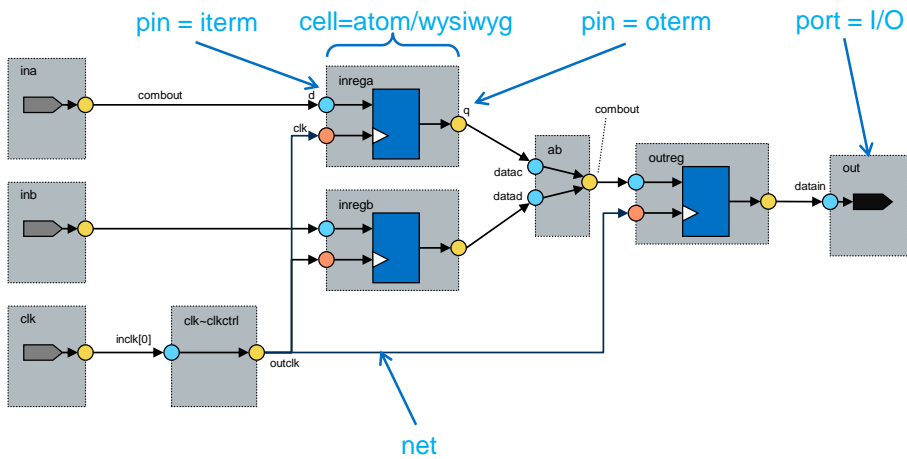


## DSP Builder Flow

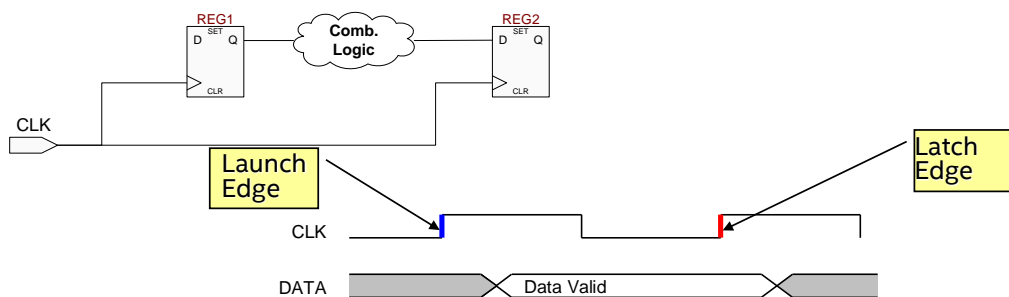


- Simulink test-bench
- Device-independent, unpipelined design
- Configuration parameters
  - Clock speed, data rate, input data width, channel count, Avalon-MM interface definition, ...
- Target device
- Cycle- and bit-accurate simulation
- Optimize design
- Fixed and Floating Point
- Optimal mapping to device features
- Timing-aware pipeline insertion, balancing and scheduling
- Automated resource sharing in IP and folded subsystems
- Component ready for integration
- ModelSim testbenches for verification

## Logic Netlist Example



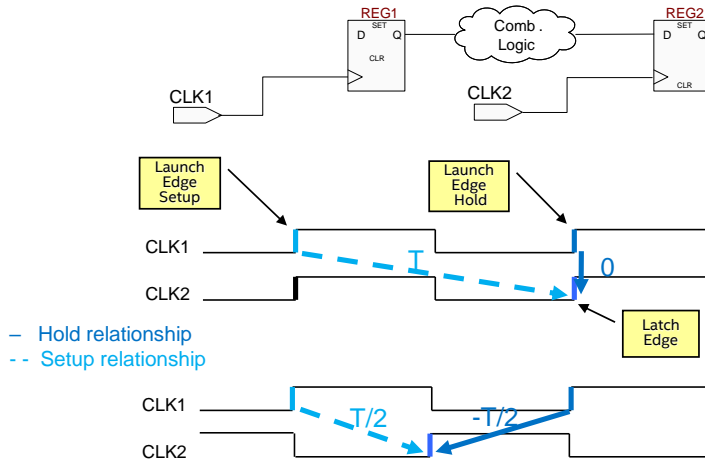
## Launch & Latch Edges



**Launch Edge:** the edge which “launches” the data from source register

**Latch Edge:** the edge which “latches” the data at destination register (with respect to the launch edge)

## Setup & Hold relationships



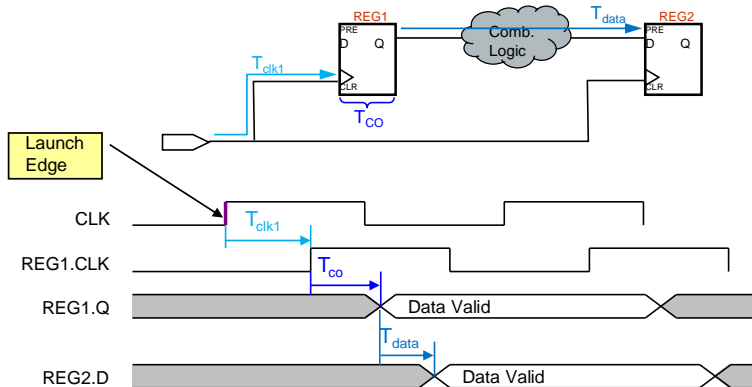
37



37

## Data Arrival Time

The time for data to arrive at destination register's D input



$$\text{Data Arrival Time} = \text{launch edge} + T_{clk1} + T_{co} + T_{data}$$

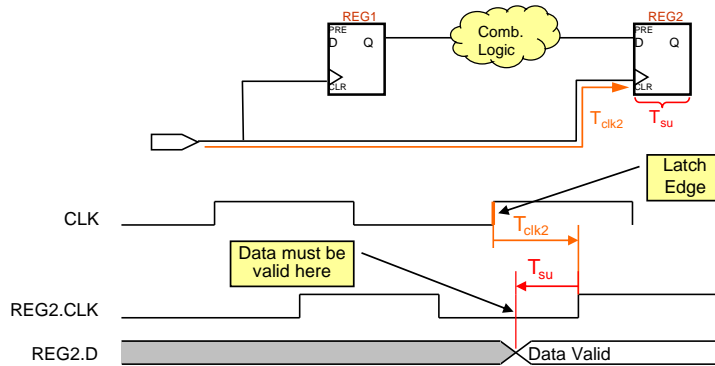
38



38

## Data Required Time - Setup

The minimum time required for the data to get latched into the destination register



Data Required Time = Clock Arrival Time -  $T_{su}$  - Setup Uncertainty

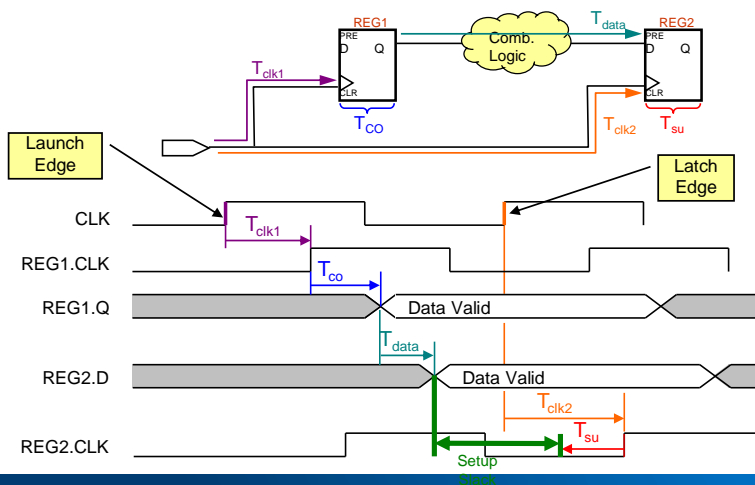
39



39

## Setup Slack

The margin by which the setup timing requirement is met. It ensures launched data arrives in time to meet the latching requirement.



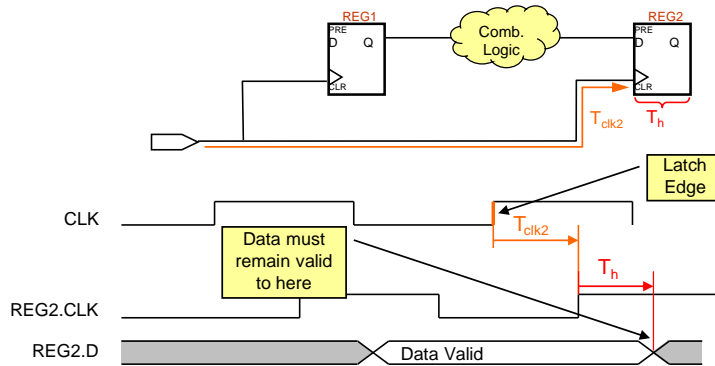
40



40

## Data Required Time - Hold

The minimum time required for the data to get latched into the destination register



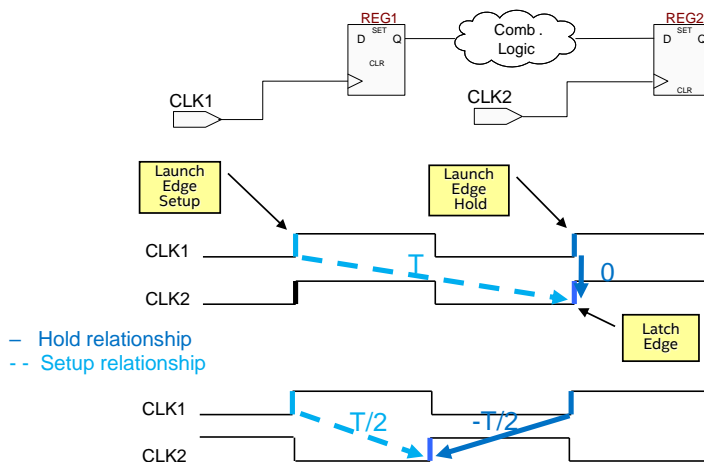
Data Required Time = Clock Arrival Time +  $T_h$  + Hold Uncertainty

41



41

## Setup & Hold relationships



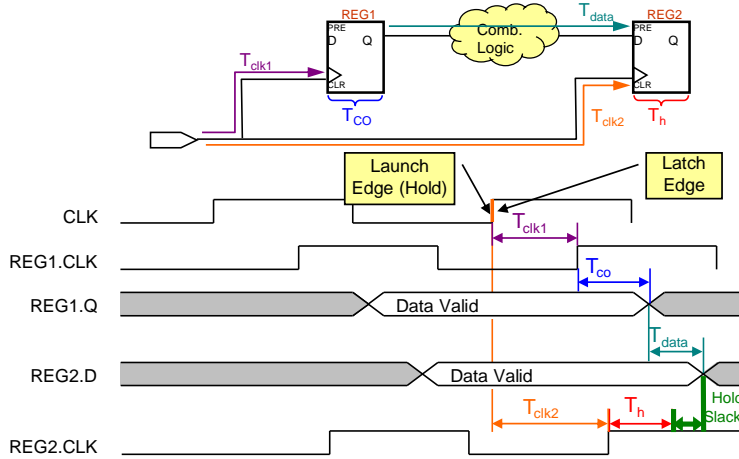
42



42

# Hold Slack

The margin by which the hold timing requirement is met. It ensures latch data is not corrupted by data from another launch edge.



# TimeQuest Prime

The screenshot shows the TimeQuest Prime interface with a report for a hold slack violation. The report includes a table of path statistics, a data arrival path table, and a data required path table. A waveform view on the right shows the timing relationship between the launch and latch clocks, highlighting the hold slack violation.

| Slack  | From Node | To Node  | Launch Clock | Latch Clock | Relationship | Clock Skew | Data Delay |
|--------|-----------|--|--------------|-------------|--------------|------------|------------|
| -4.668 | FPGA_CLK1 | sig_signaltap_auto_signaltap_block_trigger_in_reg[0] | FPGA_CLK1    | FPGA_CLK1   | 0.000        | 3.732      | -0.762     |

Path #1: Hold slack is -4.668 (VIOLATED)

| Path Summary      | Statistics | Data Path | Waveform |      |   |                    |  |  |
|-------------------|------------|-----------|----------|------|---|--------------------|--|--|
| Data Arrival Path |            |           |          |      |   |                    |  |  |
| 1                 | 0.000      | 0.000     | RR       | IC   | 1 | IDBUF_X0_Y18_N22   | FPGA_CLK1-input                                      | launch edge time                                     |
| 2                 | 0.000      | 0.000     | RR       | CELL | 3 | IDBUF_X0_Y18_N22   | FPGA_CLK1-input                                      | clock path   |
| 3                 | 0.000      | 0.000     | RR       | IC   | 1 | PLL_3              | inst1inst1PLL_ADC_intern_instatabl_0ed1p17rck[0]     | clock network delay                                  |
| 4                 | -0.762     | -0.762    | RR       | IC   | 1 | PLL_3              | inst1inst1PLL_ADC_intern_instatabl_0ed1p17rck[0]     | data path  |
| 5                 | -5.381     | -9.276    | RR       | COMP | 1 | PLL_3              | inst1inst1PLL_ADC_intern_instatabl_0ed1p17rck[0]     | inst1inst1PLL_ADC_intern_instatabl_0ed1p17rck[0]     |
| 6                 | -2.024     | 2.557     | RR       | IC   | 1 | CLKCTRL_013        | inst1inst1PLL_ADC_intern_instatabl_0ed1p17rck[0]     | inst1inst1PLL_ADC_intern_instatabl_0ed1p17rck[0]     |
| 7                 | -2.024     | 0.000     | RR       | CELL | 2 | CLKCTRL_013        | inst1inst1PLL_ADC_intern_instatabl_0ed1p17rck[0]     | inst1inst1PLL_ADC_intern_instatabl_0ed1p17rck[0]     |
| 8                 | -0.991     | 1.833     | RR       | IC   | 1 | LCCOMB_X39_Y23_N24 | auto_signaltap_block_trigger_in_reg[0]-feedercomout  | auto_signaltap_block_trigger_in_reg[0]-feedercomout  |
| 9                 | -0.832     | 0.159     | RR       | CELL | 1 | LCCOMB_X39_Y23_N24 | auto_signaltap_block_trigger_in_reg[0]-feedercomout  | auto_signaltap_block_trigger_in_reg[0]-feedercomout  |
| 10                | -0.832     | 0.000     | RR       | IC   | 1 | FF_X39_Y23_N25     | auto_signaltap_block_trigger_in_reg[0]               | auto_signaltap_block_trigger_in_reg[0]               |
| 11                | -0.762     | 0.070     | RR       | CELL | 1 | FF_X39_Y23_N25     | sig_signaltap_auto_signaltap_block_trigger_in_reg[0] | sig_signaltap_auto_signaltap_block_trigger_in_reg[0] |

