# Webinar II

## HPP-Heterogeneous Processing Platform
### *Programming Models, Performance and Power Profiling for the HPP*

Yassine Hariri
Senior Engineer, Platform Design
Hariri@cmc.ca

Hugh W. Pollitt-Smith
Senior System Design Engineer
Pollitt-smith@cmc.ca

2-12-2015

# Agenda

- Overview
- emSYSCAN Development Systems Update
- HPP: Heterogeneous processing platform
- HPP GPU
  - GPU Programming
    - Libraries, OpenACC and Programming langages
  - CUDA Development Using NVIDIA Nsight, Eclipse Edition
    - Project Management, Edit, Build, Debug and Profile
  - Power profiling using nvidia-smi
  - Live demo
- HPP FPGA
  - OpenCL for FPGA
  - The AOCL FPGA Programming Flow
  - Power and performance profiling
  - Live demo
- HPP schedule
- Heterogeneous processing workshop

# EMSYSCAN DEVELOPMENT SYSTEMS UPDATE

# Embedded Systems *Canada (emSYSCAN)*
## $54M investment in Canada's National Design Network
## 37+ universities, 250+ faculty, 5 years

Microsystems
Rapid-Prototyping,
Characterization
and Integration Labs

4 Universities:
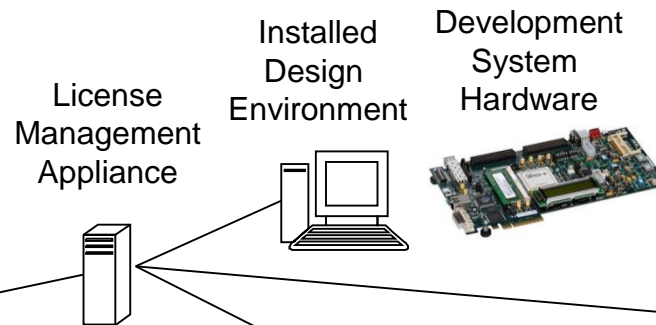• UBC
• U Manitoba
• Queen's
• École Polytechnique

Multi-Technology Design Environment
• System architecture exploration
• Multi-technology simulation
• Design of custom devices for manufacturing

Development Systems
• System validation and proof-of-concept demonstration

License Management Appliance

Installed Design Environment

Development System Hardware

Real-Time Embedded Software Lab

University of Waterloo

• Design, analysis, debug of real-time software on next-generation processor systems

Common, Shared Platforms
Interconnected Community of Users
Knowledge Repository
Centralized Management & Operations

License Management Server (LMS)

# Development Systems for Proof of Concept



Images courtesy of National Instruments, Xilinx, BEECube,, NVIDIA

# National Research Platform: Enriched Projects; Results Sooner

- <u>Common</u> set of programmable research platforms with proof of concept features

- <u>Pooled</u> equipment timeshared among users

- <u>Sharing</u> of knowledge on equipment usage

- <u>Adaptive</u> over time in terms of equipment quantities and equipment features

- <u>Large</u> community of users, institutions

- <u>Leveraged</u> industrial partners (e.g., STMicro.)

> National project scope and sizeable outcomes enabled by centralized project implementation and management by CMC Microsystems

# Installation and usage

- Shared access systems can be accessed at no charge but require Designer level subscription
  - Subscription provides access to support, tools, reference designs, forums, workshops, travel, select/swap, training, additional discounts
- Systems delivered on site, remote access
- Designated Development System coordinator(s) at each site
  - Communicate institutional needs for purchase specifications
  - Local advocate, information source
  - Encourage participation in National Project

# emSYSCAN Development Systems delivered (Gen1)

- Embedded Systems Platform:
  - Xilinx ML605, Altera DE4-530
- Advanced Processing Platform
  - BEEcube BEE3, BEE4, miniBEE
- Software-Defined Radio Platform
  - BEEcube miniBEE, RF daughtercard
- Simulation Acceleration Platform
  - Nallatech P385-D5 (Altera Stratix V, OpenCL)
- Multiprocessor Array Platform
  - NVIDIA Tesla K20 GPU
  - Intel Xeon Phi
- Microsystems Integration Platform
  - National Instruments PXI-based, FPGA, MEMS, microfluidics, RF, photonics features

# NDN Development Systems Community

https://community.cmc.ca/community/development-systems



**Development Systems National Catalog**

created by hugh on Jun 20, 2012 2:30 PM, last modified by hugh on Jan 18, 2013 12:34 PM

