

DESIGN ENVIRONMENT FOR EXTREME-SCALE BIG DATA ANALYTICS ON HETEROGENEOUS PLATFORMS

CPS Summer School – Sep 22, 2022

Tutorial

Olympus: How to use HLS for building customized memory architectures

CHRISTIAN PILATO

EVEREST Scientific Coordinator Assistant Professor, Politecnico di Milano

christian.pilato@polimi.it

STEPHANIE SOLDAVINI

PhD Student, Politecnico di Milano

stephanie.soldavini@polimi.it

About Me – Christian Pilato

Tenure-Track Assistant Professor (RTD)

Website: http://pilato.faculty.polimi.it



2 Tutorial: Olympus - Sep 22, 2022



The EVEREST Project



Tutorial: Olympus - Sep 22, 2022 3

EVEREST Approach

Big data applications with heterogeneous data sources

Three use cases







What are the relevant requirements for data, languages and applications?

How to design data-driven policies for computation, communication, and storage?

How to create FPGA accelerators and associated binaries?

How to manage the system at runtime?

How to evaluate the results?

How to disseminate and exploit the results?

EVER55K

Open-source framework to support the optimization of selected workflow tasks



4 Tutorial: Olympus - Sep 22, 2022



EVEREST Partners



IBM Reseach Lab, Zurich (Switzerland) Project administration, prototype of the target system PI: Christoph Hagleitner



Università della Svizzera italiana (Switzerland) Data security requirements and protection techniques PI: Francesco Regazzoni



Centro Internazionale di Monitoraggio Ambientale (Italy) Weather prediction models PI: Antonio Parodi



Virtual Open Systems (France)

Virtualization techniques, runtime extensions to manage heterogeneous resources PI: Michele Paolino



Numtech (France)

Application for monitoring the air quality of industrial sites PI: Fabien Brocheton

Politecnico di Milano (Italy) Scientific coordination, high-level synthesis, flexible memory managers, autotuning PI: Christian Pilato

TU Dresden (Germany) Domain-specific extensions, code optimizations and variants PI: Jeronimo Castrillon



IT4Innovations (Czech Republic)



Exploitation leaders, large HPC infrastructure, workflow libraries PI: Katerina Slaninova

Duferco Energia (Italy) Application for prediction of renewable energies PI: Lorenzo Pittaluga



Sygic A/S (Slovakia) Application for intelligent transportation in smart cities PI: Radim Cmar





EVEREST Use Cases



6 Tutorial: Olympus - Sep 22, 2022



The Case of Computational Fluid Dynamics

Numerical simulations are becoming more and more popular for many applications

- Computational Fluid Dynamics requires to solve partial differential equations
- Inverse Helmholtz operator ("Helmholtz" for the friends) is parametric with respect to polynomial degree *p*

1	va	ar	iı	npι	ıt	S		:	[1	1	11]						
2	Va	ar	iı	npι	ιt	D		:	[1	1	11	11]					
3	va	ar	iı	ıpι	ıt	u		:	[1	1	11	11]					
4	Va	ar	٥ı	ıtı	put	. 1	J	:	[1	1	11	11]					
5	Va	ar	t					:	[1	1	11	11]					
6	va	ar	r					:	[1	1	11	11]					
7	t	=	S	#	S	#	S	#	u	•	[[1	L 6]	[3	7]	[5	8]]	
8	r	=	D	*	t												
9	v	=	S	#	S	#	S	#	t	•	[[(6]	[2	7]	[4	8]]	

Final result is obtained by "small" contributions on independent data

- CFD kernel is composed of three high-level tensor operators (two contractions and one Hadamard product) repeated millions of times – good for spatial parallelism
- Each operator requires $p^2 + 2 \cdot p^3$ (double) elements as input and produces p^3 (double) elements 21.74 KB + 10.40 KB per element when p = 11
- Six tensors (p³ elements) to store intermediate results additional 62.39 KB



EVEREST Target System



Tutorial: Olympus - Sep 22, 2022 8

EVEREST System Development Kit (SDK)



Different **input flows** starting from different **input languages**

Support for multiple target boards





Collection of interoperable and open-source tools to create hardware/software systems that can adapt to the target system, the application workflow, and the data characteristics

- Compilation framework based on MLIR to unify the input languages
- > High-level synthesis and hardware generation flow to automatically create optimized architectures
- Creation of hardware and software variants to match architecture features
- Use of state-of-the-art frameworks and commercial toolchains for FPGA synthesis



Stvm





Challenges for HPC Architectures (i)

Challenge 1: Input languages and frameworks

• Application designers are usually not FPGA experts and may use high-level framework that are not supported by current HLS tools – how to talk with them?