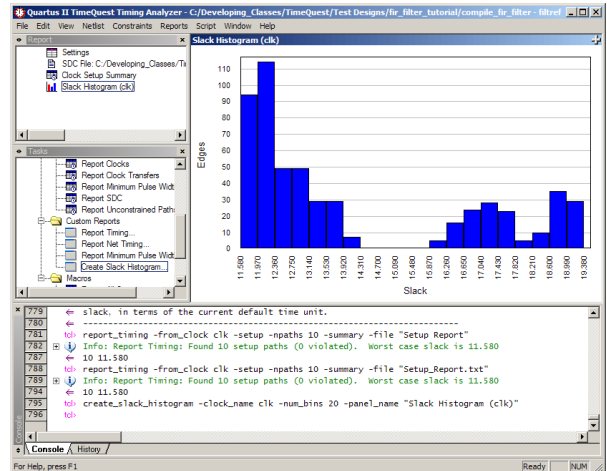
Data Required Path

| Path Summary       | Statistics | Data Path | Waveform |      |   |                     |  |  |
|--------------------|------------|-----------|----------|------|---|---------------------|--|--|
| Data Required Path |            |           |          |      |   |                     |  |  |
| 1                  | 0.000      | 0.000     | RR       | IC   | 1 | latch edge time     | latch edge time                                      |  |
| 2                  | 3.732      | 3.732     | RR       | CELL | 1 | clock path          | clock path   |  |
| 3                  | 3.732      | 3.732     | RR       | IC   | 1 | clock network delay | clock network delay                                  |  |
| 4                  | 3.732      | 0.000     | RR       | IC   | 1 | clock uncertainty   | clock uncertainty                                    |  |
| 5                  | 3.906      | 0.174     | RR       | CELL | 1 | FF_X39_Y23_N25      | sig_signaltap_auto_signaltap_block_trigger_in_reg[0] | sig_signaltap_auto_signaltap_block_trigger_in_reg[0] |

# TimeQuest Timing Analyzer

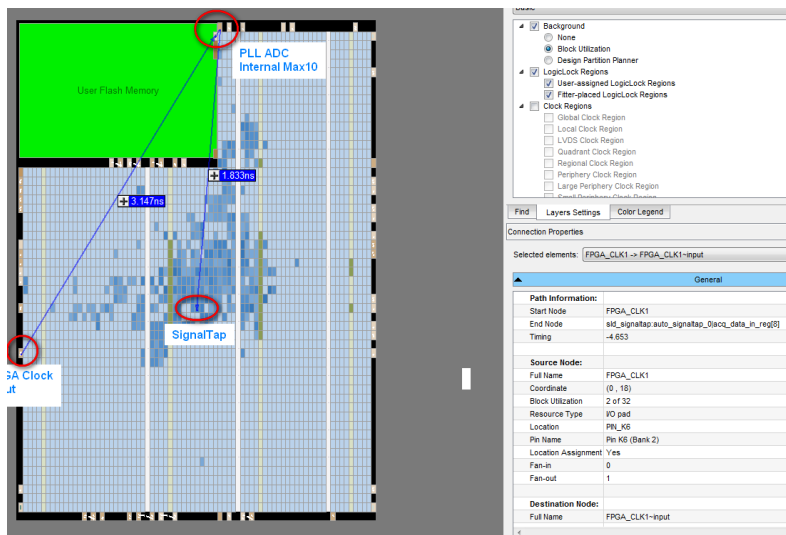
## Features

- Synopsys Design Constraints (SDC) support
  - Standardized constraint methodology
- Easy-to-use interface
  - Constraint entry
  - Standard reporting
- Scripting emphasis
  - Presentation focuses on using GUI



45

# Chip Planner



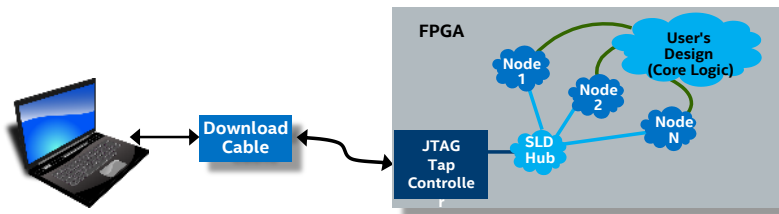
46

## On-chip Debug Technology

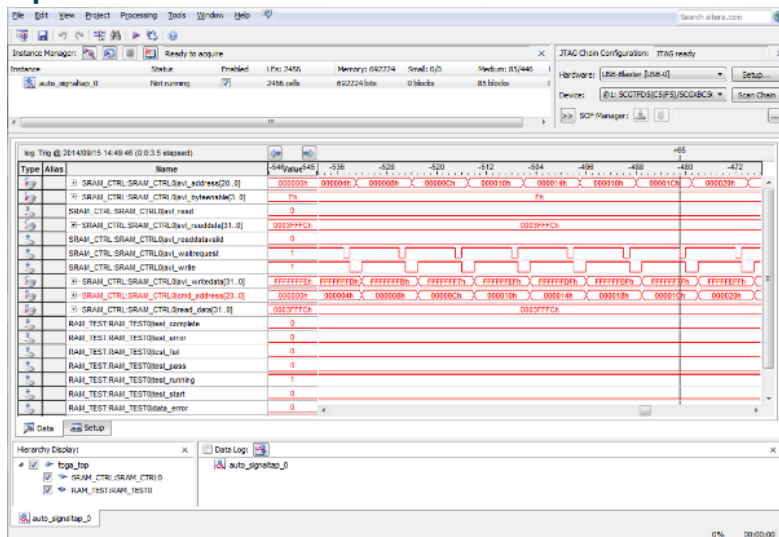
Debug tools communicate with the FPGA via standard JTAG interface

Multiple debug functions can share the JTAG interface simultaneously

- Altera's system-level debugging (SLD) hub technology makes this possible



## SignalTap Window - Data







# HW BUILDING BLOCKS



## Elementary math functions supporting floating point

Coverage of ~70 elementary math functions

Patented & published efficient mapping to FPGA hardware

- Polynomial approximation, Horner's method, truncated multipliers, ...

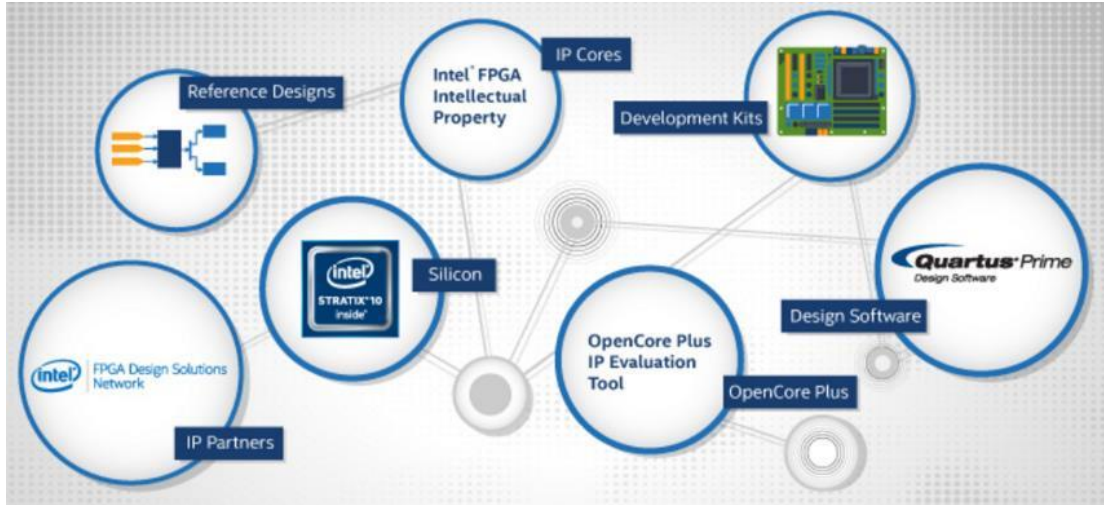
Compliant to OpenCL & IEEE754 accuracy standards

Rounding mode options for fundamental operators

Half- to Double-precision

|  |   |   |  |
|--|---|---|--|
|  | <b>Exp., Log and Power</b><br>FPLn<br>FPLn1px<br>FPLog10<br>FPLog2<br>FPExp<br>FPExpFPC<br>FPExpM1<br>FPExp2<br>FPExp10<br>FPPowr | <b>Trigonometrics of pi*x</b><br>FPSinPiX<br>FPCosPiX<br>FPTanPiX<br>FPCotPiX   | <b>Trigonometrics misc</b><br>FPHypot<br>FPRangeReduction  |
|  | <b>Trig with argument reduction</b><br>FPSinX<br>FPCosX<br>FPSinCosX<br>FPTanX<br>FPCotX  | <b>Inverse trigonometric functions</b><br>FPArcsinX\<br>FPArcsinPi<br>FPArccosX<br>FPArccosPi<br>FPArctanX<br>FPArctanPi<br>FPArctan2 | <b>Basic Floating Point</b><br>FPAdd<br>FPAddExpert<br>FPAddN<br>FPSubExpert<br>FPAddSub \<br>FPAddSubExpert<br>FPFusedAddSub<br>FPMul<br>FPMulExpert<br>FPConstMul<br>FPAcc<br>FPSqrt<br>FPDIVSqrt<br>FPRecipSqrt<br>FPCbrt<br>FPDiv<br>FPinverse<br>FPFloor<br>FPCeil<br>FPRound<br>FPRint<br>FPFrac<br>FPMod<br>FPDim<br>FPAbs<br>FPMin<br>FPMax<br>FPMinAbs<br>FPMaxAbs<br>FPMinMaxFused<br>FPMinMaxAbsFused<br>FPCompare<br>FPCCompareFused |
|  |   | <b>Conversion</b><br>FXPToFP<br>FPToFXP<br>FPToFXPExpert<br>FPToFXPFused<br>FPToFP  |  |
|  |   | <b>Macro Operators</b><br>FPFusedHorner<br>FPFusedHornerExpert<br>FPFusedHornerMulti<br>FPFusedMultiFunction                          |  |

# Intel FPGA IP ecosystem



# Intel FPGA IP Catalogue

[https://www.altera.com/content/dam/altera-www/global/en\\_US/pdfs/literature/sg/product-catalog.pdf](https://www.altera.com/content/dam/altera-www/global/en_US/pdfs/literature/sg/product-catalog.pdf)

| PRODUCT NAME   | VENDOR NAME         | PRODUCT NAME   | VENDOR NAME         |
|--|---------------------|--|---------------------|
| <b>ARITHMETIC</b>  |                     | <b>VIDEO AND IMAGE PROCESSING</b>                              |                     |
| Floating Point Megafunctions                             | Intel               | Video and Image Processing Suite <sup>1</sup>                  | Intel               |
| Floating Point Arithmetic Co-Processor                   | Digital Core Design | HD JPEG 2000 Encoders/Decoders                                 | IntoPIX             |
| Floating Point Arithmetic Unit                           | Digital Core Design | TICO Lightweight Video Compression                             | IntoPIX             |
| <b>ERROR DETECTION/CORRECTION</b>                        |                     | Multi-Channel JPEG 2000 Encoder and Decoder Cores              | Barco Silex         |
| Reed-Solomon Encoder/Decoder II <sup>1</sup>             | Intel               | VC-2 High Quality Video Decoder                                | Barco Silex         |
| Viterbi Compiler, High-Speed Parallel Decoder            | Intel               | VC-2 High Quality Video Encoder                                | Barco Silex         |
| Viterbi Compiler, Low-Speed/ Hybrid Serial Decoder       | Intel               | MPEG-2 TS Encapsulator/Decapsulator for SMPTE2022 1/2          | IntoPIX             |
| Turbo Encoder/Decoder                                    | Intel               | JPEG Encoders  | CAST, Inc.          |
| High-Speed Reed Solomon Encoder/Decoder                  | Intel               | Ultra-fast, 4K-compatible, AVC/ H.264 Baseline Profile Encoder | CAST, Inc.          |
| BCH Encoder/Decoder                                      | Intel               | Low-Power AVC / H.264 Baseline Profile Encoder                 | CAST, Inc.          |
| Low-Density Parity Check Encoder/Decoder                 | Intel               | H.265 Main Profile Video Decoder                               | CAST, Inc.          |
| Zip-Accel-C: GZIP/ZLIB/Deflate Data Compression Core     | CAST, Inc.          | H.265 Encoders   | Jointwave Group LLC |
| Zip-Accel-D: GUNZIP/ZLIB/Inflate Data Decompression Core | CAST, Inc.          | H.264 Encoders   | Jointwave Group LLC |
| <b>FILTERS AND TRANSFORMS</b>                            |                     | <b>DSP (CONTINUED)</b>   |                     |
| Transform (FFT)  |                     | Video Processor and Deinterlacer                               | UP, Inc.            |

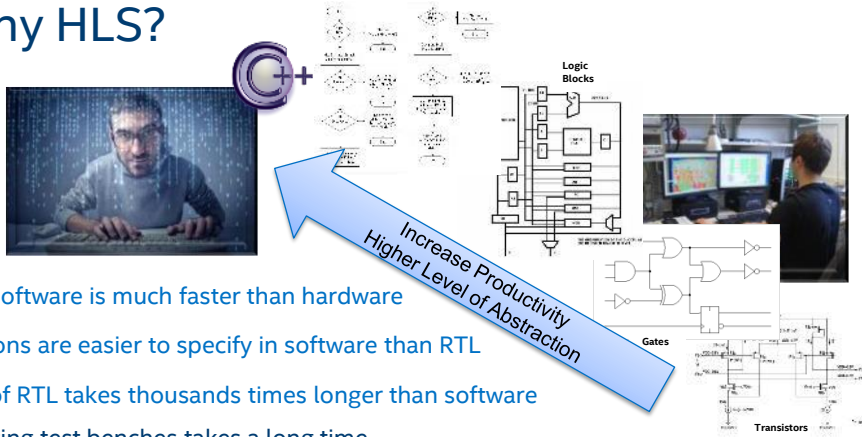
## Board Partners



# ABSTRACTING THE FPGA DEVELOPMENT FLOW



## So, Why HLS?



Debugging software is much faster than hardware

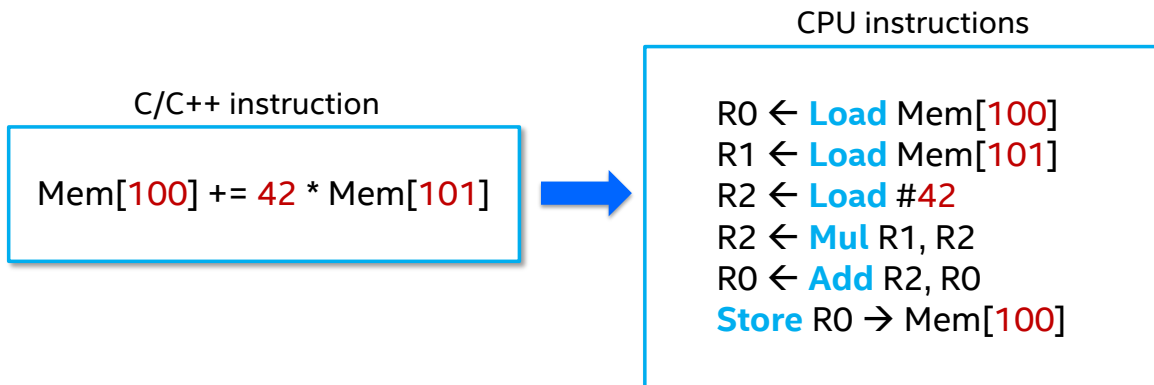
Many functions are easier to specify in software than RTL

Simulation of RTL takes thousands times longer than software

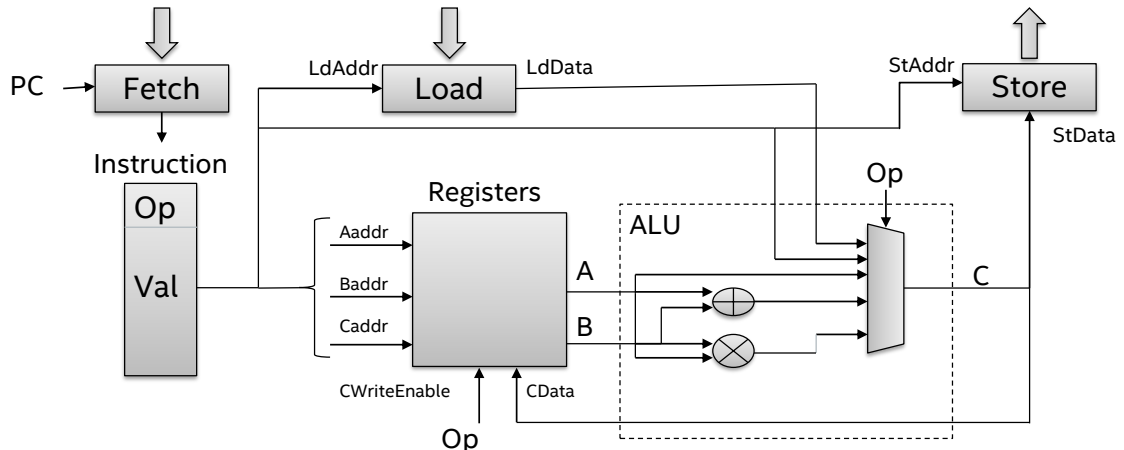
- Creating test benches takes a long time
- Updates to code requires changes to testbench to verify
- Simulation of PCIe at 1fs for 500ms = Enough said!