The following Development Systems have been delivered to the National Design Network (NDN) for shared access. The designated Coordinator/Contact can provide additional details on availability and how to access:

| System Product | Location/University | Quantity | Coordinator/Contact |
|---|---|---|---|
| **BEEcube BEE3 (Advanced Processing Platform)** | | | |
| | CMC Microsystems (online access) | 2 | ✉ Hugh Pollitt-Smith |
| | McMaster University | 1 | Dr. Nicola Nicolici |
| | University of Guelph | 1 | Dr. Stefano Gregori |
| | McGill University | 1 | Dr. Zeljko Zilic |
| | University of New Brunswick | 1 | Dr. Kenneth Kent |
| | University of Saskatchewan | 1 | Dr. Seok-Bum Ko |
| | Université du Québec à Chicoutimi | 1 | Dr. Hung-Tien Bui |
| | Université du Québec à Trois-Rivières | 1 | Dr. Adel Omar Dahmane |
| | Université du Québec à Outaouais | 1 | Dr. Ahmed Lakhssassi |
| | University of Windsor | 1 | Dr. Rashid Rashidzadeh |

**Actions**

**More Like This**

- Introductory setup instructions for the BEEcube miniBEE
- Advanced Processing Platform (BEE3) Design Environment
- Development Systems Inventory
- Advanced Processing Platform
- Setup Information for the HP Z400 workstations (ML605, DE4-530, BEE3)

# Upcoming emSYSCAN Development Systems deployments

- Embedded Systems Platform
  - Xilinx Virtex-7, Ultrascale, Zynq options
  - Altera Arria 10 and Arria 10 SoC options
  - Shipping Q1 2016
- Advanced Processing Platform
  - RFP currently in evaluation
  - Shipping Q1/Q2 2016
- Software-defined Radio
  - BEEcube nano/megaBEE (2x2 up to 16x16 MIMO options)
  - Shipping Q1 2016

## 32 Institutions

Design, compute, store on local resources; secure download

Secure, remote access to platforms, compute infrastructure (thin client)

## Benefiting Canadians

- Information and Communications Technologies
- Healthcare
- Transportation
- Energy

- Manufacturing
- Security

## ADEPT:
## Advanced Design Platform Technology

### Design Platforms

- Multi-technology Interposer
- Sensor/ Actuator
- Heterogeneous Embedded
- Silicon Photonics
- User Platform 1
  ⋮

Intellectual Property Blocks & Physical Design Kits

Advanced Design Methods

Computer-Aided Design Tools

### Fabrication Process Repository

- Equipment database
- Recipes
- TCAD

### Access Infrastructure

License Management System

Compute Servers (Compute Canada)

**User accounts, storage**   Accelerators

### Cybersecurity Testbed

Observation, analysis, delivery

## Fabrication Laboratories

- Lab Capabilities
- Recipe Development
- Prototyping

## Vendors & Partners

- CAD tools
- Intellectual Property Blocks & Physical Design Kits
- Design Methods
- Multi-project wafer services & scale-up manufacturing

⬭ Existing infrastructure

🔒 Secure links

---

*Canada's National Design Network – ADEPT Management & Operations*
Includes software procurement, configuration, installation and delivery. Access and utilization management, engineering/technical support. Cybersecurity installations, secure testbed assistance and demonstrations.      Train-the-trainer events. Advisory Group coordination. Governance, reporting, legal and financial administration.

# HPP Distribution

- Based on Development Systems Coordinator consultations in April 2014:
  - Generation 1 (2014/15): 18 systems
    - USask, UQTR, Outaouais, McGill, York, Windsor, Waterloo, Western, Ottawa, Ryerson, RMC, Victoria
  - Generation 2 (2016/17): 12 systems
    - Memorial, Guelph, McMaster, Toronto, Polytechnique, UQTR, Outaouais

# HPP: HETEROGENEOUS PROCESSING PLATFORM

# HPP main components

- The HPP workstation integrates the following main components:
    - Dual core Intel Xeon E5-2620 V3
    - NVidia GPU (Tesla K20)
    - FPGA board (Nallatech P385-A72).
    - Xeon Phi 7120A
- Key Platform Benefits
    - Customizability: Select the right mix of accelerators for your application
    - Greater flexibility for HW/SW exploration
    - Scalability: Create one node and scale up by adding more nodes
    - Fast automated setup and configuration
    - Faster path to commercialization
    - Technical support and training from CMC Microsystems

# HPP fully installed system

# HPP Pre-Installed Software Components

These software components required by the HPP are pre-installed on the workstation:

- RedHat Enterprise Linux 6.6 (Kernel version: 2.6.32-504.el6.x86_64)

