

Making the most out of Heterogeneous Chips with CPU, GPU and FPGA

Rafael Asenjo

Dept. of Computer Architecture University of Malaga, Spain.







Xilinx Zynq UltraScale+ MPSoC



Agenda

Motivation

- Hardware: iGPUs and iFPGAs
- Software
 - Programming models for chips with iGPUs and iFPGAs
- Our heterogeneous templates based on TBB
 - parallel_for
 - parallel pipeline

- A new mantra: Power and Energy saving
- In all domains





Motivation	System	Cores	(TFlop/s)	(kW)	(GFlops/watts
	TSUBAME3.0 - SGI ICE XA, IP139- SXM2, Xeon E5-2680v4 14C 2.4GHz, Intel Omni-Path, NVIDIA Tesla P100	36,288	1,998.0	142	14.110
 GPUs came to rescue: 	SXM2 , HPE GSIC Center, Tokyo Institute of Technology		Xeon +	⊦ Tes	la P100
 Massive Data Parallel Code at a low price in terms of power 	kukai - ZettaScaler-1.6 GPGPU system, Xeon E5-2650Lv4 14C 1.7GHz, Infiniband FDR, NVIDIA Tesla P10 <u>0</u> ,	10,080	460.7	33	14.046
-Supercomputers and servers: NVIDIA	ExaScalar Yahoo Japan Corporation Japan		Xeon +	- Tes	la P100
• GREEN500 Top 5:	AIST AI Cloud - NEC 4U-8GPU Server, Xeon E5-2630Lv4 10C 1.8GHz, Infiniband EDR, NVIDIA Tesla P100 SXM2 , NEC	23,400	961.0	76	12.681
– June 2017	National Institute of Advanced Industrial Science and Technology Japan		Xeon +	⊦ Tes	la P100
• Top500.org (Nov 2017):	RAIDEN GPU subsystem - NVIDIA DGX-1, Xeon E5-2698v4 20C 2.2GHz, Infiniband EDR, NVIDIA Tesla P100, Fujitsu	11,712	635.1	60	10.603
–87 systems w. NVIDIA –12 systems w. Xeon Phi	Center for Advanced Intelligence Project, RIKEN Japan		Xeon +	⊦ Tes	la P100
–5 systems w. PEZY	Wilkes-2 - Dell C4130, Xeon E5- 2650v4 12C 2.2GHz, Infiniband EDR,	21,240	1,193.0	114	10.428
6	University of Cambridge United Kingdom	\neg	Xeon +	- Tes	la P100

	T0P500		Green	500 Nov .	2017		Rmax	Power	Power Efficiency
Rank	Rank	System		Cores	(TFlop/s)	(kW)	(GFlops/watts)		
1	259	Shoubu system B - ZettaScaler-2.2, Xeon D-1571 16C 1.3GHz, Infiniband EDR, PEZY-SC2 , PEZY Computing / Exascaler Inc. Advanced Center for Computing and Communication, RIKEN Japan			794,400	842.0	50	17.009	
					Xeon + PEZY				
2	307	Suiren2 - Zett Infiniband EDI	aScaler-2.2, X R, PEZY-SC2 ,	eon D-1571 160 PEZY Computin	C 1.3GHz, g / Exascaler	762,624	788.2	47	16.759
	High Energy A Japan		ccelerator Re	search Organiza	ation /KEK	Xec	n + PEZY		
3	276	Sakura - Zetta Infiniband EDI	aScaler-2.2, Xo R, PEZY-SC2 ,	eon E5-2618Lv3 PEZY Computin	8C 2.3GHz, g / Exascaler	794,400	824.7	50	16.657
		PEZY Computing K.K. Japan			Xeon + PEZY				
4	149	DGX SaturnV 20C 2.2GHz, Ir NVIDIA Corpo United States	Volta - NVIDIA nfiniband EDR ration	, DGX-1 Volta36, , NVIDIA Tesla V	Xeon E5-2698v4 100 , Nvidia	22,440	1,070.0	97	15.113
5	4	Gyoukou - Zet 1.3GHz, Infinit Japan Agency Japan	taScaler-2.2 F band EDR, PE2 for Marine-Ea	HPC system, Xeo YY-SC2 700Mhz arth Science and	on D-1571 16C , ExaScaler d Technology	19,860,000	19,135.8 on + PEZ	1,350 Y	14.173