Design Exploration is much easier and faster in software

## Mapping a Simple Program to an FPGA



First let's take a look at execution on a *simple* CPU



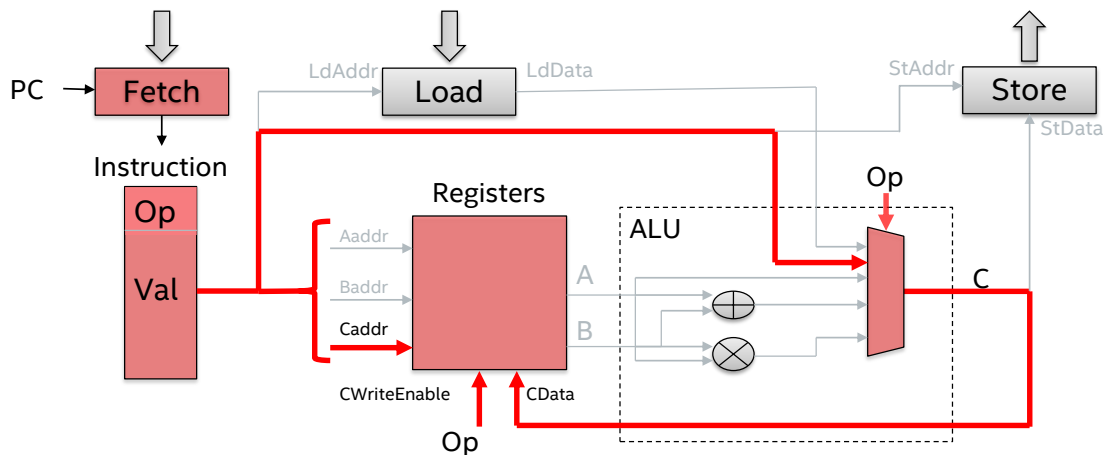
Fixed and general architecture:

- General "cover-all-cases" data-paths
- Fixed data-widths
- Fixed operations



57

Looking at a Single Instruction



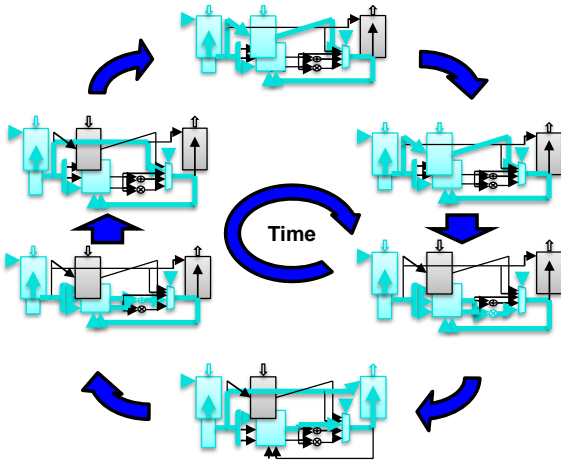
**Very inefficient use of hardware!**



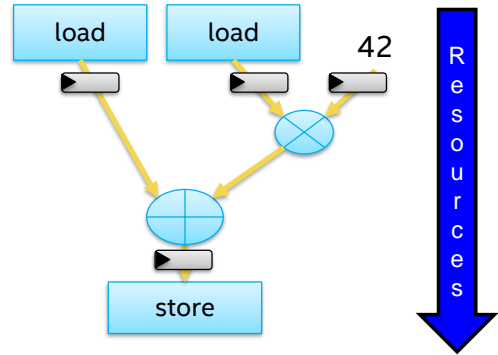
58

## Sequential Architecture vs. Dataflow Architecture

### Sequential CPU Architecture



### FPGA Dataflow Architecture

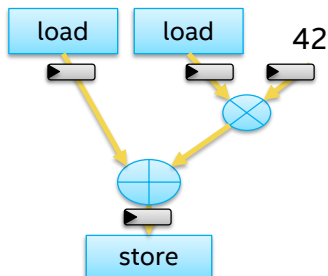


## Custom Data-Path on the FPGA Matches Your Algorithm!

High-level code

```
Mem[100] += 42 * Mem[101]
```

Custom data-path



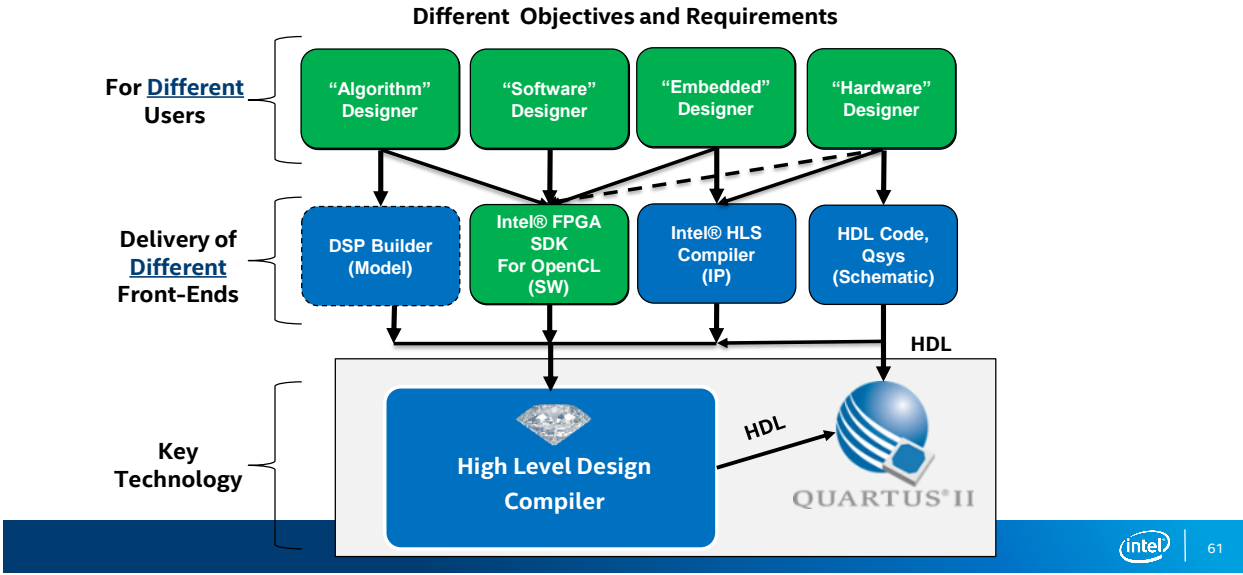
**Build exactly what you need:**

- Operations
- Data widths
- Memory size & configuration

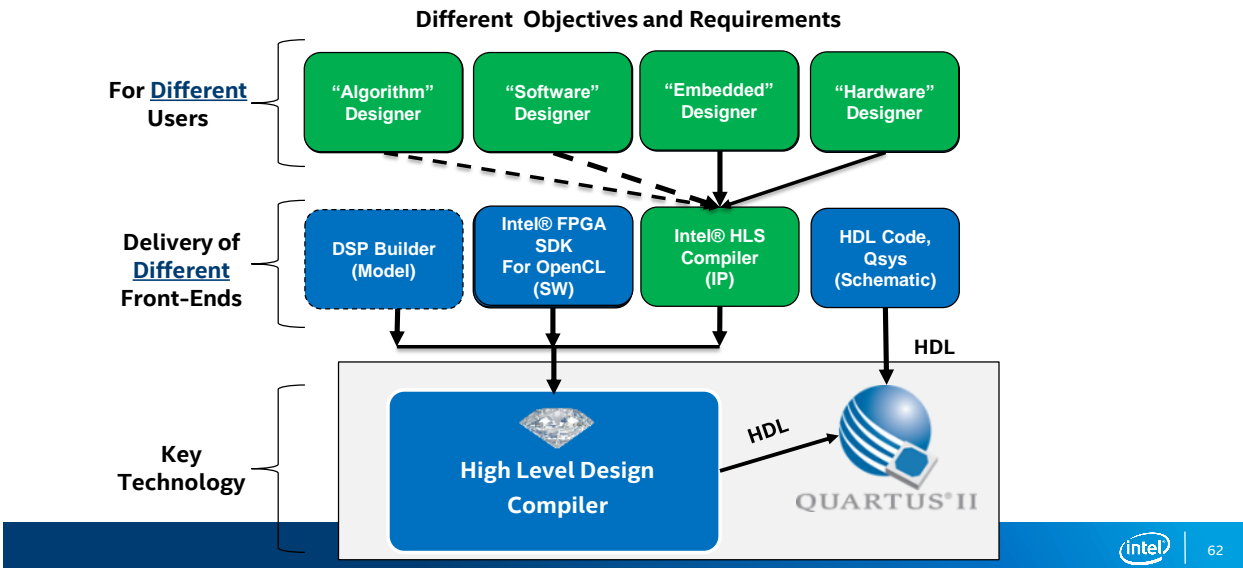
**Efficiency:**

Throughput / Latency / Power

# Different Solutions for Different Users



# Different Solutions for Different Users





# OPENCL SDK

JDEV2017

## What is OpenCL?

A software programming model for software engineers and a software methodology for system architects

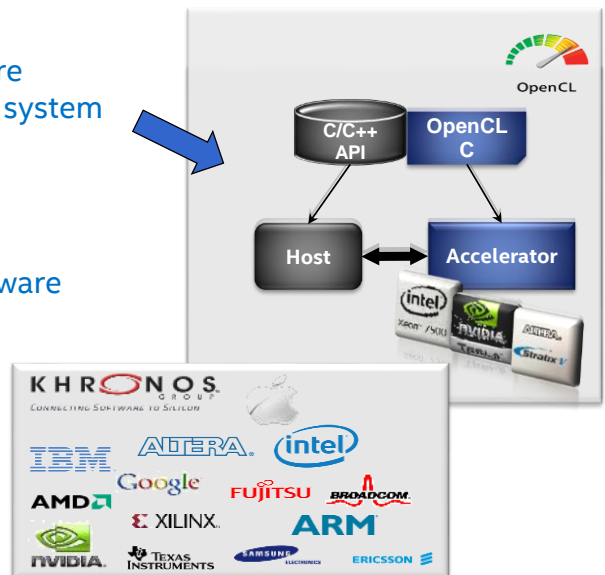
- First industry standard for heterogeneous computing

Provides increased performance with hardware acceleration

- Low Level Programming language
- C99

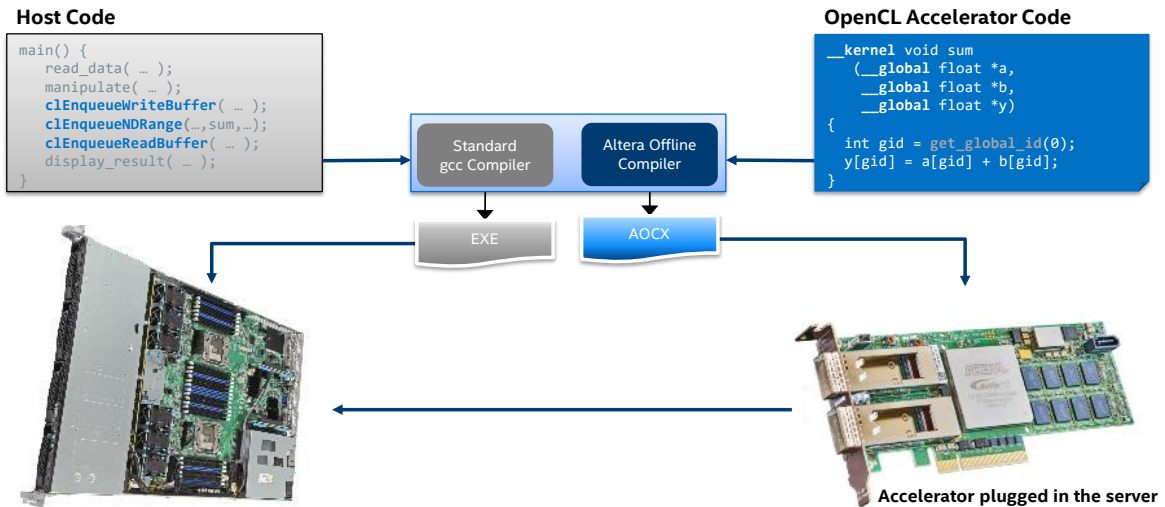
Open, royalty-free, standard

- Created by Apple
- Managed by [Khronos Group](#)
- Intel active member
- Conformant to the standard





## OpenCL Use Model: Abstracting the FPGA away

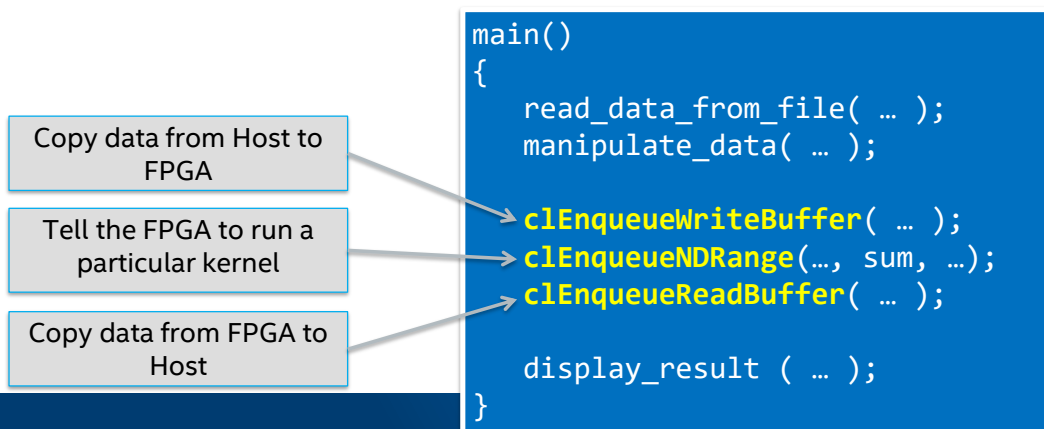


75

## Traditional OpenCL Host Program

Pure software written in standard C/C++ languages

Communicates with the accelerator devices via an API which abstracts the communication between the host processor and the kernels