- Java Runtime Environment (JRE)

- gcc compiler and toolchain

| Accelerator | Software and tools |
|---|---|
| GPU Tesla K20 | • NVIDIA CUDA 7 Toolkit |
| Nallatech FPGA P385 | • Altera Quartus 15.0, Altera SDK for OpenCL 1.0<br><br>• Nallatech FPGA P385 Board Support Package |
| Xeon Phi 7120a | • Intel Manycore Platform Software Stack (MPSS) 3.5.1 for Linux<br><br>• Intel Parallel Studio XE 2015, Professional Studio for C++, Linux version |

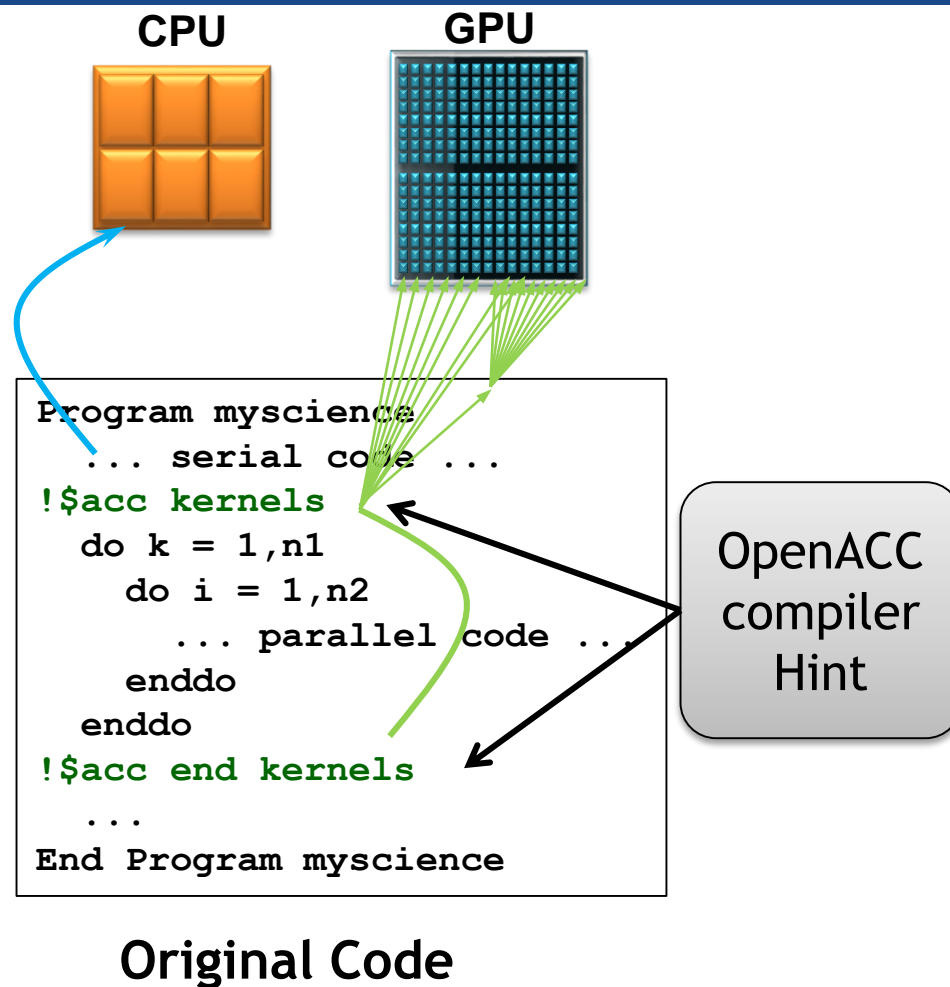# HPP: GPU

# Applications Acceleration

Applications

Libraries

OpenACC Directives

Programming Languages

# Applications Acceleration

Applications

Libraries

OpenACC Directives

Programming Languages

# Libraries

- **Ease of use:** Deep knowledge of GPU programming is not required

- **"Drop-in":** Standard APIs, minimal code changes

- **Quality:** High-quality implementations

- **Performance:** NVIDIA libraries are highly optimized

# CUDA Libraries Ecosystem

- CUDA Tools and Ecosystem described in detail on NVIDIA Developer Zone:

  developer.nvidia.com/cuda-tools-ecosystem

# Applications Acceleration

Applications

Libraries

OpenACC Directives

Programming Languages

# OpenACC Directives



CPU      GPU

```
Program myscience
  ... serial code ...
!$acc kernels
  do k = 1,n1
    do i = 1,n2
      ... parallel code ...
    enddo
  enddo
!$acc end kernels
  ...
End Program myscience
```