Challenge 2: CPU-Host Communication Cost

• FPGA logic requires the data on the board, but data transfers can be much more expensive than kernel computation – can we execute more than one point?

Challenge 3: Read/Write Burst Transactions

• We need to determine the proper size of the transactions to get the maximum performance – how to reorganize the data transfers and get the parameters?



Challenges for HPC Architectures (ii)

Challenge 4: Full Bandwidth Utilization

- AXI interfaces may be large (e.g., 256 bits on the Alveo) how to leverage them?
- HBM architectures have many channels how to parallelize data transfers?

Challenge 5: Data Allocation

 Data must be placed in memory to maximize its utilization but also to enable efficient data transfers/computation – custom data layouts?

Challenge 6: Synthesis-Related Issues

- FPGA devices are large but still not sufficient for hosting many kernels how to trade-off optimizations and parallel instances?
- FPGA architectures may be different separate platform agnostic and dependent parts?
- FPGA logic architectures are complex and may introduce performance degradation how to "guide" the synthesis process?

11Tutorial: Olympus - Sep 22, 2022https://tinyurl.com/CPS2022olympus



Hardware HPC (Memory) Architectures

What do we mean with **memory architecture**?

Every hardware module that is responsible to provide data to the accelerator kernels

Additional issues:

- BRAM resources are limited
 - Helmholtz operator requires >94 KB of local data
 - If local storage is not optimized, the number of parallel kernels can be limited
- Application-specific details can be used to optimize the data transfers
 - In Helmholtz, one of the tensors is constant over all elements how to match these details with platform characteristics?
 - Better to transfer data for a "batch" of elements and then execute them in series how many? again, limited storage



Hardware Compilation Flow



13 Tutorial: Olympus - Sep 22, 2022



From DSL to Bitstream – Focus on Memory



PLM Customization for Heterogeneous SoCs

High-Level Synthesis (HLS) to create the accelerator logic

• Definition of memory-related parameters (e.g., number of process interfaces)

Generation of specialized PLMs

- Technology-related optimizations
- Possibility of system-level optimizations across accelerators



Reuse What is not Used

Generally, we use one **PLM unit** (possibly composed of many banks) for each **data structure** (array)

Reuse the same memory IPs for several data structures

"Two data structures are compatible if they can be allocated to the same PLM unit (memory IPs)"

<u>A common case</u>: accelerator kernels never executed at the same time

- Possible only at system-level, when integrating the components
- Optimizations of accelerator logic and memory subsystem are independent



Memory Compatibility Graph (MCG)

Graph to represent the possibilities for optimizing the data structures

- Each node represents a data structure to be allocated, annotated with its data footprint (after data allocation)
- Each edge represents compatibility between the two data structures
- Can be automatically extracted from the MLIR-based compiler flow
 - Variant exploration to achieve the "best solutions"



- a) Address-space compatibility: the data structures are compatible and can use the same memory IPs
- b) Memory-interface compatibility: the ports are never accessed at the same time and the data structures can stay in the same memory IP



Clique Definition

"A clique is a subset of the vertices of the memory compatibility graph such that every two vertices are connected by an edge"



18 Tutorial: Olympus - Sep 22, 2022



PLM Controller Generation

19

A lightweight PLM controller is created for each compatibility set (clique) based on the bank configuration

- Accelerator logic is not aware of the actual memory organization
- Array offsets need to be translated into proper memory addresses





Custom logic with negligible overhead, especially when the number of banks and their size is a power of two

Tutorial: Olympus - Sep 22, 2022 https://tinyurl.com/CPS2022olympus



Creation of Parallel Architectures



K. F. A. Friebel, S. Soldavini, G. Hempel, C. Pilato, J. Castrillon. "From Domain-Specific Languages to Memory-Optimized Accelerators for Fluid Dynamics" HPCFPGA'21

20 Tutorial: Olympus - Sep 22, 2022



Olympus – Automated System-Level Integration

Complete hardware architecture generation flow from high-level specification



Olympus – System generation flow

Determines the system-level architectures based on:

- Algorithm parallelism
- Characteristics of the target platform(s)
- Interfaces of the modules (HLS tools)