76

# OpenCL Kernels

## Kernel: Data-parallel function

- Defines many parallel threads
  - Explicitly or inferred
- Keyword extensions to specify parallelism and memory hierarchy

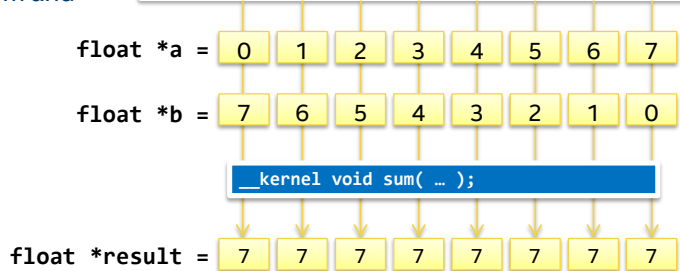
## Executed by an OpenCL device

- CPU, GPU, FPGA, DSP

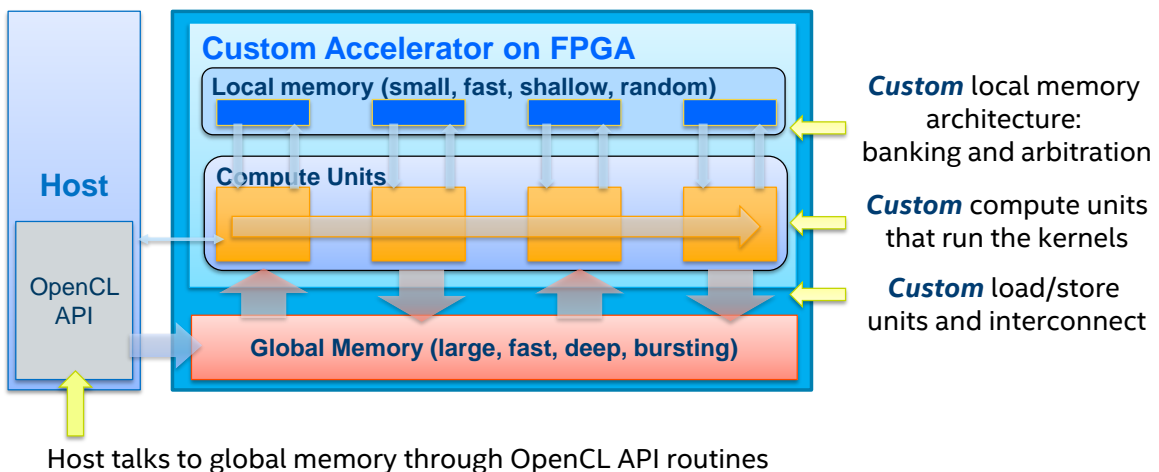
## Code portable NOT performance portable

- Between FPGAs it is!

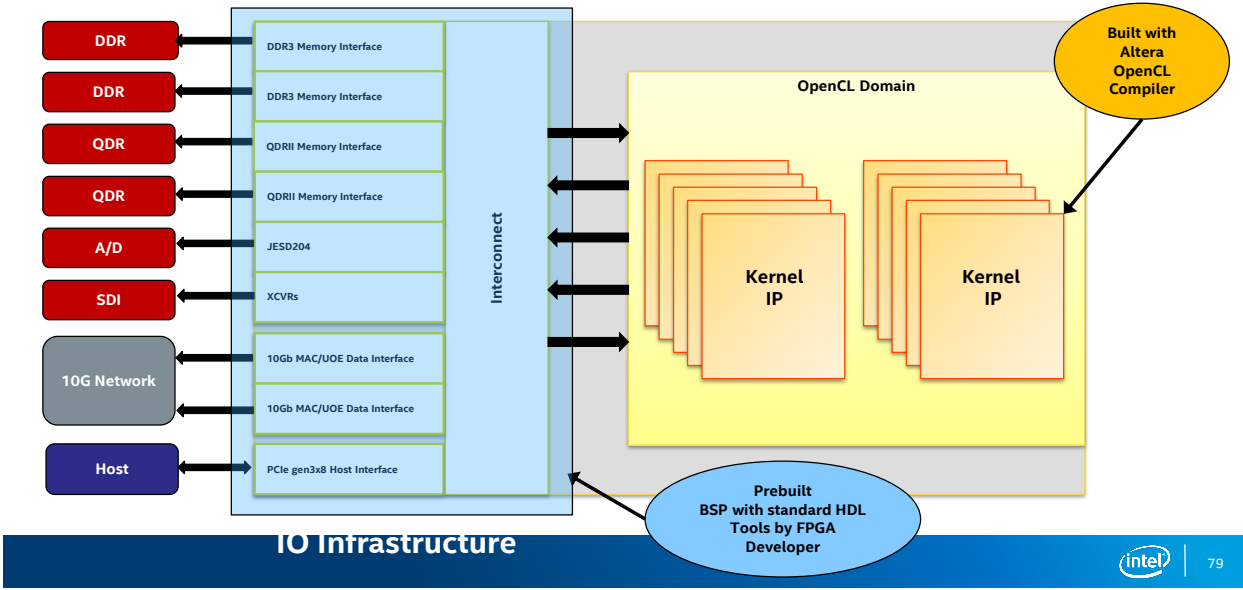
```
__kernel void sum(
  __global float *a,
  __global float *b,
  __global float *answer)
{
  int xid = get_global_id(0);
  result[xid] = a[xid] + b[xid];
}
```



## Software Programmer's View of an OpenCL Platform



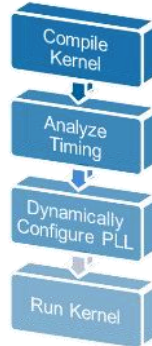
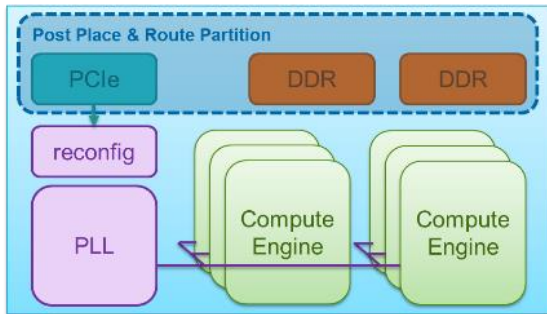
# Board Support Package



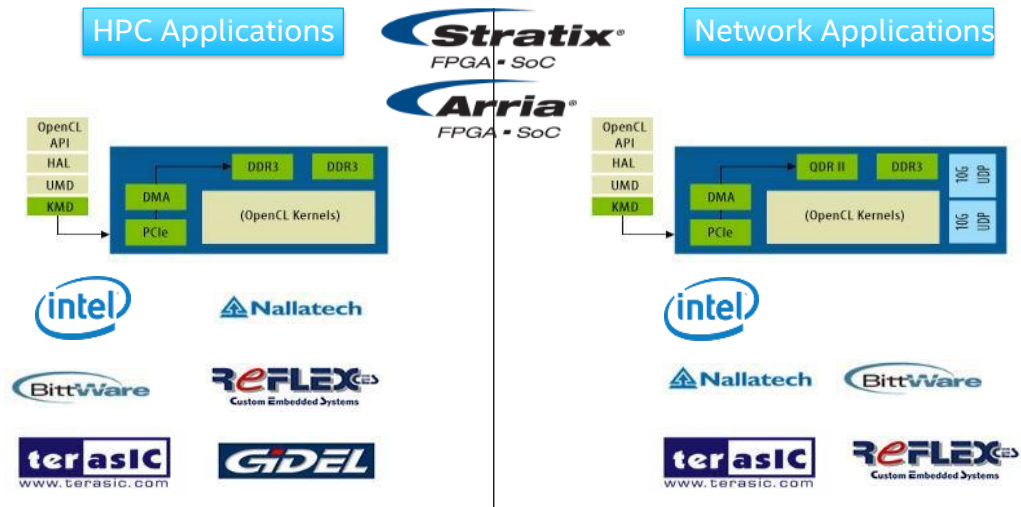
# Guaranteed Timing Closure

■ **Some interfaces have required clock frequencies**

- PCIe 125 MHz / 250 MHz
- DDR3-1600 800 MHz
- Kernel ??



## Start with OpenCL ready platforms 1/2



81

## Bittware Arria 10 Gx Specifications

### Intel Arria 10 GX FPGA - 10AX115N3F40E2SG

- Up to 1150K logic elements available
- Up to 53 Mb of embedded memory
- Up to 3,300 18x19 variable-precision multipliers

PCIe x8 interface supporting Gen1, Gen2, or Gen3

Dual QSFP cages for 2x 40GbE or 8x 10GbE

Up to 32 GBytes of DDR4 SDRAM with ECC (x72)

BMC for Intelligent Platform Management

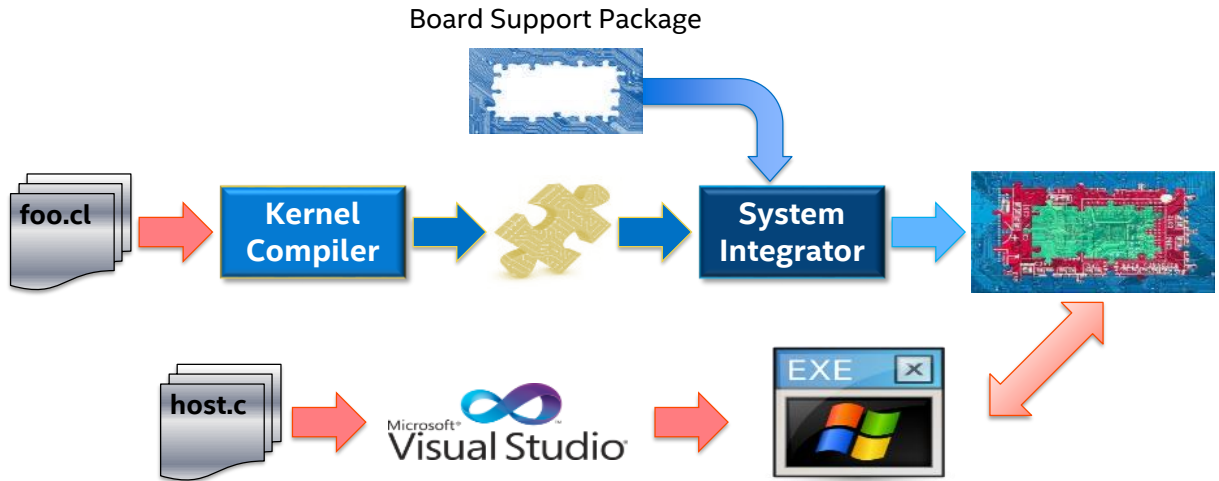
Precision clock and timing options

Utility I/O: USB 2.0

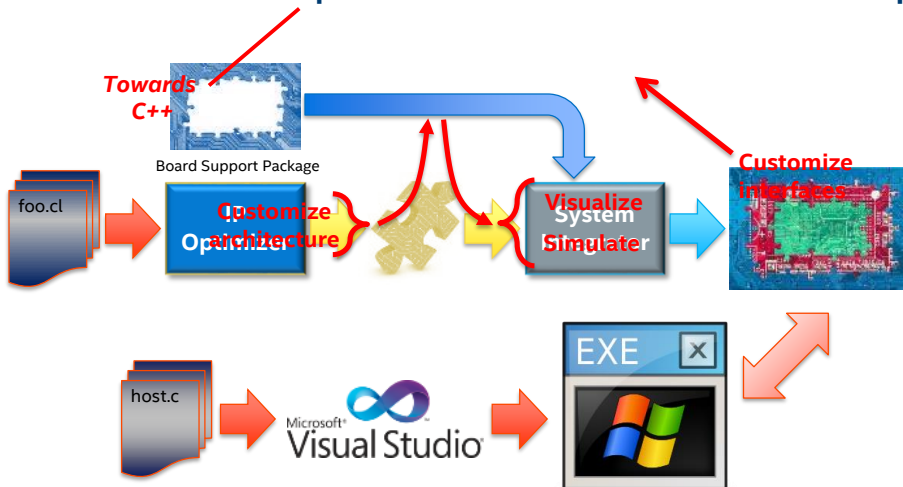


82

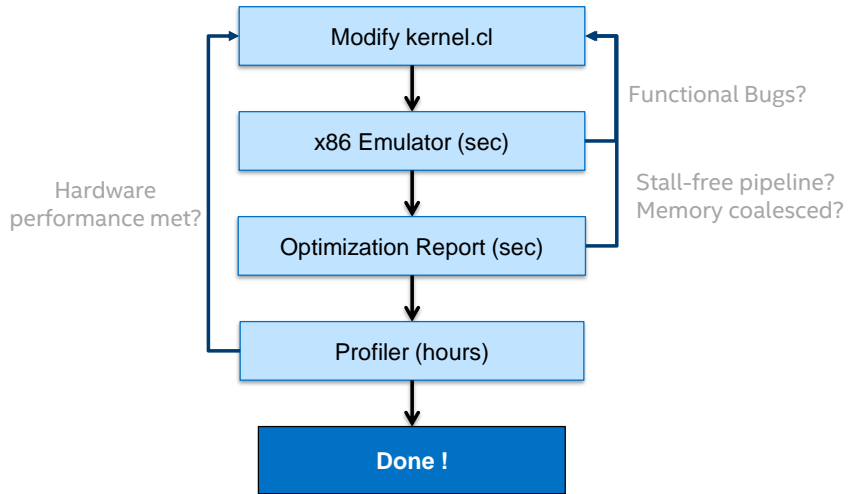
## Development Flow for FPGA



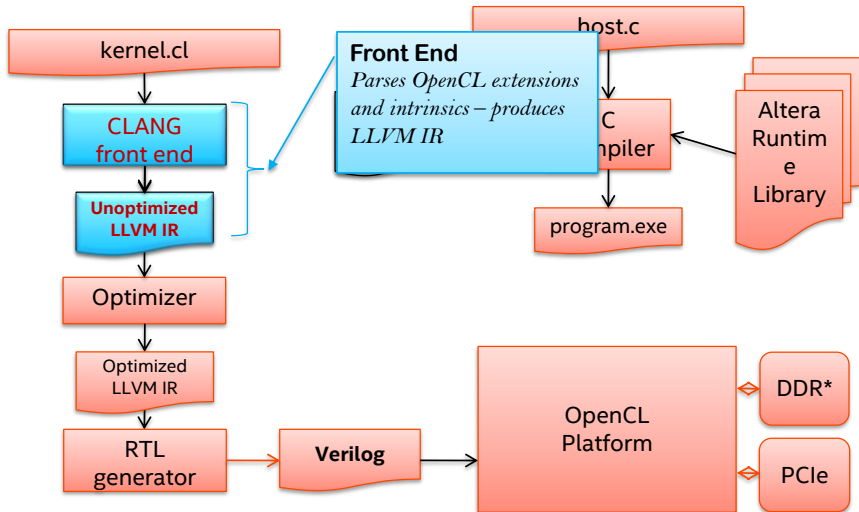
## How Does i++ Compiler for HLS Differ From OpenCL?