OpenACC compiler Hint

**Original Code**

# Applications Acceleration

Applications

Libraries

OpenACC Directives

Programming Languages

# GPU Programming Languages

- Numerical analytics: MATLAB, Mathematica, LabVIEW
- Fortran: OpenACC, CUDA Fortran
- C: OpenACC, CUDA C
- C++: Thrust, CUDA C++
- Python: PyCUDA, Copperhead

# HPP GPU :
# CUDA DEVELOPMENT USING NVIDIA NSIGHT, ECLIPSE EDITION

# NVIDIA® Nsight™ Eclipse Edition

- **CUDA Integrated Development Environment**
- **Project Management**
- **Edit**
- **Build**
- **Debug**
- **Profile**

# Powered By Eclipse

- **Extensible via robust selection of open-source and commercial plugins**

- **Revision control: CVS, SVN, Git, Perforce, …**

- **Issue tracking**

- **…**

- **Strong cross-platform support**

- **Nsight Eclipse Edition available for Linux and Mac OSX**

# Included In CUDA Toolkit

# NVIDIA® Nsight™ Eclipse Edition

- **CUDA Integrated Development Environment**
  - **Project Management**
  - **Edit**
  - **Build**
  - **Debug**
  - **Profile**

# Nsight Project Support

- **CUDA C / C++**
- **Project Types**
  - **Executable**
  - **Shared Library**
  - **Static Library**

- **New vs. Existing**
  - **New project, build managed by Nsight**
  - **Existing project, Nsight can use your Makefile**

# Creating A New CUDA Project

# Nsight main window after creating a new project

# Edit

- **CUDA Integrated Development Environment**
    - **Project Management**
    - **Edit**
    - **Build**
    - **Debug**
    - **Profile**

# CUDA Editor

- **CUDA-aware syntax highlighting**
- **Host / Device code highlighting**
- **Smart code assist**
- **As-you-type error detection**
- **CUDA API documentation pop-ups**
- **Automatic code refactoring**

```
128      // Launch the Vector Add CUDA Kernel
129      int threadsPerBlock = 256;
130      int blocksPerGrid =(numElements + threadsPerBlock - 1) / threadsPerBlock;
131      printf("CUDA kernel launch with %d blocks of %d threads\n", blocksPerGrid, threadsPerBlock);
132      vectorAdd<<<blocksPerGrid, threadsPerBlock>>>(d_A, d_B, d_C, numElements);
133      /**
134       * CUDA Kernel Device code
135       *
136       * Computes the vector addition of A and B into C. The 3 vectors have the same
137       * number of elements numElements.                                                 etErrorString(err));
138       */
139      __global__ void
140      vectorAdd(const float *A, const float *B, float *C, int numElements)
141      {
142          int i = blockDim.x * blockIdx.x + threadIdx.x;
143
144          if (i < numElements)
145
146
```

# Edit

- **CUDA Integrated Development Environment**
  - **Project Management**
  - **Edit**
  - **Build**
  - **Debug**
  - **Profile**

# CUDA Builder

- **Full CUDA toolchain support**
  - **All nvcc features**
  - **Debug, release, and custom build configurations**
- **Dependent project support**
  - **Static libraries**
  - **Shared libraries**
  - **Manages all build dependencies**
- **Source-correlated error reporting**