• There is (parallel) life beyond supercomputers:



Source: IDC. HPC forecast Intersect360

- Plenty of GPUs elsewhere:
 - Integrated GPUs on more than 90% of shipped processors



- Plenty of GPUs on desktops and laptops:
 - Desktops (35 130 W) and laptops (15 60 W):



Intel Kaby Lake

http://www.anandtech.com/show/10610/intel-announces-7th-genkaby-lake-14nm-plus-six-notebook-skus-desktop-coming-injanuary/2

AMD APU Kaveri



http://www.techspot.com/photos/article/770-amd-a8-7600-kaveri/

• Plenty of integrated GPUs in mobile devices too.

Huawei HiSilicon Kirin 960 (2 - 6 W)

Huawei Mate 9





• And user-programmable DSPs as well:

Qualcomm Snapdragon 820/821 (2 - 6 W)











https://www.qualcomm.com/products/snapdragon/processors/820

- And FPGAs too...
 - ARM + iFPGA
 - Altera Cyclon V
 - Xilinx **Zynq**-7000
 - Intel + Altera's FPGA
 - Heterogeneous Architecture Research Platform
- CPU + iGPU + iFPGA
 - Xilinx Zynq UltraScale+
 - 4 Cortex A53
 - 2 Cortex R5
 - GPU Mali 400
 - FPGA



Note: Illustration not drawn to scale.



- Plenty of room for improvements
 - Want to make the most out of the CPU, the GPU and the FPGA
 - Taking productivity into account:

Productivity = Performance - Pain

- High level abstractions in order to hide HW details
- Lack of high-level productive programming models
- "Homogeneous programming of heterogeneous architectures"
- Taking energy consumption into account:
 Performance → Performance / Watt

Performance → Performance / Joule

Productivity = Performance / Joule - Pain

Agenda

- Motivation
- Hardware: iGPUs and iFPGAs
- Software
 - Programming models for chips with iGPUs and iFPGAs
- Our heterogeneous templates based on TBB
 - parallel_for
 - parallel pipeline

Hardware: iGPUs



Intel Skylake – Kaby Lake

ULT 2+3 <u>2+2</u> 4+2 Modular design PCIE FDI PCIE FDI LPT-LP DDI SA DDI -2 or 4 cores SA OPIO Core Core - GPU Core DDI OPIO Core SA Core GT-2 • GT-1: 12 EU Core Core GT-2 Core • GT-2: 24 EU • GT-3: 48 EU GT-3 • GT-4: 72 EU ULT 2+2 4+3 SKU options LPT-LP eDRAM - 2+2 (4.5W, 15W) OPIO - 4+2 (35W -- 91W) FDI OPIO PCIE DDI OPIO DDI SA SA - 2+3 (15W, 28W,...) 1.5M Core Core 00 - 4+4 (65W) Core Core GT-3 1.5N Core GT-2

Core

http://www.anandtech.com/show/6355/intels-haswell-architecture

•

Intel Skylake – Kaby Lake



https://software.intel.com/sites/default/files/managed/c5/9a/The-Compute-Architecture-of-Intel-Processor-Graphics-Gen9-v1d0.pdf

Intel Processor Graphics



20

Intel Processor Graphics



Intel Processor Graphics



Intel Virt. Technology for Directed I/O (Intel VT-d)



AMD Kaveri

• Steamroller microarch (2 – 4 "Cores") + 8 GCN Cores.



http://wccftech.com/

AMD Graphics Core Next (GCN)

- In Kaveri,
 - 8 Compute Units (CU)
 - Each CU: 4 x 16-wide SIMD
 - Total: 512 FPUs
 - 866 MHz
- Max GFLOPS=
 - 0.86 GHz x
 - 512 FPUs x
 - 2 fmad =
 - 880 GFLOPS



OpenCL execution on GCN



HiSilicon Kirin 960



• ARM64 big.LITTLE (Cortex A73+Cortex A53) + Mali G71

Mali G71

- Supports OpenCL 2.0
- 4, 6, 8, 16, 32 Cores
- Each core:
 - 3 x 4-wide SIMD units
- 32 x 12 x 2 fmad x
 0.85GHz=652 GFLOPS





HSA (Heterogeneous System Architecture)