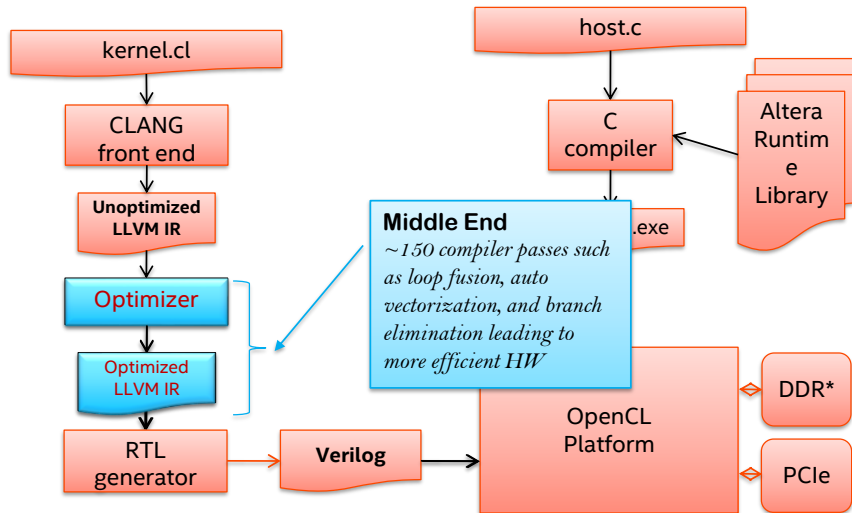
# Multi-Step AOC Compilation



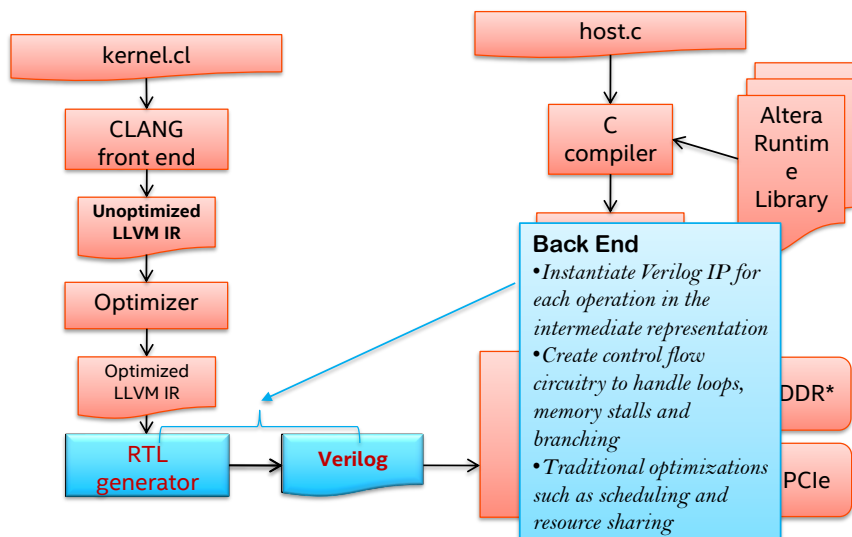
# OpenCL Compiler Flow



## OpenCL Compiler Flow



## OpenCL Compiler Flow



# Architectural Viewer (Stall Points)

Local memory is 2D [1024][4]. Compiler can prove, by access analysis, that accesses to each lower dimension are disjoint.

Local memory is well-configured into multiple banks. All loads and all stores can be serviced every clock cycle. The banking is also described in words in Area Report but not how each load/store connected to each bank.

Compile-time known address.

Selected store

Kernel Code:

```

1 #define NUM_READS 4
2 attribute((reqd_work_group_size(1024,1,1)))
3 kernel void big_lmem (global int* restrict in,
4   global int* restrict out)
5 {
6   local int lmem[1024][NUM_READS];
7   int gi = get_global_id(0);
8   int gs = get_global_size(0);
9   int li = get_local_id(0);
10  int ls = get_local_size(0);
11
12  int res = in[gi];
13  #pragma unroll
14  for (int i=0; i<NUM_READS; i++) {
15    lmem[(li-i) % ls][i] = res;
16  }
17
18  barrier(CLK_GLOBAL_MEM_FENCE);
19
20  res = 0;
21  #pragma unroll
22  for (int i=0; i< NUM_READS; i++) {
23    res ^= lmem[(li-i) % ls][i];
24  }
25
26  out[gi] = res;
27 }

```

Stall Points Graph:

Hardware Details:

|            |           |
|------------|-----------|
| Width      | 32 bits   |
| Type       | Pipelined |
| Stall free | Yes       |

## Area Report

Shows estimated resource usage for each line of code

Provides details explaining reason for inefficiencies

Information on how to improve your design

Enhanced accuracy of estimates and relationship to code

Kernel Code:

```

1 #define NUM_READS 8
2 #define NUM_WRITES 8
3 #define NUM_BARRIERS 2
4
5 __attribute__((reqd_work_group_size(1024,1,1)))
6 kernel void big_lmem (global int* restrict in,
7   global int* restrict out) {
8
9   local int lmem[1024];
10  int gi = get_global_id(0);
11  int gs = get_global_size(0);
12  int li = get_local_id(0);
13  int res = in[gi];
14  #pragma unroll
15  for (int i=0; i<NUM_WRITES; i++) {
16    lmem[li - i] = res;
17    res >>= 1;
18  }
19  barrier(CLK_GLOBAL_MEM_FENCE);
20  res = 0;
21  #pragma unroll
22  for (int i=0; i< NUM_READS; i++) {
23    res ^= lmem[li - i];
24  }
25
26  out[gi] = res;
27 }

```

Area Report Table:

| Resource             | Usage      | Peak        | Max     | Min    | Details  |
|----------------------|------------|-------------|---------|--------|--|
| big_lmem (Logic: 2K) | 8532 (14%) | 10000 (14%) | 78 (0%) | 0 (0%) | Kernel object logic.   |
| Function overhead    | 1000       | 1795        | 0       | 0      | Local memory: Primarily inefficient configurations. Local memory is used up to maximum period of 23. Implemented as 4 local banks, reconfigured in time. |
| MEMS (MEMS)          | 152        | 1024        | 0       | 0      | Local memory: Primarily inefficient configurations. Local memory is used up to maximum period of 23. Implemented as 4 local banks, reconfigured in time. |
| MEMS (2) (MEMS)      | 4487 (14%) | 7487 (14%)  | 27 (0%) | 0 (0%) | Local memory: Primarily inefficient configurations. Local memory is used up to maximum period of 23. Implemented as 4 local banks, reconfigured in time. |

Details for big\_lmem (Logic: 2K):

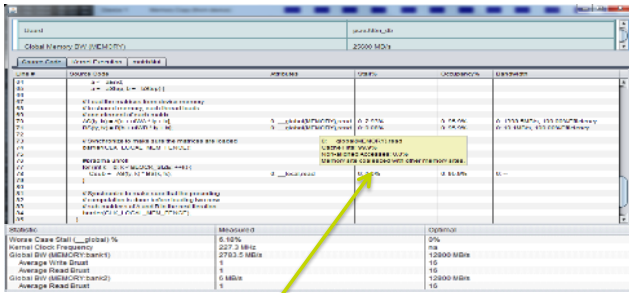
- Kernel object logic.
- Local memory: Primarily inefficient configurations. Local memory is used up to maximum period of 23. Implemented as 4 local banks, reconfigured in time.
- MEMS (MEMS): Primarily inefficient configurations. Local memory is used up to maximum period of 23. Implemented as 4 local banks, reconfigured in time.
- MEMS (2) (MEMS): Primarily inefficient configurations. Local memory is used up to maximum period of 23. Implemented as 4 local banks, reconfigured in time.



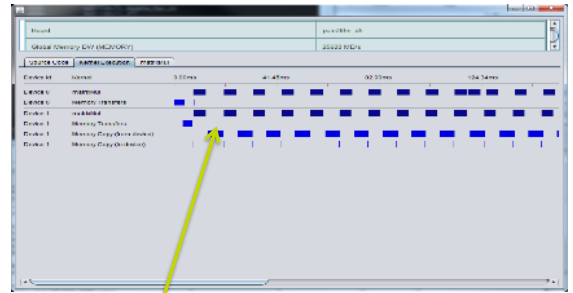
# Dynamic Profiler

Intel FPGA SDK for OpenCL enables users to get runtime information about their kernel performance

Bottlenecks, bandwidth, saturation, pipeline occupancy



Performance Stats



Execution Times



## Key Optimization Methodology is Still Data Localization

10 TFLOP floating point performance in Startix 10

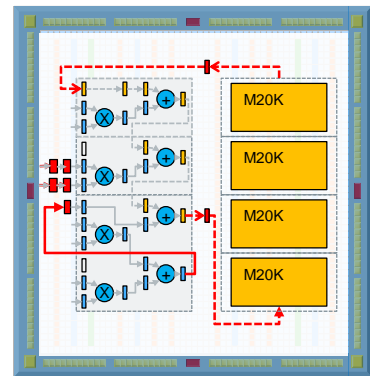
- Use every DSP, every clock cycle compute spatially

>8 TB/s memory bandwidth: keep state on chip!

- Exceeds available external bandwidth by orders of magnitude
- Random access, low latency (2 clks)
- Place all data in on-chip memory compute temporally

### Avoid costly data movement

- Highest performance/W of Intel's programmable offerings



Fine-grained & low latency between compute and memory





# OPENCL EXECUTION MODELS

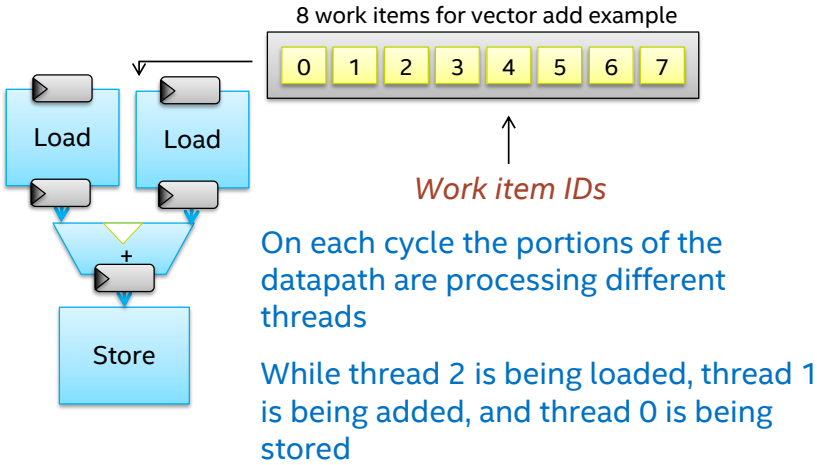
Harnessing Pipeline Parallelism



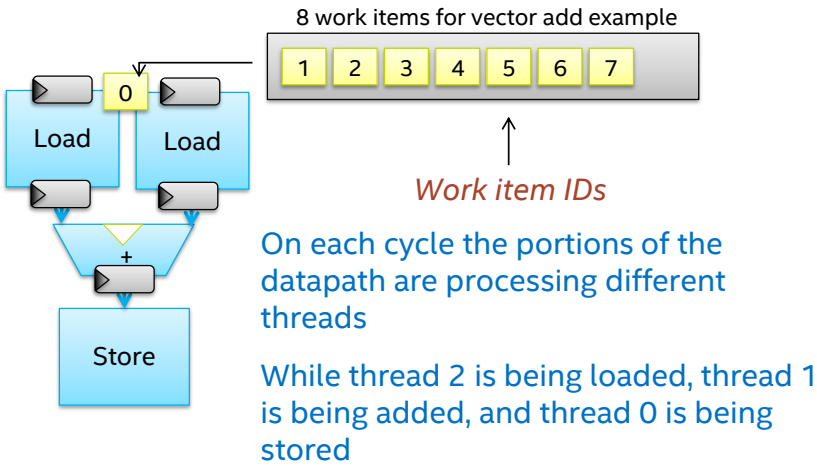
## Mapping Multithreaded Kernels to FPGAs

- **The most simple way of mapping kernel functions to FPGAs is to replicate the unrolled hardware for each thread**
  - Inefficient and wasteful
- **Better method involves taking advantage of *pipeline parallelism***
  - Attempt to create a deeply pipelined representation of a kernel
  - On each clock cycle, we attempt to send in input data for a new thread
  - Method of mapping coarse grained thread parallelism to fine-grained FPGA parallelism