# Build Error Reporting

# Run Application

# Edit

- **CUDA Integrated Development Environment**
  - **Project Management**
  - **Edit**
  - **Build**
  - **Debug**
  - **Profile**
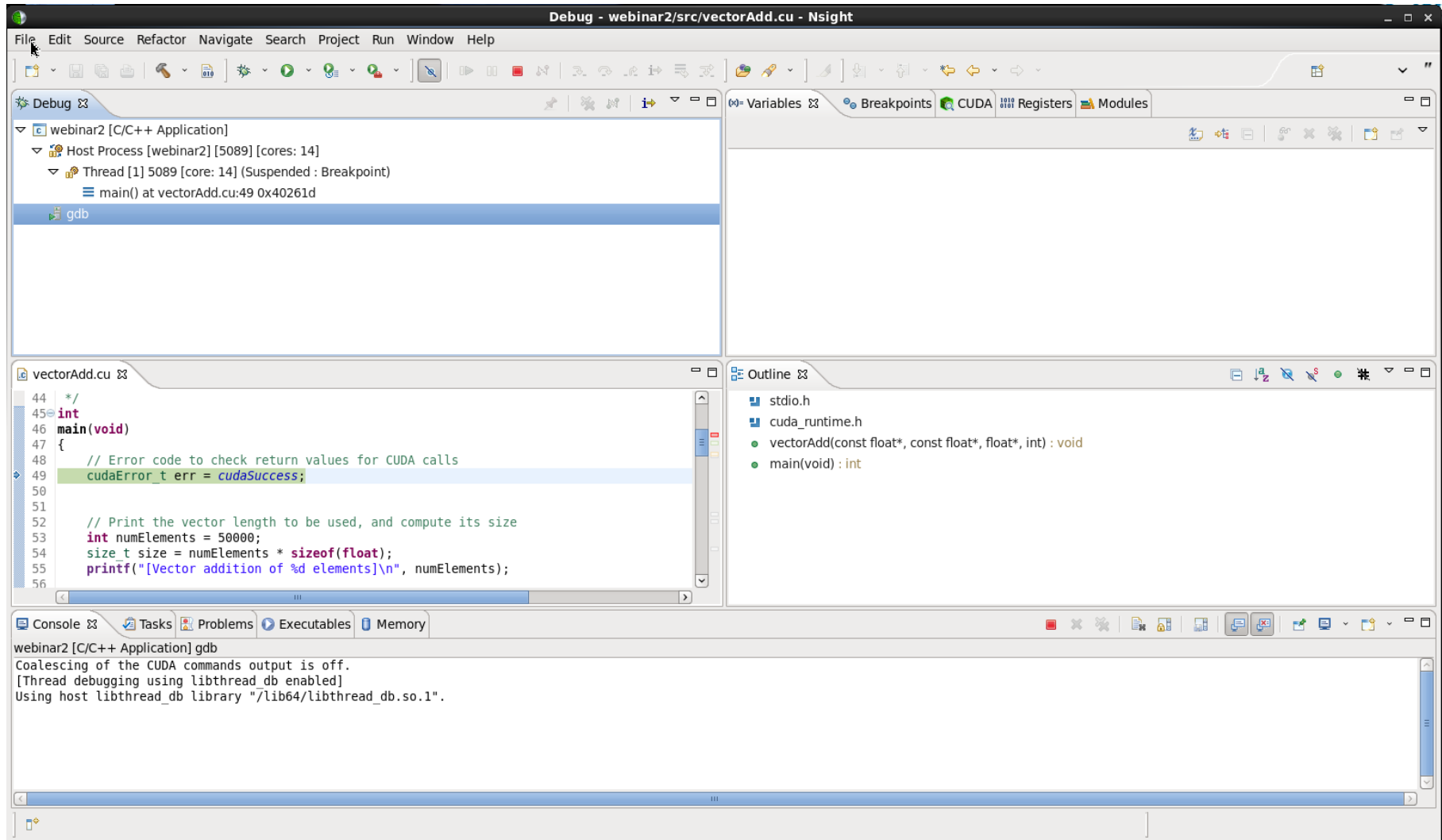
# CUDA Debugger

- **Unified CPU / GPU Debugging**
  - **Simultaneous visibility into both CPU and GPU state**
  - **Multi-GPU support**
- **Full GPU debugging**
  - **Set kernel breakpoints**
  - **Single-step, run until, etc.**
  - **View variables, registers, and expression values across multiple GPU threads at the same time**
  - **Examine thread, warp, block state**
  - **Source and assembly level debugging**

# Add a break point in the code

# GPU / CPU Threads, Call Stacks

# Stepping

# GPU / CPU Threads, Call Stacks

# Edit

- **CUDA Integrated Development Environment**
    - **Project Management**
    - **Edit**
    - **Build**
    - **Debug**
    - **Profile**