Produces

22

- Synthesizable C++ code that includes:
 - Accelerators and PLM generated with HLS
 - Communication modules to match interfaces
 - Standard AXI interfaces to the system (either cloudFPGA SHELL or HBM channels)
 - May include "intelligent" policies to coordinate (or protect) data transfers
- System configuration file to create the overall architecture
 - Support for multiple computing units executing in parallel
 - Interfacing with Xilinx HLS and synthesis tools





EVERES

From MLIR to System Architecture

Automatic integration of memory optimizations for **high-performance data transfers**, such as:

- **Double buffering** to hide latency of host-FPGA data transfers
- Bus optimization (and data interleaving) for maximizing bandwidth (e.g., 256-bit AXI channels) algorithms for efficient data layout on the bus
- Dataflow execution model to enable kernel pipelining automatic (pre-HLS) code transformations



23 Tutorial: Olympus - Sep 22, 2022

https://tinyurl.com/CPS2022olympus

EVEREST

Results on HBM FPGA



24 Tutorial: Olympus - Sep 22, 2022



About Me – Stephanie Soldavini

- BS & MS in Computer Engineering from Rochester Institute of Technology (2014-2019)
- PhD Student at Politecnico di Milano (2020-Now)





Problem

- **Productivity**: Application designers usually do not have the necessary hardware design knowledge to create efficient hardware architectures
- **Performance**: Fine-tuned hardware descriptions are required to efficiently coordinate data transfers and execution



26 Tutorial: Olympus - Sep 22, 2022



Target Boards: Alveo Data Center Accelerator Cards with HBM



Resources	SLR0	SLR1	SLR2
HBM	32×256MB	-	-
DDR4	16GB	16GB	-
PLRAM	$2 \times 4 MB$	$2 \times 4 MB$	$2 \times 4 MB$
CLB LUT	369K	333K	367K
CLB Register	746K	675K	729K
Block RAM tile	507	468	512
UltraRAM	320	320	320
DSP	2,733	2,877	2,880





Olympus Flow Diagram



28 Tutorial: Olympus - Sep 22, 2022



Inputs: C HLS Kernel

- C code of application kernel to be accelerated
- Standard C array interfaces

```
input/kernel_body.cpp
```

```
1 #include "kernel_body.h"
     3 void kernel_body(double S[121], double D[1331], double u[1331], double v[1331],
                                                    double t[1331], double r[1331], double t1[1331],
                                                    double t3[1331], double t0[1331], double t2[1331])
     5
    6 {
     7
                                 for (int c1 = 0; c1 \le 10; c1 += 1)
                                                    for (int c2 = 0; c2 \le 10; c2 += 1)
     8
     9
                                                                       for (int c3 = 0; c3 \ll 10; c3 \leftrightarrow 1) {
10
                                                                                          // stmt0
                                                                                          t1[121 * c1 + 11 * c2 + c3] = 0;
11
                                                                                          for (int c8 = 0; c8 \le 10; c8 += 1)
12
13
                                                                                                             // stmt0
                                                                                                            t1[121 * c1 + 11 * c2 + c3] = t1[121 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 * c1 + 11 * c2 + c3] + S[11 *
14
15
                                                                     }
                                 for (int c1 = 0; c1 <= 10; c1 += 1)
16
17
                                                    for (int c2 = 0; c2 \ll 10; c2 \ll 1)
                                                                       for (int c3 = 0; c3 \ll 10; c3 \leftrightarrow 1) {
18
                                                                                          // stmt1
19
                                                                                          t0[121 * c1 + 11 * c2 + c3] = 0;
20
21
                                                                                          for (int c8 = 0; c8 \le 10; c8 += 1)
22
                                                                                                             // stmt1
                                                                                                            t0[121 * c1 + 11 * c2 + c3] = t0[121 * c1 + 11 * c2 + c3] + S[11 * c3] 
23
24
                                                                       }
                                 for (int c1 = 0; c1 <= 10; c1 += 1)
25
26
                                                    for (int c2 = 0; c2 \ll 10; c2 \leftrightarrow 1)
27
                                                                       for (int c3 = 0; c3 \ll 10; c3 \leftrightarrow 1) {
28
                                                                                          // stmt2
                                                                                          t[121 * c1 + 11 * c2 + c3] = 0;
29
                                                                                          for (int c8 = 0; c8 \le 10; c8 += 1)
30
31
                                                                                                             // stmt2
                                                                                                            t[121 * c1 + 11 * c2 + c3] = t[121 * c1 + 11 * c2 + c3] + S[11 * c1]
32
33
```