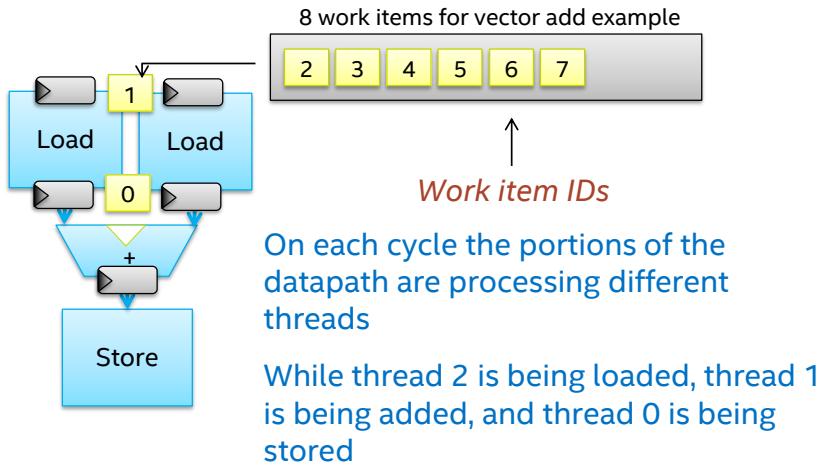
# Example Datapath for Vector Add



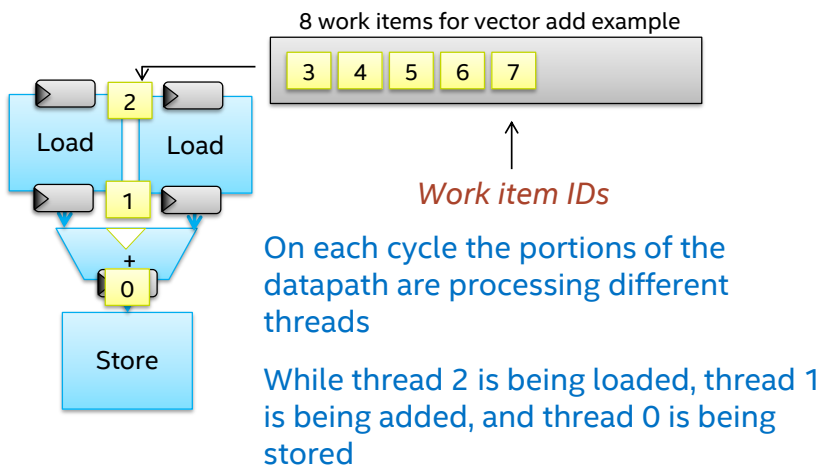
# Example Datapath for Vector Add



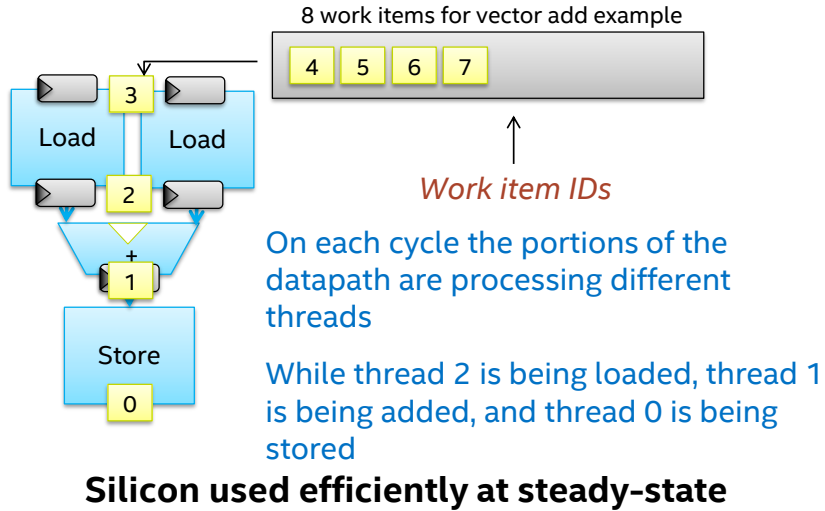
## Example Datapath for Vector Add



## Example Datapath for Vector Add



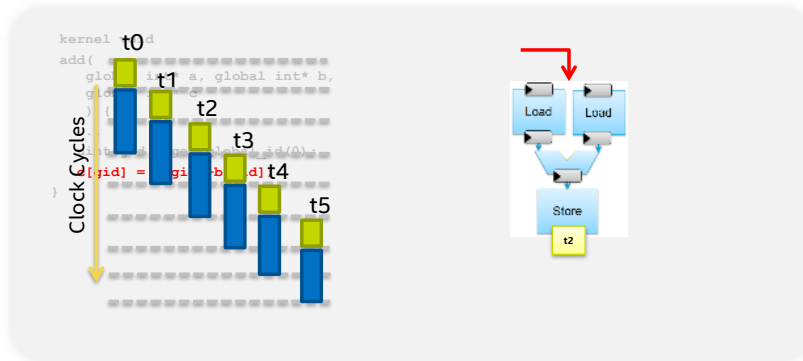
# Example Datapath for Vector Add

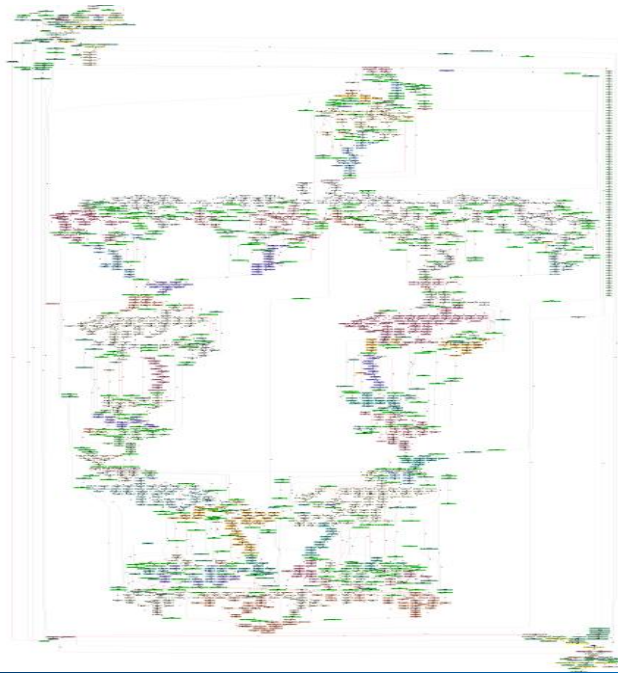


# Execution of Threads on FPGA

Better method involves taking advantage of *pipeline parallelism*

- **Throughput = 1 thread per cycle**





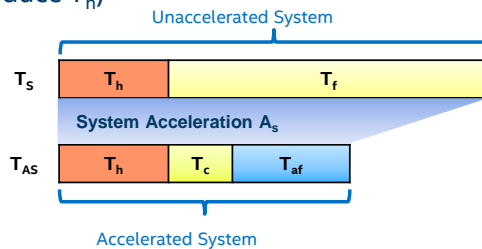
# OPTIMIZING KERNEL FOR FPGA



## Goals of OpenCL Optimization

Keep as much FPGA resources busy as possible doing useful stuff!

- Increase the number of parallel operations (Reduce  $T_{af}$ )
- Reduce communication latency of the accelerated system (Reduce  $T_c$ )
- Increase efficiency of operations (Reduce  $T_{af}$ )
- Minimize overhead (Reduce  $T_h$ )



## Optimization technics

Loop Unrolling

Vectorization

Kernel Compute duplication

Task Kernel: Single Work-Item Kernels

Channels

NDRange Execution

Memory Optimizations: Global, Local & Private

Floating Point Optimizations



# LOOP UNROLLING

JDEV2017

## Loop Unrolling

By default, loops can hinder performance

Loop unrolling replicate hardware to execute multiple loop iterations at once

Increase performance by decreasing number of iteration through loop

- Structure of the compute units built will be significantly altered
- More resource consumption
- More local memory ports may be required

Simple loops will be unrolled automatically

- AOC will generate a warning



## unroll kernel pragma

`#pragma unroll <N>` instructs AOC to attempt to unroll a loop <N> times

- Without <N>, AOC will attempt to unroll the loop fully
- Warning issued if AOC unable to unroll

```
#pragma unroll 2
for (size_t k=0; k<4; k++)
{
    mac += data_in[(gid*4)+k] * coeff[k];
}
```

Control the amount of hardware used for loops



## Loop Unrolling Example

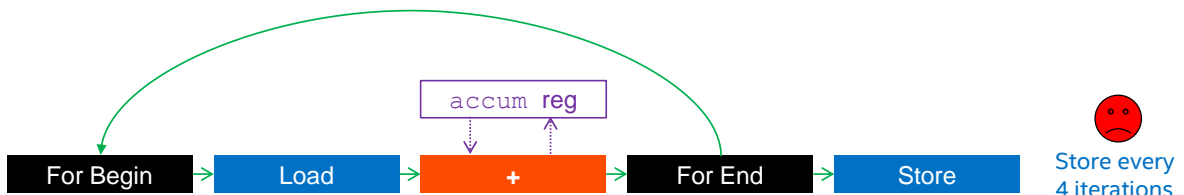
Sum of 4 values for every work-item

Store a new result every 4 iterations

```
accum = 0;

for (size_t i=0; i<4; i++)
{
    accum += data_in[(gid*4)+i];
}

sum_out[gid] = accum;
```

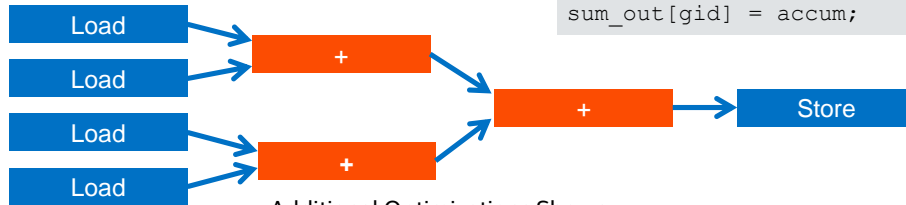


## Loop Unrolling Example: Fully Unrolled

Unroll every iteration of the loop

Store a new result every clock cycle

```
accum = 0;
#pragma unroll
for (size_t i=0; i<4; i++)
{
    accum += data_in[(gid*4)+i];
}
sum_out[gid] = accum;
```



Additional Optimizations Shown:

1. accum register removed
2. Order of operation optimization done if allowed
3. Operators removed if not needed
  - There would be 4 adders created if initial value of accum is not 0.



110

## Loop Unrolling in the Optimization Report

Loop unrolling reported in the <kernel file>.log

Reported information

- Loop location
- Nesting relationship
- Requested unroll factor
- Achieved unroll factor

```
Loop Report:
+ Fully unrolled loop (file nd_full_nested.cl line 4)
  Loop was automatically and fully unrolled.
  Add "#pragma unroll 1" to prevent automatic unrolling.
-+ Fully unrolled loop (file nd_full_nested.cl line 6)
  Loop was fully unrolled due to "#pragma unroll" annotation.
```



111



# VECTORIZATION



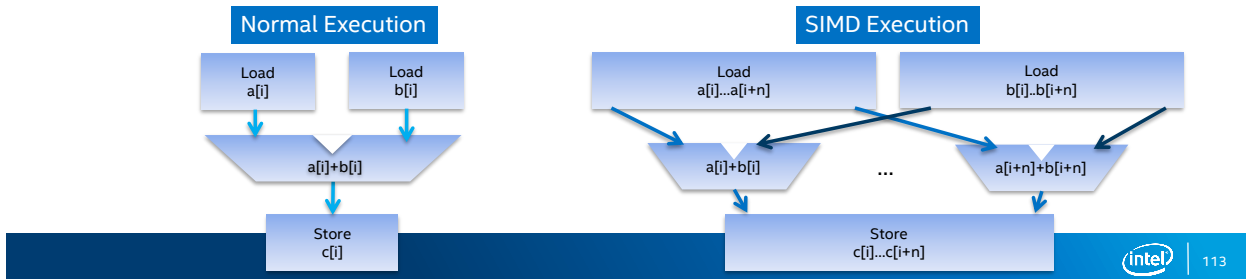
## Kernel Vectorization

Widen the pipeline to achieve higher throughput

- Allow multiple work-items from the same workgroup to execute in Single Instruction Multiple Data (SIMD) fashion

Translate scalar operations into SIMD operations

- E.g. multiplication or addition
- Can be done manually or automatically



## Vectorize Kernel Code Manually

Replicate operations in the kernel manually

Must also adjust NDRange in host application

```
__kernel void mykernel (...)
{
    size_t gid = get_global_id(0);
    result[gid] = in_a[gid] + in_b[gid];
}
```

Original  
Kernel

```
__kernel void mykernel (...)
{
    size_t gid = get_global_id(0);
    result[gid*4+0] = a[gid*4+0] + b[gid*4+0];
    result[gid*4+1] = a[gid*4+1] + b[gid*4+1];
    result[gid*4+2] = a[gid*4+2] + b[gid*4+2];
    result[gid*4+3] = a[gid*4+3] + b[gid*4+3];
}
```

Manually  
Vectorized  
Kernel



## Automatic Kernel Vectorization

Use attribute to enable automatic kernel compute unit vectorization

- Memory accesses automatically coalesced
- No need to adjust NDRange in host application

`num_simd_work_items` attribute

- Specify the number of work-items within a workgroup to be executed in a SIMD manner
  - Hardware operators automatically vectorized
- Vectorization only takes affect in the X dimension of the workgroup
- Must use with `reqd_work_group_size` attribute
  - `reqd_work_group_size` must be evenly divisible by `num_simd_work_items` in the X dimension
- Factor of 2,4,8,16

```
__attribute__((num_simd_work_items(4)))
__attribute__((reqd_work_group_size(64,1,1)))
__kernel void mykernel (...)
{...}
```



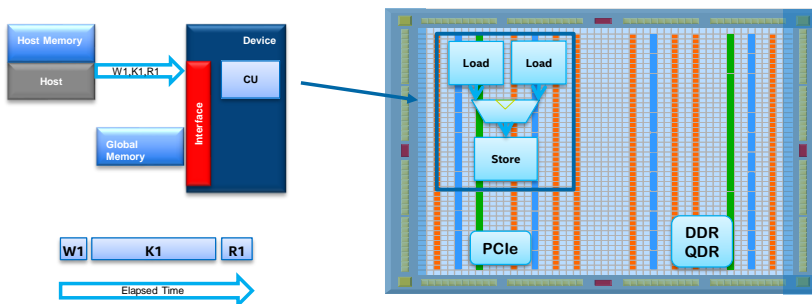


# KERNEL COMPUTE DUPLICATION



## Default Compute Units

- Only one compute unit per kernel created by default
- Workgroups distributed to compute unit in sequence



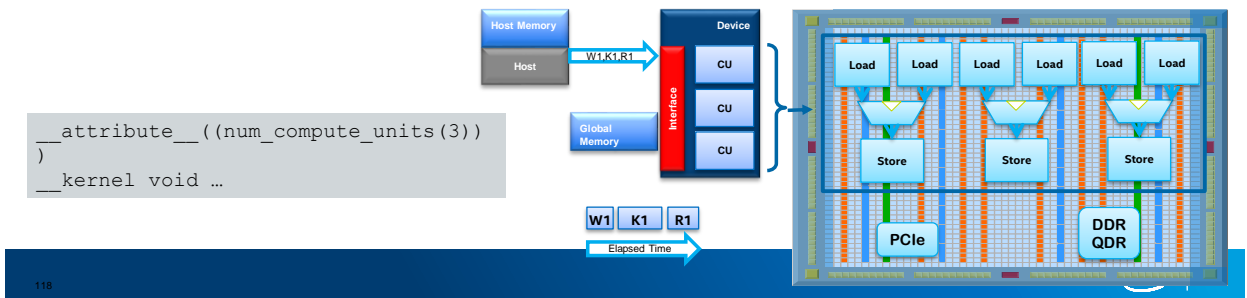
## Multiple Compute Units

Use `num_compute_unit` attribute to specify number of kernel compute units to

- Entire compute unit including all local memory, control logic, and operators replicated

Workgroups from the same NDRange kernel launch are distributed to available compute units and processed in parallel

- Need at least three times as workgroups as compute units to effectively utilize all hardware



## Example: Combining Replication and Vectorization

Resource estimates of 16 SIMD lanes indicate “no fit”

Resource estimates of 8 SIMD lanes suggest 12 lanes may fit

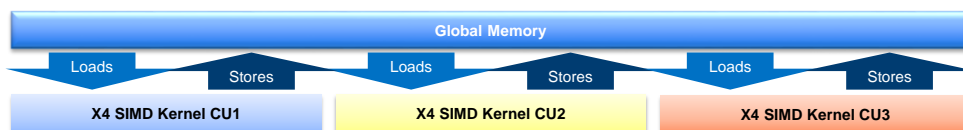
- Automatic vectorization only supports 2, 4, 8 and 16 lane configurations

Generate 12 lanes by combining `num_simd_work_items` and `num_compute_units`

```

__attribute__((num_simd_work_items(4)))
__attribute__((num_compute_units(3)))
__attribute__((reqd_work_group_size(8,8,1)))
__kernel void mykernel (...)
{ ...

```



# Kernel Attributes in the Optimization Report

Effects of all kernel attributes are reported in the <kernel file>.log file

```
=====  
Kernel: MovAu  
=====  
The kernel is compiled as an ND-Range.  
  
The kernel was vectorized for 2 SIMD work-items along the lowest dimension.  
  
The kernel was replicated 2 times due to num_compute_units attribute.  
You should have at least three times as many work-groups as the number of compute units  
to efficiently utilize all compute units.  
  
The kernel has a required work-group size of (128, 1, 1).  
=====
```



# TASK KERNEL: SINGLE WORK ITEM

## Single Work-Item Execution

### Launching kernels with NDRange of (1,1,1)

- A kernel executed on a compute unit consisting of one work-item
- Defined as a Task in OpenCL

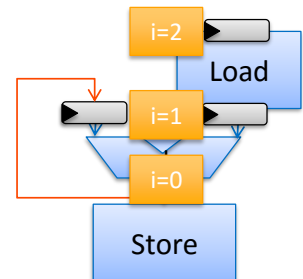
### Why?