# CUDA Profiler

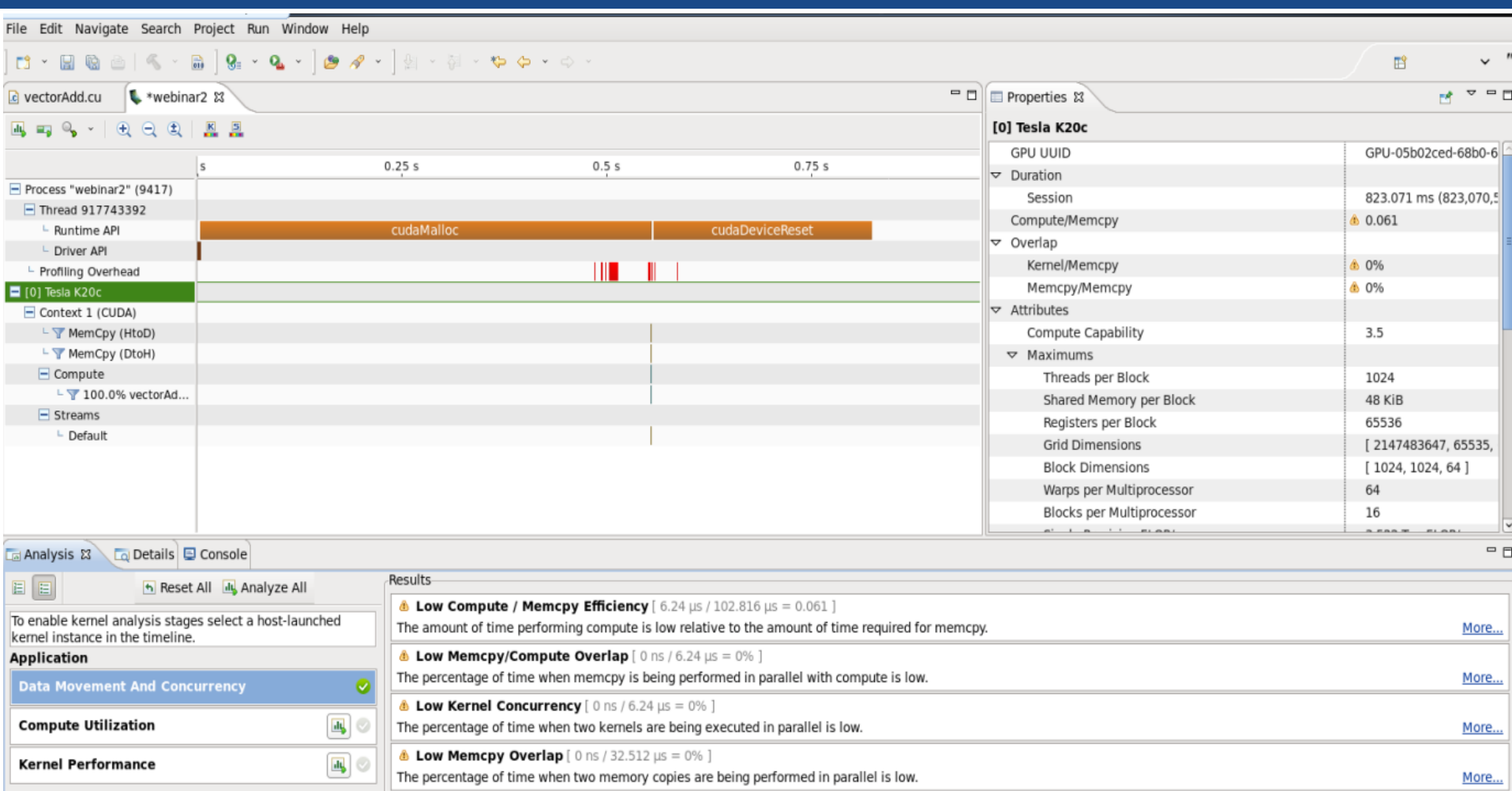- **Unified GPU / CPU profiler**
- **Visualize GPU / CPU interactions Identify GPU utilization and efficiency bottlenecks**
- **View low-level counters and metrics**
- **Multi-GPU support**
- **Automated application analysis identifies optimization opportunities**
- **Online documentation gives direction on how to exploit opportunities to get performance improvement**

# Unified GPU / CPU Timeline

# Analysis Documentation



Results

⚠ **Low Compute / Memcpy Efficiency** [ 6.24 µs / 102.816 µs = 0.061 ]
The amount of time performing compute is low relative to the amount of time required for memcpy.                    More...

⚠ **Low Memcpy/Compute Overlap** [ 0 ns / 6.24 µs = 0% ]
The percentage of time when memcpy is being performed in parallel with compute is low.                    More...

⚠ **Low Kernel Concurrency** [ 0 ns / 6.24 µs = 0% ]
The percentage of time when two kernels are being executed in parallel is low.                    More...

⚠ **Low Memcpy Overlap** [ 0 ns / 32.512 µs = 0% ]
The percentage of time when two memory copies are being performed in parallel is low.                    More...

## 9.1. Data Transfer Between Host and Device

The peak theoretical bandwidth between the device memory and the GPU is much higher (177.6 GB/s on the NVIDIA Tesla M2090, for example) than the peak theoretical bandwidth between host memory and device memory (8 GB/s on the PCIe x16 Gen2). Hence, for best overall application performance, it is important to minimize data transfer between the host and the device, even if that means running kernels on the GPU that do not demonstrate any speedup compared with running them on the host CPU.