Inputs: JSON Kernel Specification

- Format derived from the JSON required for the Vitis RTL Blackbox flow¹
- Required info:
 - c_function_name
 - the name of the kernel function AND filename (cpp & h)
 - param_type
 - "mm" : memory mapped C arrays
 - "stream" : Xilinx hls::stream<>
 - c_parameters
 - Details on each port interface (must be in the same order as in the C source
 - c_name : Name of the array
 - c_type : Data type of a single element
 - c_port_direction : in, out, or inout
 - depth : Number of elements in the array

input/helmholtz.json

1 {				
2	"c_function_name": "kernel_body",	35	5	
3	"param_type": "mm",	36	Ľ	"c name": "r"
4	"c_parameters": [37		"c_type" "double"
5	{	38		"c port direction": "inout".
6	"c_name": "S",	39		"depth": 1331
7	<pre>"c_type": "double",</pre>	40	3.	
8	"c_port_direction": "in",	41	ŝ	
9	"depth": 121,	42		"c_name": "t1",
10	},	43		"c_type": "double",
11	{	44		<pre>"c_port_direction": "inout",</pre>
12	"c_name": "D",	45		"depth": 1331
13	<pre>"c_type": "double",</pre>	46	},	
14	<pre>"c_port_direction": "in",</pre>	47	{	
15	"depth": 1331,	48		"c_name": "t3",
16	},	49		"c_type": "double",
17	{	50		"c_port_direction": "inout",
18	"c_name": "u",	51		"depth": 1331
19	<pre>"c_type": "double",</pre>	52	},	
20	<pre>"c_port_direction": "in",</pre>	53	- {	
21	"depth": 1331,	54		"c_name": "t0",
22	},	55		"c_type": "double",
23	{	56		"c_port_direction": "inout",
24	"c_name": "v",	57		"depth": 1331
25	<pre>"c_type": "double",</pre>	58	3,	
26	"c_port_direction": "out",	59	1	"c name". "+2"
27	"depth": 1331,	61		C_name tz
28	},	62		"c port direction": "inout"
29	{	63		"denth" 1331
30	"c_name": "t",	64	3	ucpen . 1551
31	<pre>"c_type": "double",</pre>	65 1		
32	"c_port_direction": "inout",	66 }		
33	"depth": 1331	,,		
34	}.			



Makefile





Test Kernel: Inverse Helmholtz

- Tensor operator commonly used in computational fluid dynamics (CFD) applications
- Implemented in C as 7 loop nests of depth 3-4
 - input/kernel_body.cpp

$$t_{ijk,e} = \sum_{l=0}^{p} \sum_{m=0}^{p} \sum_{n=0}^{p} S_{li}^{T} \cdot S_{mj}^{T} \cdot S_{nk}^{T} \cdot u_{lmn,e} = (S \otimes S \otimes S \otimes u_{e})_{iljmkn}^{lmn}$$
(1a)

$$r_{ijk,e} = D_{ijk,e} \cdot t_{ijk,e}$$

$$v_{ijk,e} = \sum_{l=0}^{p} \sum_{m=0}^{p} \sum_{n=0}^{p} S_{li} \cdot S_{mj} \cdot S_{nk} \cdot r_{lmn,e} = (S \otimes S \otimes S \otimes r_e)_{limjnk} {}^{lmn}$$
(1c)



(1b)



Basic Implementation



make olympus

- Generate the code (default values in Makefile yield basic implementation)
- Sources: ~/alveo_tests/helmholtz_autogen/RB1_BW64_S0-student/krnl_helm/CLEAN/
 - src/ : kernel sources, CU.cpp is the generated Compute Unit wrapper
 - host/: host sources
 - HostImpl.gen.cpp: Implementation of driver functions (moving data to global mem, invoking CU)
 - HostSampleTop.gen.cpp: A sample test bench main file
- make chost chw POINTS=16 run
 - Build the sw_emu versions of the host and hardware and run with 16 iterations
- make hls TARGET=hw
 - Run HLS



Basic - Results

source /tools/Xilinx/Vitis/2021.1/settings64.sh

3rd-to-last line of "make hls" output: vitis_analyzer [path]/[kernel name].xo.compile_summary