- Data parallel (multiple work-items) execution may not be suitable for certain situations
  - Difficulties partitioning data among parallel works-items
  - Dependency across work-items
  - Data not available prior to kernel execution
  - Data not easily divided into workgroups
- Sequential programming model of tasks more similar to C programming
  - Certain usage scenario more suited for sequential programming model
  - Easier to port

## Tasks and Loop-pipelining

Allow users to express programs as a single-thread kernel

```
for (int i=1; i < n; i++) {
    c[i] = c[i-1] + b[i];
}
```



Compiler will infer parallel pipelined execution across loop iterations

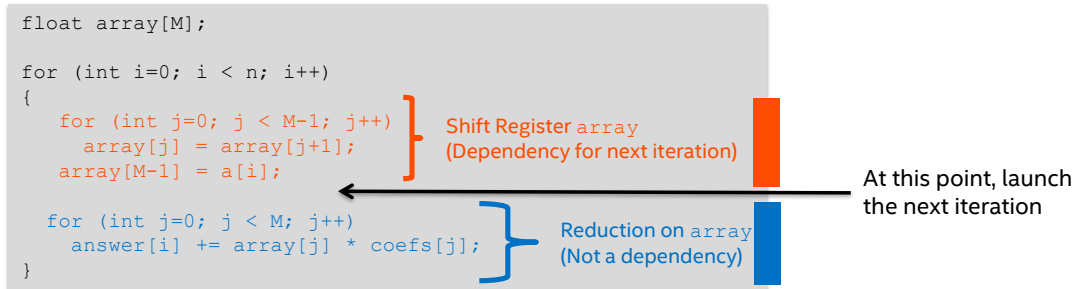
- Pipeline parallelism still leveraged to efficiently execute loops in Altera's OpenCL solution
- Dependencies resolved by the compiler
- Values transferred between loop iterations with FPGA resources
  - No need to buffer up data
  - Easy and cheap to share data through feedbacks in the pipeline



# Loop Pipelining

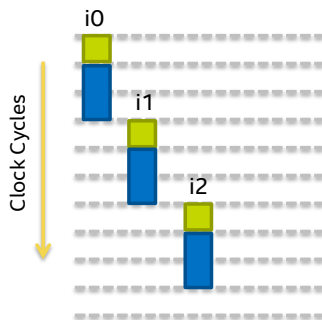
AOC will pipeline each iteration of the loop for acceleration

- Analyze any dependencies between iterations
- Schedule these operations
- Launch the next iteration as soon as possible



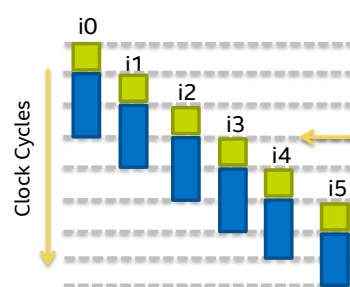
## Loop Pipelining Example

### No Loop Pipelining



*No Overlap of Iterations!*

### With Loop Pipelining

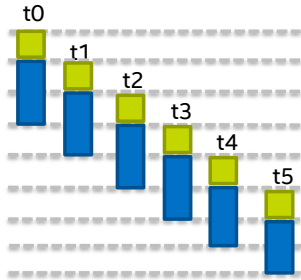


Looks like multi-threaded execution!

*Finishes Faster because Iterations Are Overlapped*

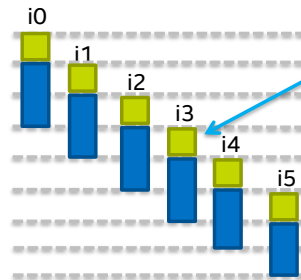
## Parallel Threads vs Loop Pipelining

### Parallel Threads



Parallel threads launch 1 thread per clock cycle in pipelined fashion

### Loop Pipelining



Loop dependencies may not be resolved in 1 clock cycle

Loop Pipelining enables Pipeline Parallelism AND the communication of state information between iterations.

- If dependency can be resolved in 1 clock cycle, then the resulting computational throughput is the same



# CHANNEL

## Communication: Channels/Pipes

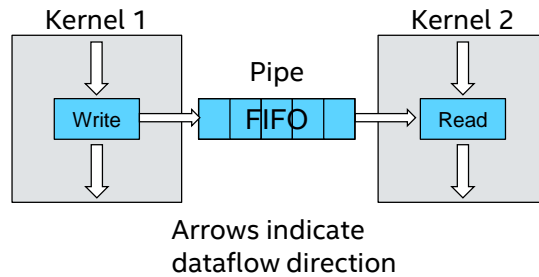
Key for systolic array architectures

Enables much lower latency data movement through data processing paths (result reuse)

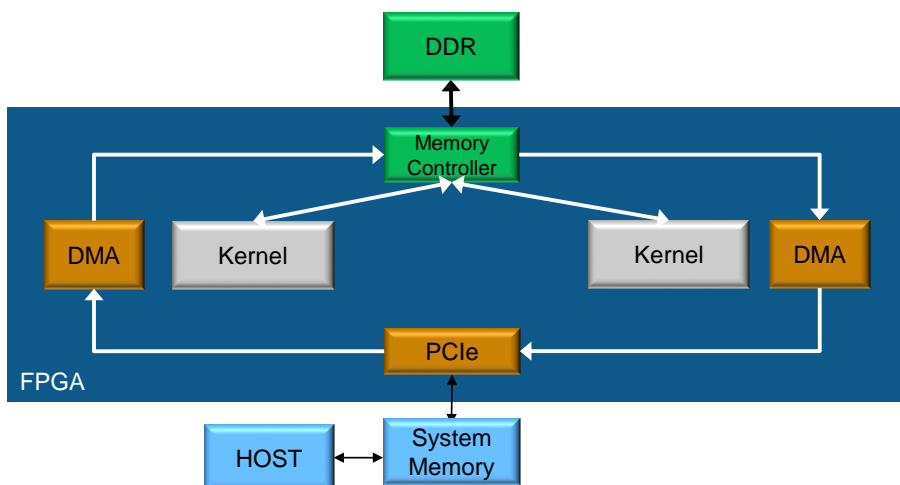
- No need to move data back and forth to external memory or cache

Provide unidirectional, FIFO-like communication into or out of kernels, or into or out of FPGA

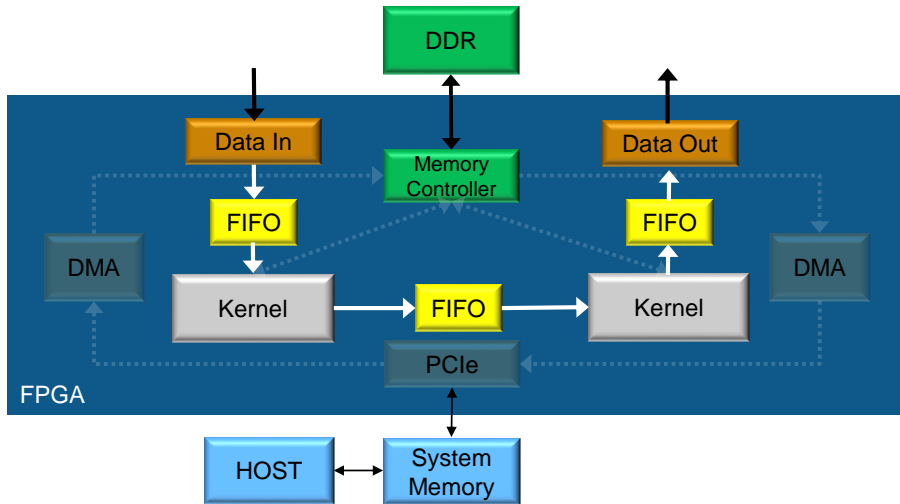
- Pipes are built on top of channels and only support kernel to kernel data movement



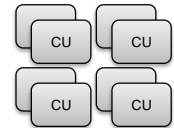
## Traditional Data Movement Without Channels



## Data Movement Using Channels



## Systolic Array of Compute Units



Replicate kernel hardware with `num_compute_units(X,Y,Z)` attribute

- Creates  $X*Y*Z$  copies of kernel pipeline
  - Increases throughput
  - When applied to NDRange kernels, these copies used to execute multiple workgroups in parallel
    - More on this in the Optimizing NDRange kernels section
  - Consumes  $X*Y*Z$  times more resources for that kernel compute unit

With single work-item kernels, AOC allows customization of kernel compute units using the `get_compute_id()` function

- Create compute ID dependent logic

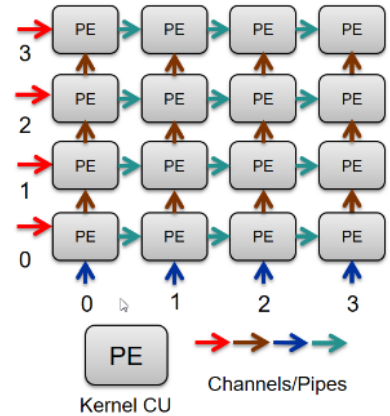
## Example with num\_compute\_units

Using compute ID to determine channel usage

```
channel float4 ch_PE_row[3][4];
channel float4 ch_PE_col[4][3];
channel float4 ch_PE_row_side[4];
channel float4 ch_PE_col_side[4];

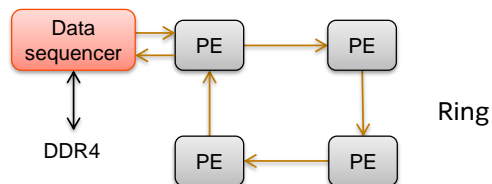
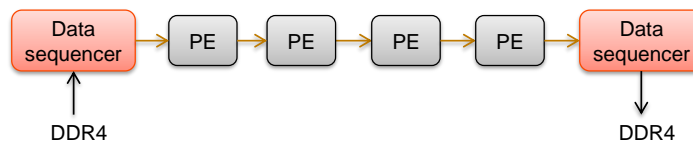
__attribute__((autorun))
__attribute__((max_global_work_dim(0)))
__attribute__((num_compute_units(4,4)))
kernel void PE() {
    float4 a,b;
    if (get_compute_id(0)==0) //First PE of row
        a = read_channel(ch_PE_col_side[col]);
    else
        a = read_channel(ch_PE_col[row-1][col]);

    if (get_compute_id(1)==0)
        ...
}
```

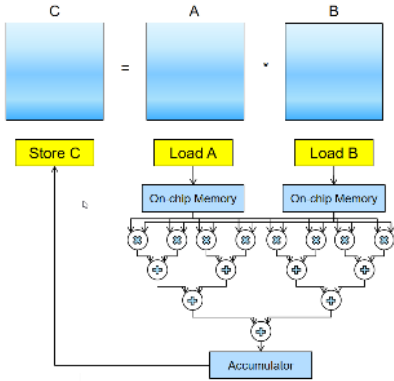


## Other Topologies Possible in FPGA Too

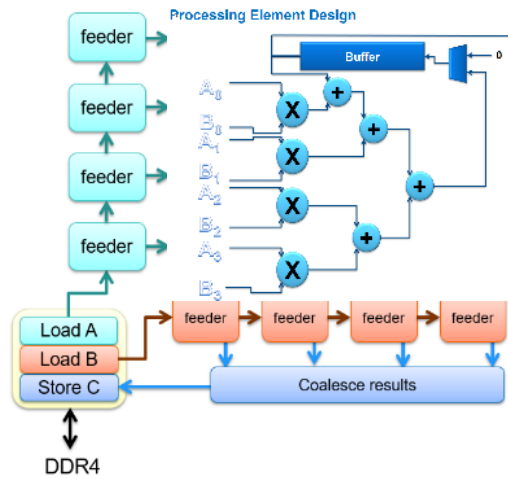
Linear Pipeline / Daisy chain



# Use Case: Matrix Multiplication



Performance: ~1 TFLOPs



# MEMORY OPTIMIZATIONS

Global Memory & Local Memory



# OpenCL Global Memory

## Global Memory Overview

- Programmers view of global memory in OpenCL
- Compiler view of global memory

## Global Memory Architecture

- Load-Store Units
- LSU Arbitration
- Const Cache

## Optimizations



## Global Memory in OpenCL

### 'global' address space

- Used to share data between host and device
- Shared between workgroups

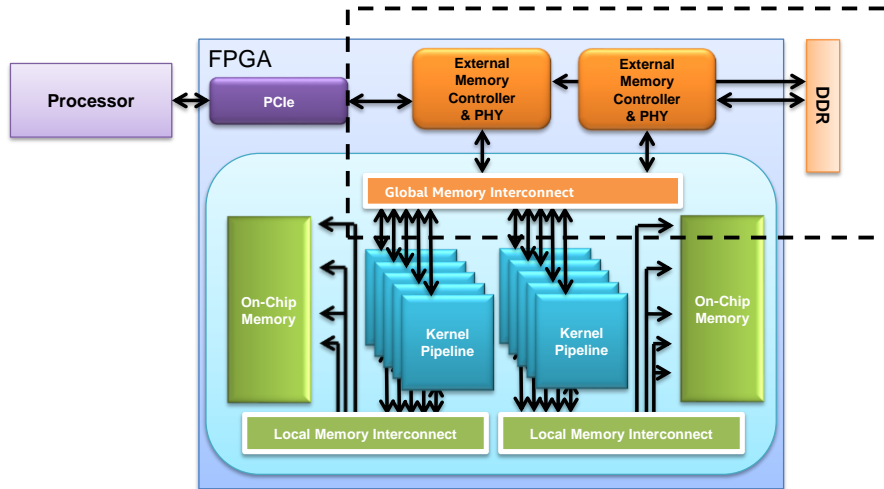
```
__kernel void Add(__global float* a,
                 __global float* b,
                 __global float* c)
{
    int i = get_global_id(0);
    c[i] = a[i] + b[i];
}
```

### Generally allocated on host as cl\_mem object

- Created with clCreateBuffer
- Data transferred with clRead/clWrite Buffer
- cl\_mem object assigned to global pointer argument in kernel



## OpenCL BSP Global Memory



## Compiler's View of Global Memory

Agnostic to the memory technology itself

- DDR, QDR, HMC, QPI

Only a few pertinent parameters (provided by BSP)

- How many interfaces
- Width of the bus
- Burst size (affinity for linear access)
- Latency
- Bandwidth

Exposed as Avalon Memory Mapped Slave

- Compiler builds datapath and interconnect to communicate to fixed IP
- BSP developer can create BSP variants for different memory configurations



# OpenCL Global Memory

## Global Memory Overview

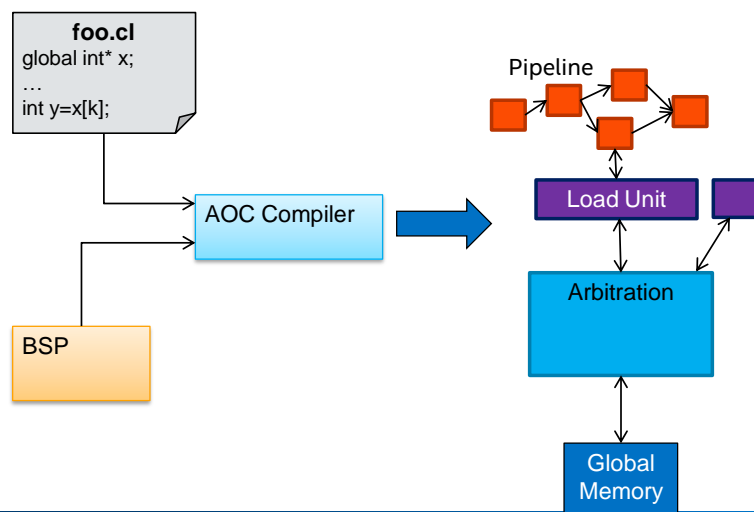
- Programmers view of global memory in OpenCL
- Compiler view of global memory

## Global Memory Architecture

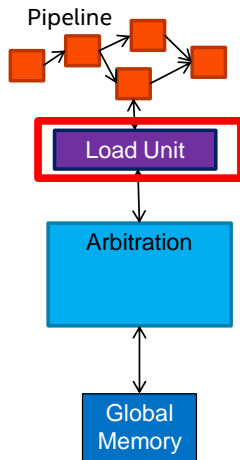
- Load-Store Units
- LSU Arbitration
- Const Cache