**Note: High Priority:** Minimize data transfer between the host and the device, even if it means running some kernels on the device that do not show performance gains when compared with running them on the host CPU.

Intermediate data structures should be created in device memory, operated on by the device, and destroyed without ever being mapped by the host or copied to host memory.

Also, because of the overhead associated with each transfer, batching many small transfers into one larger transfer performs significantly better than making each transfer separately, even if doing so requires packing non-contiguous regions of memory into a contiguous buffer and then unpacking after the transfer.

Finally, higher bandwidth between the host and the device is achieved when using *page-locked* (or *pinned*) memory, as discussed in the *CUDA C Programming Guide* and the Pinned Memory section of this document.

### 9.1.1. Pinned Memory

Page-locked or pinned memory transfers attain the highest bandwidth between the host and the device. On PCIe x16 Gen2 cards, for example, pinned memory can attain roughly 6GB/s transfer rates.

Pinned memory is allocated using the cudaHostAlloc() functions in the Runtime API. The bandwidthTest CUDA Sample shows how to use these functions as well as how to measure memory transfer performance.

For regions of system memory that have already been pre-allocated, cudaHostRegister() can be used to pin the memory on-the-fly without the need to allocate a separate buffer and copy the data into it.

Pinned memory should not be overused. Excessive use can reduce overall system performance because pinned memory is a scarce resource, but how much is too much is difficult to know in advance. Furthermore, the pinning of system memory is a heavyweight operation compared to most normal system memory allocations, so as with all optimizations, test the application and the systems it runs on for optimal performance parameters.

# Power profiling using nvidia-smi

- **nvidia-smi -q -d POWER -i 0 -l 1 -f out.log**
    - ➤ **−q, −−query** Display GPU or Unit info
    - ➤ **−d TYPE, −−display=TYPE** Display only selected information: MEMORY, UTILIZATION, ECC, TEMPERATURE, POWER, CLOCK, COMPUTE…
    - ➤ **−i, −−id=ID** Display data for a single specified GPU or Unit.
    - ➤ **−f FILE, −−filename=FILE** Redirect query output to the specified file in place of the default stdout. The specified file will be overwritten.

Power Consumption GPU (Watts)

# Live Demo

# HPP: FPGA

# FPGA pains

- Programmability is an issue
  - Hardware Description Languages (HDLs) are complexe
  - At Register transfer level (RTL) abstraction
  - Designer needs to :
    - Create Custom Memory Hierarchies
    - Manage Communication with PC
    - Create PC side SW that plays nice with all this

# OpenCL for FPGA

- Unified programming model
    - More accessible to the software developpers
    - Host CPU-FPGA communication
    - C based programming language
    - Memory Hierarchy auto generated
    - Potential to integrate existing VHDL/Verilog IP

# OpenCL for FPGA



White Paper: Implementing FPGA Design with the OpenCL Standard (Figure 4)

# The AOCL FPGA Programming Flow



Altera SDK for OpenCL Programming Guide

# AOCL design flow steps

1. **Intermediate compilation (aoc –c [-g] *&lt;your_kernel_filename&gt;*.cl)**
   – Checks for syntatic errors
   – Generates a **.aoco** file without building the hardware configuration file
   – Generate estimated resource usage summary *&lt;your_kernel_filename&gt;*.**log**

2. **Emulation (aoc -g &lt;your_kernel_filename&gt;.cl)**
   – The AOCL Emulator generates a **.aocx** file that executes on x86-64 Windows or Linux host
   – Assess the functionality of your OpenCL kernel

3. **Profiling (aoc --profile &lt;your_kernel_filename&gt;.cl)**
   – aocl report *&lt;your_kernel_filename&gt;*.aocx profile.mon
   – Instruct the Altera Offline Compiler to instrument performance counters in the Verilog code in the **.aocx** file
   – During execution, the performance counters collect performance information which you can then review in the Profiler GUI.