- HSA Foundation's goal: Productivity on heterogeneous HW
 CPU, GPU, DSPs,...
- Founders: AMDIN ARM Comagination



Advantages of iGPUs

- Discrete and integrated GPUs: different goals
 - NVIDIA Pascal: 3584 CUDA cores, 250W, 10 TFLOPS
 - Intel Iris Pro Graphics 580: 72EU x 8 SIMD, ~15W, 1.15 TFLOPS
 - Mali-G71: 32 EU x 3 x 4 SIMD, ~ 2W, 0.652 TFLOPS
- CPU and GPU are both first-class citizens.
 - SVM and platform atomics allows for closer collaboration
 - Avoid PCI data transfer and associated POWER dissipation
 - Operating System doesn't get in the way \rightarrow less overhead
- CPU and GPU may have similar performance
 - It's more likely that they can collaborate
- Cheaper!

Hardware: iFPGAs



FPGA Architecture



Courtesy: Javier Navaridas, U. Manchester.

CLB: Configurable Logic Block



Courtesy: Jose Luis Nunez-Yanez, U. Bristol.



How do we "mold this liquid silicon"?

- Hardware overlays
 - Map a CPU or GPU onto the FPGA
 - ✓ 100s of simple CPU can fit on large FPGAs
 - ✓ The CPU or GPU overlay can change/adapt
 - Resulting overlay is not as efficient as a "real" one
 - Same drawbacks as "general purpose" processors



FPGAs for Software Programmers, Dirk Kock, Daniel Ziener

How do we "mold this liquid silicon"?



✓ No FETCH nor DECODE of instructions → already hardwired
 ✓ Data movement (Exec. Units ⇔ MEM) reduction → Power save
 ✓ Not constrained by a fixed ISA:

Bit-level parallelism; shift, bit mask, ...; variable precision arith.
 Less frequency (hundreds of MHz)
 In order execution (proposal from David Kaeli's group to OoO)
 Programming effort

Hardware thread reordering to boost OpenCL throughput on FPGAs, ICCD'16
Xilinx Zynq UltraScale+ MPSoC

- Available:
 - − SW coherence (cache flush) \rightarrow Slow
 - HW coherence (MESI coherence protocol) \rightarrow Fast. Ports:
 - Accelerator Coherency Port, ACP
 - Cache-coherent interconnect (CCI) ACE-Lite ports
 - Part of ARM® AMBA® AXI and ACE protocol specification



GP

I2C UART USB 2.0 SPI ad SPI NOt NAND

Processing System

4 x

A53

Intel Broadwell+Arria10

- Cache Coherent Interconnect (CCI)
 - New IP inside the FPGA: Intel Quick Path Interconnect with CCI
 - AFU (Accelerated Function Unit) can access cached data



Towards CCI everywhere

- CCIX consortium (<u>http://www.ccixconsortium.com</u>)
 - New chip-to-chip interconnect operating at 25Gbps
 - Protocol built for FPGAs, GPUs, network/storage adapters, ...
 - Already in Xilinx Virtex UltraScale+



Which architecture suits me best?

- All of them are Turing complete, but
- Depending on the problem:
 - CPU:
 - For control-dominant problems
 - Frequent context switching
 - GPU:
 - Data-parallel problems (vectorized/SIMD)
 - Little control flow and little synchronization
 - FPGA:
 - Stream processing problems
 - Large data sets processed in a pipelined fashion
 - Low power constraints







Agenda

- Motivation
- Hardware: iGPUs and iFPGAs

Software

- Programming models for chips with iGPUs and iFPGAs
- Our heterogeneous templates based on TBB
 - parallel_for
 - parallel pipeline

http://imgs.xkcd.com/comics/standards.png

HOW STANDARDS PROLIFERATE

see: A/C chargers, Character encodings, Instant Messaging, etc)