## Optimizations

## Compiler Generated Hardware



## Global Load/Store Unit (LSU)



### Width Adaptation

- User data (32-bit int) **to** memory word (512-bit DRAM word)
- Coalesces to avoid wasted bandwidth

### Burst coalescing

- Coalesces consecutive memory transactions into large burst transaction

## LSU Types (1)

### Pipelined

- Used for local memory

### Simple

- Passes transactions to interconnect from pipeline
- Used for loads/stores used very infrequently

### Burst-Coalesced

- Most common global memory LSU
- Specialized LSU to groups loads/stores into bursts
- Load LSU can cache/re-use data
  - Private caching is applied heuristically

### Streaming

- Simplified LSU used if compiler can determine access pattern is completely linear

## LSU Types (2)

### Semi-streaming

- Load LSU Specialized for nearly linear accesses
- Instantiated heuristically

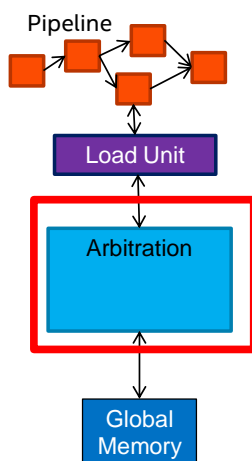
### Burst-non-aligned

- Wrapper around Burst-coalesced LSU, instantiated if access alignment may not be aligned to LSU width
- Tries to minimize overhead of non-aligned accesses by coalescing

### Wide LSU

- Wrapper around other LSU types, converts wider than global memory accesses into global memory sized accesses
- More area efficient than multiple memory width LSUs

## Arbitration



### Arbitrate to physical interfaces

- Tree interconnect (high bandwidth)
- Ring interconnect (high fmax)
  - Increase reliance on large bursts
- Arbitration type chosen base on # of LSUs

Distribute (load balance) across physical interfaces

## Const Cache

Constant buffer resides in global memory but accessed via on-chip cache shared by all work-groups

- Constant cache optimized for high cache hit performance

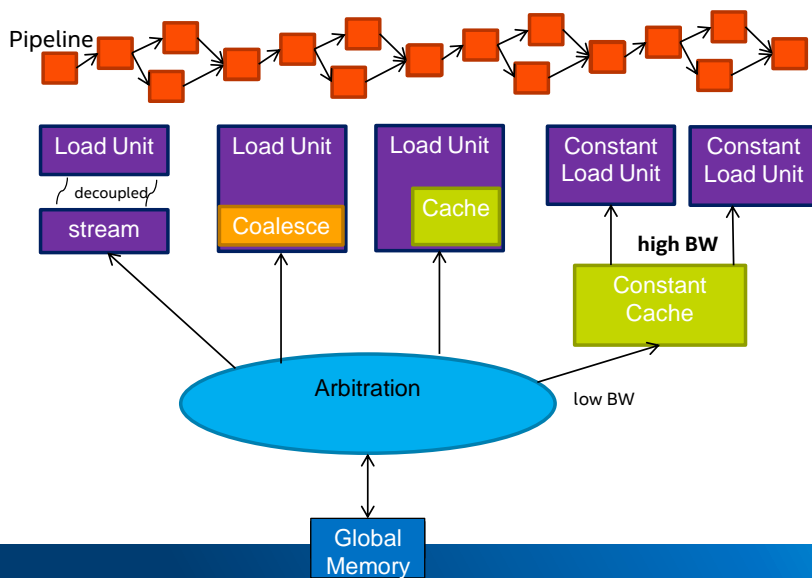
Use for read-only data that all work-groups access

- E.g. high-bandwidth table lookups

Constant cache default size is 16kB

- Uses on-chip RAM blocks that are shared with local memory

## Complete Picture



# OpenCL Global Memory

## Global Memory Overview

- Programmers view of global memory in OpenCL
- Compiler view of global memory

## Global Memory Architecture

- Load-Store Units
- LSU Arbitration
- Const Cache

## Optimizations



# Contiguous Memory Accesses

## Interleaved memory suitable for contiguous access

- Sequential accesses to global memory are the ideal access pattern to increase memory efficiency

## AOC will analyze access pattern of load and stores

Looks for sequential load and store operations for the entire kernel and directs kernel to access consecutive locations in global memory

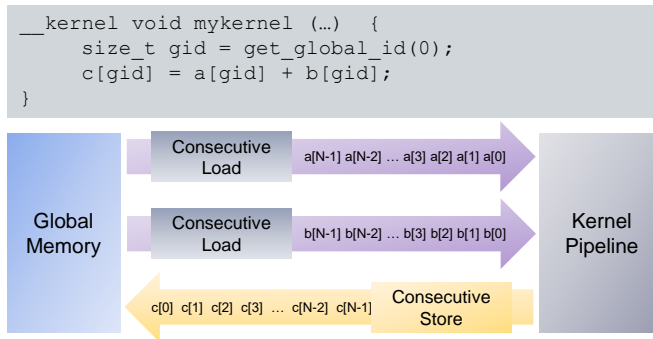
Leads to increased access speeds and reduced hardware resource needs



## Array Indexing and Contiguous Memory Accesses

Basing array index on the work-item global ID leads to efficient contiguous load and store operations

- Data sent to and received from the kernel pipeline as needed
- Computation and memory accesses can happen simultaneously



## Contiguous Memory Accesses (Tasks)

For Task kernels, memory should be indexed with an increasing loop counter for contiguous accesses

```

__kernel void mykernel (...) {
    for(int i = 0; i<BUFFER_SIZE; i++)
    {
        c[i] = a[i] + b[i];
    }
}

```

## Ensure 4-Byte Alignment for Data Structures

Struct alignments smaller than 4 bytes result in larger and slower hardware

```
typedef struct {
    char r,g,b,alpha;
} Pixel;
```

1-byte aligned  
struct

Force 4-byte alignment

Create a union of the structure  
and an integer

```
typedef struct {
    char r,g,b,alpha;
} Pixel_s;

typedef union {
    Pixel_s p;
    int not_used;
} Pixel;
```

or

Use aligned attribute with  
GCC or in the kernel

```
typedef struct {
    char r,g,b,alpha;
} __attribute__((aligned(4))) Pixel;
```

## Using \_\_constant Buffers

Constant cache default size is 16kB

- Uses on-chip RAM blocks that are shared with local memory

Manually specify cache size with AOC option

- Specify in bytes, AOC will round up to closest power of 2
- Have no effect if no kernels use \_\_constant address space

Constants suffer huge penalties for cache misses

- Use \_\_global const instead if constant argument can't fit in the cache

```
aoc --const-cache-bytes 32768 mykernel.cl
```

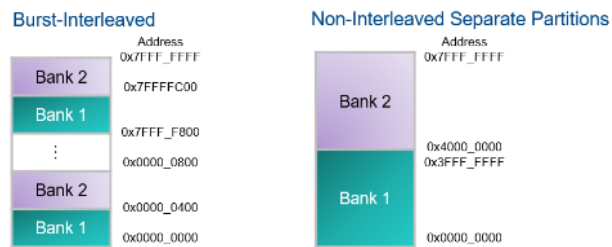
## Global Memory Banking Optimizations

Address space can be partitioned between banks or finely interleaved

Default: Configures global memory as burst-interleaved

- Best for sequential traffic
- Best for load balancing between memory banks

Burst-interleaving granularity determined by board vendor



## Manually Partitioned Global Memory

Turn off interleaving and assign data manually into memory banks

- Control memory bandwidth across a group of buffers

Advantages

- Better performance when accessing multiple pointers assigned to multiple banks
- Leads to more deterministic behavior
  - Designer knows the access pattern

In majority of use cases, manually partition of global memory leads to improved performance



## Manual Partitioning Mechanism

```
aoc --sw-dimm-partition <kernel file>.cl
```

- Configure memory banks as non-interleaved address spaces

Use `CL_MEM_BANK` flags to allocate memory buffer to one of the banks

- Allocate each buffer to a single memory bank only

|                                   |   |
|-----------------------------------|---|
| <code>CL_MEM_BANK_1_ALTERA</code> | Allocates to lowest available memory region                             |
| <code>CL_MEM_BANK_2_ALTERA</code> | Allocates to the second memory bank                                     |
| <code>CL_MEM_BANK_n_ALTERA</code> | Allocates to the n <sup>th</sup> bank, as long as the board supports it |

```
clCreateBuffer(context, CL_MEM_BANK_2_ALTERA |
                 CL_MEM_READ_WRITE, size, 0, 0);
```

## Avoiding False Memory Dependencies

Avoid using pointers that alias other pointers

- Only one pointer variable is used to access the contents

Use the `restrict` keyword whenever possible

- Specify to the compiler no aliasing for a pointer
- Prevents AOC from creating memory dependencies between non-conflicting load and store operations
- May cause functional errors if used with pointers that aliases other pointers

```
__kernel void my_kernel ( __global int * restrict A,
                          __global int * restrict B)
{
    ...
}
```

## Heterogeneous Memory

Some boards offer more than one type of global memory

- Eg. DDR and QDR
- One memory will be used by default

Memory location can be assigned per kernel argument

- `__attribute__((buffer_location("MEMORY_NAME")))`

Host will move memory to correct memory type upon kernel invocation

```
__kernel void foobar (
    __global uint *src,
    __global uint *dst,
    __global __attribute__((buffer_location("QDR"))) char* g_tables)
{ ...
```



# FLOATING POINT OPTIMIZATIONS

# Floating-Point Optimizations

Apply to `float` and `double` data types

AOC has the ability to create more efficient hardware for floating-point operations

Optimizations will cause small differences in floating-point results

- **Not** IEEE Standard for Floating-Point Arithmetic (IEEE 754-2008) compliant

AOC floating-point optimizations

- Needs to be manually enabled
- Tree Balancing
- Reducing Rounding Operations

Other optimizations

- Floating-point vs. fixed-point representations
- Use a device with hard floating point



# Arithmetic Order of Operation Rules

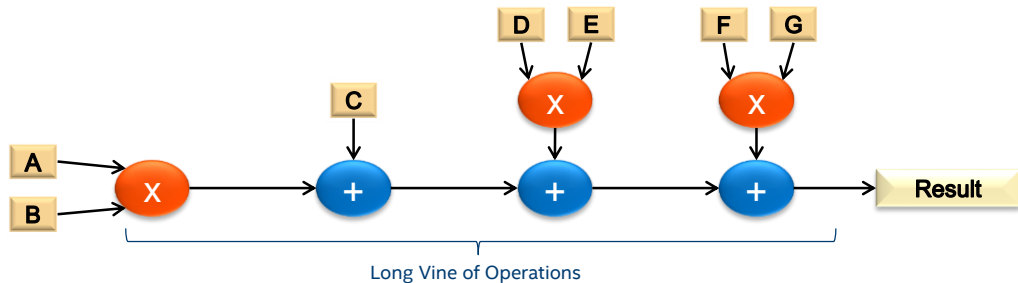
Strict order of operation rules apply in OpenCL

By default, AOC honors those rules

- May lead to long, unbalanced, slower, less-efficient floating-point operations

Example

```
Result = (((A * B) + C) + (D * E)) + (F * G)
```



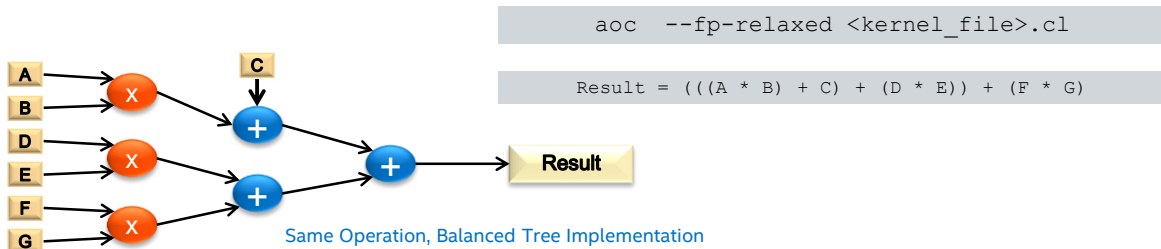
## Tree Balancing

Allow AOC to reorder arithmetic operations to convert into a tree pipeline structure

- Possibly affects the precision, not consistent with IEEE 754

Enable AOC tree balancing with `-fp-relaxed` option

- Design needs to tolerate the small differences in floating-point results



## Rounding Operations

For a series of floating-point operations, IEEE 754 require multiple rounding operation

Rounding can require significant amount of hardware resources

Fused floating-point operation

- Perform only one round at the end of the tree of the floating-point operations
- Leads to more accurate results
- Other processor architectures support certain fused instructions such as fused multiply and accumulate (FMAC)
- AOC can fuse any combination of floating-point operators

## Reducing Rounding Operations

AOC will not reduce rounding operations by default

Enable AOC rounding reduction with `--fpc` option

- Not IEEE 754 compliant
- Use when program can tolerate these differences in floating-point results

```
aoc --fpc <kernel_file>.cl
```

1. Removes floating-point rounding operations whenever possible  
Round floating-point operation only once at the end of the tree of operations  
Applies to \*, +, and -
2. Carry additional mantissa bits to maintain precision  
Carries additional bits through calculations, removes them at the end of the tree of operations
3. Changes rounding mode to round toward zero

## Floating-Point vs. Fixed-Point Representation

Fixed-point implementation always use less logic compared to floating-point representation

- Save hardware by converting operations to fixed-point

No variable width fixed-point representation in OpenCL

- Use `char` (8-bit), `short` (16-bit), `int` (32-bit), or `long` (64-bit)

If data resolution required is not one of the default supported ones (8, 16, 32, or 64), use appropriate masking operations to save hardware

- Savings are relatively small but can be significant if applied across a large design

## Fixed-Point Example

### 17-bit fixed-point data resolution needed

- Use 32-bit data type to store the value
- Avoid generating hardware for the upper 15 bits by using static bit masks
- Result: 17 bit addition implemented instead of 32-bit addition

```

__kernel fixed_point_add(__global const unsigned int * restrict a,
                        __global const unsigned int * restrict b,
                        __global unsigned int * restrict result)
{
    size_t gid = get_global_id(0);
    unsigned int temp;
    result[gid] = 0x3_FFFF & ((0x1_FFFF & a[gid]) + (0x1_FFFF & b[gid]));
}

```

32-bit data types

Mask away upper 14 bits of results

Mask away upper 15 bits of inputs

## Want to learn more?

### Developer Zone

- [OpenCL online demos](#)
- [OpenCL design examples](#)
- [Application Developer Partners](#)

### White papers, Publications and Optimization Tips for OpenCL

### Instructor-Led training

- [Parallel Computing with OpenCL Workshop by Altera](#)
- [Optimization of OpenCL for Altera FPGAs Training by Altera](#)
- [Building Custom Platforms](#)

### Online training

- [Introduction to Parallel Computing with OpenCL](#)
- [Writing OpenCL Programs for Altera FPGAs](#)
- [Running OpenCL on Altera FPGAs](#)
- [FPGAs vs GPUs](#)
- [Single-Threaded vs. Multi-Threaded Kernels](#)
- [Building Custom Platforms for Altera SDK for OpenCL](#)
- [OpenCL Optimization Techniques: Secure Hash Algorithm \(SHA-1\) Example](#)
- [OpenCL Optimization Techniques: Image Processing Algorithm Example](#)