Name	Latency (cycles)	Iteration Latency	Interval
✓ ● krnl_helm			
> 🔵 krnl_helm_Pipeline_1	124		124
> 🔵 krnl_helm_Pipeline_2	1334		1334
> 🔵 krnl_helm_Pipeline_3	1334		1334
v loopute_1	9744		9744
> 🔵 compute_1_Pipeline_VITIS_LOOP_7_1_VITIS_LOOP_8_2_VITIS_LOOP_9_3	1398		1398
> 🔵 compute_1_Pipeline_VITIS_LOOP_16_5_VITIS_LOOP_17_6_VITIS_LOOP_18_7	1398		1398
> 🔵 compute_1_Pipeline_VITIS_LOOP_25_9_VITIS_LOOP_26_10_VITIS_LOOP_27_11	1398		1398
> 🔵 compute_1_Pipeline_VITIS_LOOP_34_13_VITIS_LOOP_35_14_VITIS_LOOP_36_15	1343		1343
> 🔵 compute_1_Pipeline_VITIS_LOOP_39_16_VITIS_LOOP_40_17_VITIS_LOOP_41_18	1398		1398
> 🔵 compute_1_Pipeline_VITIS_LOOP_48_20_VITIS_LOOP_49_21_VITIS_LOOP_50_22	1398		1398
> 🔵 compute_1_Pipeline_VITIS_LOOP_57_24_VITIS_LOOP_58_25_VITIS_LOOP_59_26	1398		1398
> 🔵 krnl_helm_Pipeline_4	1334		1334
C VITIS_LOOP_41_1		14159	

BRAM	BRAM (%)	DSP	DSP (%)	FF	FF (%)	LUT	LUT (%)
10	~0	135	1	40078	1	27879	2



Ping Pong Buffers



- make olympus RING_BUF=2
 - RING_BUF sets the degree of ring buffering, 2 => ping pong
 - Sources: ~/alveo_tests/helmholtz_autogen/RB2_BW64_S0-student/krnl_helm/CLEAN/
- make chost chw RING_BUF=2 POINTS=16 run
 - Build the sw_emu versions of the host and hardware and run with 16 iterations
- make hls RING_BUF=2 TARGET=hw
 - Run HLS



Ping Pong Buffers - Results

Name	Latency (cycles)	Iteration Latency	Interval
∼ ● krnl_helm			
> 🔵 krnl_helm_Pipeline_4	124		124
> 🔵 krnl_helm_Pipeline_5	1334		1334
> 🔵 krnl_helm_Pipeline_6	1334		1334
> 🔵 krnl_helm_Pipeline_1	124		124
> 🔵 krnl_helm_Pipeline_2	1334		1334
> 🔵 krnl_helm_Pipeline_3	1334		1334
v longute_1	9744		9744
> compute_1_Pipeline_VITIS_LOOP_7_1_VITIS_LOOP_8_2_VITIS_LOOP_9_3	1398		1398
> compute_1_Pipeline_VITIS_LOOP_16_5_VITIS_LOOP_17_6_VITIS_LOOP_18_7	1398		1398
> ompute_1_Pipeline_VITIS_LOOP_25_9_VITIS_LOOP_26_10_VITIS_LOOP_27_11	1398		1398
> compute_1_Pipeline_VITIS_LOOP_34_13_VITIS_LOOP_35_14_VITIS_LOOP_36_15	1343		1343
> compute_1_Pipeline_VITIS_LOOP_39_16_VITIS_LOOP_40_17_VITIS_LOOP_41_18	1398		1398
> compute_1_Pipeline_VITIS_LOOP_48_20_VITIS_LOOP_49_21_VITIS_LOOP_50_22	1398		1398
> compute_1_Pipeline_VITIS_LOOP_57_24_VITIS_LOOP_58_25_VITIS_LOOP_59_26	1398		1398
> 🔵 krnl_helm_Pipeline_8	1334		1334
> 🔵 krnl_helm_Pipeline_7	1334		1334
C VITIS_LOOP_62_1		14159	

Drovious	BRAM	BRAM (%)	DSP	DSP (%)	FF	FF (%)	LUT	LUT (%)	Interface	Data Width (SW->HW)
Previous.	10	~0	135	1	40078	1	27879	2	m_axi_gmem0	64 -> 64
									m_axi_gmem1	64 -> 64
Now:	BRAM	BRAM (%)	DSP	DSP (%)	FF	FF (%)	LUT	LUT (%)	m_axi_gmem2	64 -> 64
	18	~0	135	1	42747	1	32108	2	m_axi_gmem3	64 -> 64