Programming models for heterogeneous: GPUs

- Targeted at exploiting **one device** at a time
 - CUDA (NVIDIA)
 - **OpenCL** (Khronos Group Standard)
 - OpenACC (C, C++ or Fortran + Directives \rightarrow OpenMP 4.0)
 - C++AMP (Windows' extension of C++. Recently HSA announced own ver.)
 - Sycl (Khronos Gruop's specification for "Single-source C++ progr.)
 - ROCm (Radeon Open Compute, HCC, HIP, for AMD GCN HSA–)
 - RenderScript (Google's Java API for Android)
 - ParallDroid (Java + Directives from ULL, Spain)
 - Many more (Numba Python, IBM Java, Matlab, R, JavaScript, ...)
- Targeted at exploiting **both devices** simultaneously (**discrete** GPUs)
 - Qilin (C++ and Qilin API compiled to TBB+CUDA)
 - OmpSs (OpenMP-like directives + Nanos++ runtime + Mercurium compiler)
 - ХКаарі
 - StarPU
- Targeted at exploiting **both devices** simultaneously (**integrated** GPUs)
 - Qualcomm Symphony
 - Intel Concord

Programming models for heterogeneous: GPUs



Programming models for heterogeneous: FPGAs

- HDL-like languages
 - SystemVerilog, Bluespec System Verilog
 - Extend RTL languages by adding new features
- Parallel libraries
 - CUDA, OpenCL
 - Instrument existing parallel libraries to generate RTL code
- C-based languages
 - **SDSoC** LegUp, ROCCC, Impulse C, Vivado, Calypto Catapult C
 - Compile* C into intermediate representation and from there to HDL
 - Functions become Finite State Machines and variables mem. blocks
- Higher level languages
 - Kiwi (C#), Lime (Java), Chisel (Scala)
 - Translate* input language into HDL

* Normally a subset thereof

Courtesy: Javier Navaridas, U. Manchester.

OpenCL 2.0

- OpenCL 2.0 supports:
 - Coherent SVM (Shared Virtual Memory) & Platform atomics
 - Example: allocating an array of atomics:

```
clContex = ...
                                                               CPU Host Code
clCommandQueue = ...
clKernel = ...
cl svm mem flags flags = CL MEM SVM FINE GRAIN BUFFER | CL MEM SVM ATOMICS;
atomic int * data;
data = (atomic int *) clsvMAlloc(clContext, flags, NUM * sizeof(atomic int), 0);
clSetKernelArgSVMPointer(clKernel, 0, data);
clSetKernelArg(clKernel, 1, sizeof(int), &NUM);
clEnqueueNDRangeKernel(clCommandQueue, clKernel, ....);
atomic_store_explicit ((atomic int*)&data[0], 99, memory order release);
kernel void myKernel(global int *data, global int* NUM){
                                                            GPU Kernel Code
int i=atomic load explicit((global atomic int *)&data[0], memory order acquire);
. . . .
}
```

C++AMP

- C++ Accelerated Massive Parallelism
- Example: Vector Add (C[i]=A[i]+B[i])



SYCL Parallel STL

- Parallel STL aimed at C++17 standard:
 - std::sort(vec.begin(), vec.end()); //Vector sort
 - std::sort(seq, vec.begin(), vec.end()); // Explicit SEQUENTIAL
 - std::sort(par, vec.begin(), vec.end()); // Explicit PARALLEL
- SYCL exposes the execution policy (user-defined)



https://github.com/KhronosGroup/SyclParalleISTL

- Point Kernel: write once, run everywhere
 - Pattern tuning extends to heterogeneous load hints

```
SYMPHONY POINT KERNEL 3(vadd, int, i, float*, a, int, na,
                        float*, b, int, nb, float*, c, int, nc,
 { c[i] = a[i] + b[i];}); ...
 symphony::buffer_ptr buf_a(1024), buf_b(1024), buf_c(1024);
 symphony::range<1>(1024) r;
 symphony::pfor each(r, vadd pk, buf a, buf b, buf c,
           symphony::pattern::tuner().set_cpu_load(10)
                                       .set_gpu_load(50)
                                       .set_dsp_load(40)
 A kernel for
                   Executed across CPU(10%) + GPU(50%) + DSP(40%)
each device is
automatically
51 generated
```

Intel Concord

• C++ heterogeneous programming framework



- Papers:
 - Rashid Kaleem et al. Adaptive heterogeneous scheduling on integrated GPUs. PACT 2014.
 - Intel Heterogeneous Research Compiler (iHRC) at https://github.com/IntelLabs/iHRC/
 - Naila Farooqui et al. Affinity-aware work-stealing for integrated CPU-GPU processors. PPoPP '16. PhD dissertation: https://smartech.gatech.edu/bitstream/handle/1853/54915/FAROOQUI-DISSERTATION-2016.pdf