4. **Full deployment**
   – Execute the .aocx file on the FPGA

# OpenCL Overview



White Paper: Implementing FPGA Design with the OpenCL Standard (Page: 7)

# Nallatech 385 Hardware Overview



Fan power connector · CPLD header · DDR3 SDRAM banks · Power Supply Inductors · 2 RGB FPGA LEDs · Dual SFP+ cage and light pipes · 4 bi-colour FPGA LEDs · Altera Stratix-V user FPGA · Altera MAX-II CPLD · 1Gb Flash

- The key features of the 385 include:
    - PCI Express form factor
        - 8-lane PCI Express up to 3.0 host interface1
    - FPGA options
        - Altera Stratix-V 5SGXMA7H2F35C2N
        - Altera Stratix-V 5SGSMD5H2F35C2N
    - Two 10G LAN/WAN/FC Ethernet channels accessed via two SFP+ ports2
    - Two banks of SDRAM memory, each bank with 4GByte, x72, DDR3 SDRAM running at 1600 MT/s

# BSP and Altera SDK for OpenCL

- Four layers of the Altera Software Development Kit (SDK) for OpenCL :

| Runtime (OpenCL API) |
| HAL for memory transfers and kernel launches |
| MMD layer for raw read and write operations |
| Kernel mode driver for accessing communication medium |

| Board Hardware |

- The Nallatech OpenCL BSPs provide a memory-mapped device (MMD) layer necessary for communication with the accelerator board and the lower level kernel mode driver.

- All other upper software layers are provided by the Altera SDK for OpenCL installation.

# Linux Vector Addition Walkthrough

- Step 1 Compile
  - aoc –v --board p395_hpc_ab vectorAdd.cl
- Step 2 Build the Host code
  - make
- Step 3 Execute
  - ./vector_add

- Generated files:
  - vectorAdd.log – kernel compilation with estimated resource usage and Qsys generation.
  - quartus_sh_compile.log – Quartus tools build log for generating the actual FPGA hardware aocx file.
  - acl_quartus_report.txt – Final results for the generated design including final resource usage figures.

# Power consumption reading

- Note that this requires that you have already opened the card using the normal opencl SDK routines (i.e. in init_opencl in this case).

- **aocl_mmd_card_info**("acl0", AOCL_MMD_POWER, sizeof(float),(void*) &power, &returnedSize);

- **printf**("before run Power = %f W\n", power);

- **aocl_mmd_card_info**("acl0", AOCL_MMD_BOARD_UNIQUE_ID, sizeof(int), (void*) &uniqueId, &returnedSize);

- **printf**("Board Unique ID = %d\n", uniqueId);

```
[root@HPPPrototype bin]# ./matrix_mult
Matrix sizes:
  A: 2048 x 1024
  B: 1024 x 1024
  C: 2048 x 1024

MMD ERROR: MMD must be opened before calling aocl_mmd_card_infoInitializing OpenCL
Platform: Altera SDK for OpenCL
Using 1 device(s)
  p385_hpc_a7 : PCIe385n
Using AOCX: matrixMult.aocx
Reprogramming device with handle 1
Generating input matrices
before run Power = 17.894531 W
Board Unique ID = 7802811
About to call run
Launching for device 0 (global size: 1024, 2048)
During run Power = 18.947571 W

Time: 38.713 ms
Kernel time (device 0): 38.594 ms

Throughput: 110.94 GFLOPS

Computing reference output
Verifying
Verification: PASS
after run Power = 18.063171 W
```

Live Demo

# Project status

Status (12 universities selected 18 systems G1, 7 universities selected 8 G2)

o Assembled, cloned, tested and shipped 18/18 units

o (1) Quick start guide : Heterogeneous Parallel Platform (HPP)

o Introduction to the HPP-Heterogeneous Parallel Platform: A combination of Multicores, GPUs, FPGAs and Many-cores accelerators (August 26th)

In progress

o (2) User Guide: Performance and Power profiling for the HPP

o Programming models, performance and power profiling for the HPP-Heterogeneous Parallel Platform (December 2nd)

Next (Webinars series for the HPP)

o Computer vision using OpenCV/OpenCL targeting the HPP- Heterogeneous Processing Platform (January 13th)

# Workshop on Heterogeneous Computing Platforms (Scope)

Location: Toronto

Date: TBD (February 15 - February 19, 2016)

Objective

- Bring researchers from academia and expert from industry to discuss about heterogeneous computing and explore collaboration opportunities

Technical session: presentations

- Applications and Algorithms
- Software stack and tools
- Heterogeneous systems Architecture

Breakout session: Discuss and report

- What are the next generation platforms specifications?
  - What is the ranked list of platform features that would enable this research?
- What are the key research problems that needs acceleration?
- What are the top three barriers preventing access and effective use of these platforms?
- What are the challenges and Opportunities
- What barriers could CMC help lower ?
- What would enable and encourage a vibrant community of researchers sharing platform configuration files and other knowledge?
- Any other feedback and guidance

# Q&A

Yassine Hariri
Hariri@cmc.ca