Wide Bus



- make olympus RING_BUF=2 BUS_WIDTH=256
 - BUS_WIDTH sets the width of the bus to global memory. 256/sizeof(double) => 4 "lanes"
 - Sources: ~/alveo_tests/helmholtz_autogen/RB2_BW256_S0-student/krnl_helm/CLEAN/
- make chost chw RING_BUF=2 BUS_WIDTH=256 POINTS=16 run
 - Build the sw_emu versions of the host and hardware and run with 16 iterations
- make hls RING_BUF=2 BUS_WIDTH=256 TARGET=hw
 - Run HLS





Wide Bus - Results

Name	Issue Type	Latency (cycles)	Iteration Latency	Interval
✓ ● krnl_helm				
> 🔵 read_data_1		3009		3009
🗸 🔵 compute		49659		49659
> compute_Pipeline_VITIS_LOOP_7_1_VITIS_LOOP_8_2_VITIS_LOOP_9_3	II Violation	8051		8051
> compute_Pipeline_VITIS_LOOP_16_5_VITIS_LOOP_17_6_VITIS_LOOP_18_7	II Violation	8051		8051
> compute_Pipeline_VITIS_LOOP_25_9_VITIS_LOOP_26_10_VITIS_LOOP_27_11	II Violation	8051		8051
> compute_Pipeline_VITIS_LOOP_34_13_VITIS_LOOP_35_14_VITIS_LOOP_36_15		1343		1343
> compute_Pipeline_VITIS_LOOP_39_16_VITIS_LOOP_40_17_VITIS_LOOP_41_18	II Violation	8050		8050
> compute_Pipeline_VITIS_LOOP_48_20_VITIS_LOOP_49_21_VITIS_LOOP_50_22	II Violation	8050		8050
> compute_Pipeline_VITIS_LOOP_57_24_VITIS_LOOP_58_25_VITIS_LOOP_59_26	II Violation	8050		8050
> 🔵 krnl_helm_Pipeline_VITIS_LOOP_133_2		1334		1334
> 🔵 krnl_helm_Pipeline_VITIS_LOOP_117_1		1334		1334
C VITIS_LOOP_188_1			54077	

Drovious	BRAM	BRAM (%)	DSP	DSP (%)	FF	FF (%)	LUT	LUT <mark>(%</mark>)
Plevious.	18	~0	135	1	42747	1	32108	2
Now:	BRAM	BRAM (%)	DSP	DSP (%)	FF	FF (%)	LUT	LUT (%)
NOW.	72	1	142	1	109117	4	66257	5

Interface	Data Width (SW->HW)				
m_axi_gmem0	256 -> 256				
m_axi_gmeml	256 -> 256				
m_axi_gmem2	256 -> 256				
m_axi_gmem3	256 -> 256				



Streaming (1 compute)



- make olympus RING_BUF=2 BUS_WIDTH=256 STREAMS=1
 - STREAMS=1 turns on the HLS dataflow pragma and uses the hls::stream data type
 - Sources: ~/alveo_tests/helmholtz_autogen/RB2_BW256_S1-student/krnl_helm/CLEAN/
- make chost chw RING_BUF=2 BUS_WIDTH=256 STREAMS=1 POINTS=16 run
 - Build the sw_emu versions of the host and hardware and run with 16 iterations
- make hls RING_BUF=2 BUS_WIDTH=256 STREAMS=1 TARGET=hw
 - Run HLS