Xilinx SDSoC



Altera OpenCL



Altera OpenCL optimizations

- Locality:
 - Shift-Register Pattern (SRP)
 - Local Memory (LM)
 - Constant Memory (CM)
 - Load-Store Units Buffers (LSU)
- Loop Unrolling (LU)
 - Deeper pipelines
 - Optimized tree-like reductions
- Kernel vectorization (SIMD)
 - Wider ALUs
- Compute Unit Replication (CUR)
 - Replicate Pipelines
 - Downsides: memory divergence $\mathbf{1}$, complexity $\mathbf{1}$, frequency $\mathbf{1}$







Shift Register Pattern

• Example: 1D stencil computation



Agenda

- Motivation
- Hardware: iGPUs and iFPGAs
- Software
 - Programming models for chips with iGPUs and iFPGAs

Our heterogeneous templates based on TBB

- parallel_for
- parallel pipeline

Our work: parallel_for

- Heterogeneous parallel for based on TBB
 - Dynamic scheduling / Adaptive partitioning for irregular codes

Time

- CPU and GPU chunk size re-computed every time



New scheduler and API



New API

```
#include ''parallel_for.h''
1
   using namespace hbb;
\mathbf{2}
   int main(int argc, char* argv[]){
3
     // Start task scheduler
                                       Use the GPU
4
     HInit HInit (numcpus, true);
5
6
     KernelInfo k(kernelFile, kernelName);
7
     HBuffer<int> * a = new HBuffer<int>(N,USE_ZCB);
8
     Body body (k, a);
9
                                  Exploit Zero-Copy-Buffer
10
     parallel_for(begin, end, body, new LogFit());
11
12
13
                                       Scheduler class
```

```
1 #include ''HTask.h''
2 #include ''HBuffer.h''
3 using namespace hbb;
4
5 class Body : public HTask{
      HBuffer<int> * buf a;
6
   public:
7
     Body(KernelInfo k, HBuffer * buf_a) : HTask(k) {...}
8
     void operatorCPU(int begin, int end) {
9
        int *a = buf a->getHostPtr(BUF RW);
10
        for(i=begin; i!=end; i++) a[i] = a[i] + 1;
11
12
     void operatorGPU() (int begin, int end) {
13
        //Setting kernel arguments
14
        setKernelArg(0, sizeof(int), &begin);
15
        setkernelArg(1, sizeof(int), &end);
16
        setKernelBuf(2, sizeof(BUF), buf a->getDevicePtr(BUF RW));
17
        //Launching kernel
18
        launchKernel(begin, end);
19
\mathbf{20}
     }
21
   };
```

Related work: Qilin

Static Approach: Bulk-Oracle

Offline search (11 executions)

- Work partitioning between CPU and GPU



- One big chunk with 70% of the iterations: on the GPU
- The remaining 30% processed on the cores

Related work: Intel Concord

- Papers:
 - Rashid Kaleem et al.
 Adaptive heterogeneous scheduling on integrated GPUs. PACT 2014.
 - Naila Farooqui et al.
 Affinity-aware work-stealing for integrated CPU-GPU processors. PPoPP '16 (poster).



Related Work: HDSS



 Belviranli et al, A Dynamic Self-Scheduling Scheme for Heterogeneous Multiprocessor Architectures, TACO 2013

Debunking big chunks

• Example: irregular Barnes-Hut benchmark



Debunkign big chunks



A larger chunk causes more memory traffic

Debunking big chunks

- BarnesHut: GPU throughput along the iteration space
 - For two different time-steps
 - Different GPU chunk-sizes



LogFit main idea



Time

Parallel_for: Performance per joule



- Static: Oracle-Like static partition of the work based on profiling
- Concord: Intel approach: GPU size provided by user once
- HDSS: Belviranli et al. approach: GPU size computed once
- LogFit: our work

Antonio Vilches et al. Adaptive Partitioning for Irregular Applications on Heterogeneous CPU-GPU Chips, **Procedia Computer Science**, 2015.

LogFit: better for irregular benchmarks



LogFit working on FPGA

• Energy probe: AccelPower Cap (Luis Piñuel, U. Complut.)



LogFit working on FPGA



- Terasic DE1: Altera Cyclone V (FPGA+2 Cortex A9)
- Energy probe: AccelPower Cap (Luis Piñuel, U. Complut.)