Streaming (1 compute) - Results

Name	Latency (cycles)	Iteration Latency	Interval
✓ ● krnl_helm			
✓	14088		13877
> 🔵 read_data_l	3009		3009
> 🔵 compute_11	13876		13876
> Compute_12	13876		13876
> 🔵 compute_13	13876		13876
✓ ● compute_1	13876		13876
> 🔵 compute_1_Pipeline_VITIS_LOOP_162_1	123		123
> 🔵 compute_1_Pipeline_VITIS_LOOP_167_2	1333		1333
> 🔵 compute_1_Pipeline_VITIS_LOOP_172_3	1333		1333
> 🔵 compute_1_Pipeline_VITIS_LOOP_7_1_VITIS_LOOP_8_2_VITIS_LOOP_9_3	1398		1398
> compute_1_Pipeline_VITIS_LOOP_16_5_VITIS_LOOP_17_6_VITIS_LOOP_18_7	1398		1398
> 🔵 compute_1_Pipeline_VITIS_LOOP_25_9_VITIS_LOOP_26_10_VITIS_LOOP_27_11	1398		1398
> 🔵 compute_1_Pipeline_VITIS_LOOP_34_13_VITIS_LOOP_35_14_VITIS_LOOP_36_15	1343		1343
> 🔵 compute_1_Pipeline_VITIS_LOOP_39_16_VITIS_LOOP_40_17_VITIS_LOOP_41_18	1398		1398
> 🔵 compute_1_Pipeline_VITIS_LOOP_48_20_VITIS_LOOP_49_21_VITIS_LOOP_50_22	1398		1398
> 🔵 compute_1_Pipeline_VITIS_LOOP_57_24_VITIS_LOOP_58_25_VITIS_LOOP_59_26	1398		1398
> 🔵 compute_1_Pipeline_VITIS_LOOP_178_4	1333		1333
entry_proc	0		0
> 🔵 write_data_1	1407		1407
C VITIS_LOOP_229_1			

Previous "Compute" Latency: 49659 Previous Iteration Latency: 54077

Previous:	BRAM	BRAM (%)	DSP	DSP (%)	FF	FF (%)	LUT	LUT (%)
	72	1	142	1	109117	4	66257	5
Now:	BRAM	BRAM (%)	DSP	DSP (%)	FF	FF (%)	LUT	LUT (%)
	104	2	540	5	158002	6	105005	8

40



Streaming (1 compute) - Results





Streaming (multi compute)



- Multi compute is not controlled by Olympus, but using the INLINE pragma any submodules in the original HLS kernel will be used in the pipeline
- Edit the Makefile:
 - KERNEL_BODY ?= kernel_body_df
 - KERNEL_MODEL ?= kernel_body_sw
 - KERNEL_JSON ?= input/helmholtz_df.json

- Same commands as before:
 - make olympus RING_BUF=2 BUS_WIDTH=256 STREAMS=1
 - Sources: ~/alveo_tests/helmholtz_autogen/RB2_BW256_S1-student/krnl_helm/DF/
 - make chost chw RING_BUF=2 BUS_WIDTH=256 STREAMS=1 POINTS=16 run
 - make hls RING_BUF=2 BUS_WIDTH=256 STREAMS=1 TARGET=hw



Streaming (multi compute) - Results

Name	Latency (cycles)	Iteration Latency	Interval
✓ ● krnl_helm			
✓	3429		3010
> 🔵 read_data_1	3009		3009
> 🔵 gemm5	2861		2861
> 🔵 gemm	2861		2861
> 🔵 mmult8	1340		1340
> 🔵 mmult15	1340		1340
> 🔵 mmult22	1340		1340
> 🔵 mmult	1340		1340
> 🔵 gemm_inv9	2861		2861
> 🔵 gemm_inv_last	2861		2861
entry_proc	0		0
> 🔵 write_data_1	1407		1407

Previous Dataflow Interval: 13877

Provious	BRAM	BRAM (%)	DSP	DSP (%)	FF	FF (%)	LUT	LUT (%)
FIEVIOUS.	104	2	540	5	158002	6	105005	8
Now:	BRAM	BRAM (%)	DSP	DSP (%)	FF	FF (%)	LUT	LUT (%)
	144	3	2964	32	344814	13	279029	21



43 Tutorial: Olympus - Sep 22, 2022

C VITIS_LOOP_199_1

Streaming (multi compute) - Results







Conclusions

Data management optimizations are becoming the key for the creation of **efficient FPGA architectures** (... more than pure kernel optimizations)

HLS is now used not only to create accelerator kernels but also to generate the system-level architecture

• Portable solutions across multiple target platforms

Novel **HBM architectures** offer high bandwidth (that's why they are called *high-bandwidth memory* architectures... ③) but their design is complex:

- Necessary to match application requirements and technology characteristics
- We propose an MLIR-based compilation flow that directly interfaces with commercial HLS tools





ON HETEROGENEOUS PLATFORMS

Thank you!

CHRISTIAN PILATO

EVEREST Scientific Coordinator Assistant Professor, Politecnico di Milano

christian.pilato@polimi.it

STEPHANIE SOLDAVINI

PhD Student, Politecnico di Milano

stephanie.soldavini@polimi.it





This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 957269