Our work: pipeline

- ViVid, an object detection application
- Contains three main kernels that form a pipeline



- Would like to answer the following questions:
 - **Granularity**: coarse or fine grained parallelism?
 - Mapping of stages: where do we run them (CPU/GPU)?
 - Number of cores: how many of them when running on CPU?
 - **Optimum**: what metric do we optimize (time, energy, both)?
Granularity

Coarse grain:
CG

• Medium grain:

– MG

• Fine grain:



Fine grain also in the CPU via AVX intrinsics

Our work: pipeline



Accounting for all alternatives

• In general: **nC** CPU cores, **1** GPU and **p** pipeline stages



alternatives = $2^{p} x (nC + 2)$

• For Rodinia's SRAD benchmark (p=6,nC=4) → 384 alternatives

Mapping streaming applications on commodity multi-CPU and GPU onchip processors, **IEEE Tran. Parallel and Distributed Systems**, 2016.

Framework and Model

- Key idea:
 - 1. Run only on GPU
 - 2. Run only on CPU
- collect λ and E (homogeneous values)
- 3. Analytically extrapolate for heterogeneous execution
- 4. Find out the best configuration \rightarrow RUN



Environmental Setup: Benchmarks

Inpt St 1 St 2 St 3 Out St

St 5

Reduct. Statist. Comp.1 Comp. 2

St 6

Out =>

Filter Histogram Classifier

Four Benchmarks

- SRAD

- Tracking

- ViVid (Low and High Definition inputs)

St 1 St 2 St 3 St 4

Prep.



Extrac.



St 1

Track.

Does Higher Throughput imply Lower Energy?



SRAD on Haswell



Results on big.LITTLE architecture

• On Odroid XU3: 4 Cortex A15 + 4 Cortex A7 + GPU Mali



Pipeline on CPU-FPGAs

- Change the GPU by a FPGA
 - Core i7 4770k 3.5 GHz Haswell quad-core CPU
 - Altera DE5-Net FPGA from Terasic
 - Intel HD6000 embedded GPU
 - Nvidia Quadro K600 discrete GPU (192 cuda cores)
- Homogeneous results:

Device	Time (sec.)	Throughput (fps)	Power (watts)	Energy (Jules)
CPU (1 core)	167.610	0.596	37	6201
CPU (4 cores)	50.310	1.987	92	4628
FPGA	13.480	7.418	5	67
On-chip GPU	7.855	12.730	16	125
Discrete GPU	13.941	7.173	47	655

Pipeline on CPU-FPGA

- Heterogeneous results based on our TBB scheduler:
 - Accelerator + CPU
 - Mapping 001



4

0

Throughput (fps) Throughput (fps) Throughput (fps)

Conclusions

- Plenty of heterogeneous on-chip architectures out there
- Why not use all devices?
 - Need to find the best mapping/distribution/scheduling out of the many possible alternatives.
- Programming models and runtimes aimed at this goal are in their infancy: striving to solve this.
- Challenges:
 - Hide HW complexity providing a productive programming model
 - Consider energy in the partition/scheduling decisions
 - Minimize overhead of adaptation policies

Future Work

- Exploit CPU-GPU-FPGA chips and coherency
- **Predict energy** consumption on heterogeneous CPU-GPU-FPGA chips
- **Consider energy** consumption in the partitioning heuristic
- Rebuild the scheduler engine on top of the new TBB library (OpenCL node)
- Tackle other irregular codes

Future Work: TBB OpenCL node



Can express pipelining, task parallelism and data parallelism

Future Work: other irregular codes



Collaborators

- M^a Ángeles González Navarro (UMA, Spain)
- Andrés Rodríguez (UMA, Spain)
- Francisco Corbera (UMA, Spain)
- Jose Luis Nunez-Yanez (University of Bristol)
- Antonio Vilches (UMA, Spain)
- Alejandro Villegas (UMA, Spain)
- Juan Gómez Luna (U. Cordoba, Spain)
- Rubén Gran (U. Zaragoza, Spain)
- Darío Suárez (U. Zaragoza, Spain)
- Maria Jesús Garzarán (Intel and UIUC, USA)
- Mert Dikmen (UIUC, USA)
- Kurt Fellows (UIUC, USA)
- Ehsan Totoni (UIUC, USA)
- Luis Remis (Intel, USA)

Questions

