High-Performance Computing for Low-Power Systems

Ghent University, Belgium

Erik D'Hollander

erik.dhollander@ugent.be



Advanced HPC Systems workshop Cetraro, June 27-29, 2011





Overview

- GPU's vs. FPGA's
 - Power gain
 - Architecture
 - Programming
- Application
 - Wavelet CAD
 - GPU performance
 - FPGA simulation
- Conclusion



David & Goliath





Power gain: Miles per Gallon Iterations per Joule...







Architecture Commonalities



- Accelerators
- Parallel processing
- Streaming data









The weakest link...

GPU architectu. Fast thread scheduler!

• Fermi GPU



x16 crossbar







GPU programming

Number of threads (MxN)





CUDA



- Send stream \rightarrow Process \rightarrow Receive stream
- Between CPU and GPU:
 - DMA
- Within GPU:
 - Parallel threads
 - SPMD







FPGA as accelerator

- Send stream \rightarrow Process \rightarrow Receive stream ۲
- Between CPU and FPGA: ۲
 - DMA
- Within FPGA: •
 - Pipelined datapaths
 - Parallel datapaths





CPU







GPU \leftrightarrow **FPGA** operation

- GPU = Von Neumann: instruction fetch, decode... IF ID OF EX WB
 - instructions and operands in memory
 - operand address calculation requires instructions
- FPGA = Dataflow: instructions (operations) pipelined in hardware
 IF ID OF EX WB
 - 40% gain
 - instructions in pipeline, operands in memory
 - operand address calculation requires address generator (AG)
 - loops require loop controller (to feed data back into pipeline)







$\textbf{GPU} \rightarrow \textbf{FPGA ports?}$

Trials

- FCUDA^{*} : CUDA port to FPGA
- PGI CUDA : CUDA port to x86 platforms

Problems

- CUDA cannot express pipelines
- CUDA requires $C \rightarrow VHDL$ compiler







ROCCC*: a direct approach

- C→VHDL compiler
 - Creates pipelined datapaths
 - Creates parallel loop control
 - Supports streaming paradigm
 - Hierarchical design (modules, systems)
 - Eclipse programming environment
 - VHDL Testbench generation
 - Integrates with existing IP-cores and Xilinx toolchain
- Limitations with respect to FPGA environment
 - No use of pointers
 - Array index expressions limited to loop index + constant e.g. a[i+6]







ROCCC*: modules and systems

Modules

- Computational datapath
- Translates into pipelined logic blocks
- Scalar arguments, scalar results
- Loops unrolled
- E.g. FIR filter (5 taps);

void FIR(int A0, int A1, int A2, int A3, int A4, int& result)
 { const int T[5] = { 3, 5, 7, 9, 11 } ;
 result = A0*T[0] + A1*T[1] + A2*T[2] + A3*T[3] + A4*T[4];



*Jason R. Villarreal, Adrian Park, Walid A. Najjar, Robert Halstead: Designing Modular Hardware Accelerators in C with ROCCC 2.0. FCCM 2010: 127-134



ROCCC*: modules and systems

Systems

- Loops on streaming data I/O \rightarrow loop controller
- Operates on streams of arrays \rightarrow address generator
- Compiler organizes reads/writes of data in each clock cycle
- E.g. FIR array



*Jason R. Villarreal, Adrian Park, Walid A. Najjar, Robert Halstead: Designing Modular Hardware Accelerators in C with ROCCC 2.0. FCCM 2010: 127-134



ROCCC: Smart Buffer

Objective: minimize memory accesses to FPGA BlockRAM

- Fetched min. data in each cycle
- Stores data in registers
- Reuses data in registers
- Optimizes bandwidth
- Minimizes pipeline stalls
- E. g. FIR filter

 $\begin{array}{ll} a[i] &\to r_{(i \mod 5)} \\ b[i] &\leftarrow FIR(r0,r1,r2,r3,r4) \\ b[i+1] &\leftarrow FIR(r1,r2,r3,r4,r0) \\ b[i+2] &\leftarrow FIR(r2,r3,r4,r0,r1) \\ b[i+3] &\leftarrow FIR(r3,r4,r0,r1,r2) \\ b[i+4] &\leftarrow FIR(r4,r0,r1,r2,r3) \end{array}$





ROCCC: Smart Buffer

• Extended to multiple streams, multiple windows

• E.g. Edge detection

B[i-1][j-1] = maskV*maskV + maskH*maskH;

Replace array elements by smart buffer regs Remove loop headers (replace by loop control)



ROCCC: Smart Buffer

- Extended to multiple streams, multiple windows
- Resulting program: loop body = scalar operations on registers

 $B_0_0 = maskV*maskV + maskH*maskH;$





Smart buffer windows







ROCCC: Other optimizations

*Stream management

- Exploit wide memory bus (e.g. 64 bit) to feed multiple // streams
- Expandable internally with very wide dual port BlockRAM





*ROCCC tutorials



ROCCC: Other optimizations

- Loop unrolling
- Loop tiling
- Loop interchange
- If predication
- Scalar replacement
- Sequential loop pipelining
- Systolic array wavefront generation
- Pipelining retiming







FPGA programming toolchain

 $C \rightarrow VHDL \longleftarrow ROCCC$ **FPGA** optimizing compiler 1 void MatrixMultiplication(int** A, int** B, int*** C0) 000 0 0000 2 { 0000000 0000 з int i ; 00000000 int j; 4 int k ; int currentSum ; for (i = 0; i < 10; ++i)8 9 for (j = 0; j < 10; ++j)10 11 12 currentSum = 0 ; for (k = 0; k < 10; ++k)13 14 15 currentSum += A[i][k] * B[k][j] ; 16 17 CO[i][j][k] = currentSum ; 18 19 3 20 }

Generate circuit







Application: DWT

- Discrete wavelet transformation for Computer Aided Diagnosis
- **Objective**: detect microcalcifications (small granularity) in X-ray
- Solution:
- DWT of image
- Dismiss low frequencies
- IDWT of image
- \Rightarrow enhanced image









Image enhancement

Steps

- 1. Discrete Wavelet Transform
- 2. Remove LL part
- 3. Inverse DWT
- 4. Normalize
- 5. Stretch















GPU speed-up

- Using CUDA, parallelizing 5 phases^{*}
 - GT-330 NVIDIA (desktop GPU)

Start	Processed on GPU Execution Time (ms): 176	Processed on CPU Execution Time (ms): 4385	
			Decomposition Leve 1 2 3 4 5
			Getal tussen 0 en 1 Threshold : 0
			Wavelet daub 1 daub 2 daub 3 daub 4
			daub5 daub6 daub7 daub8 daub9



*Dries Rosseel, Generating speed-up for the optimization of RX-images, Msc thesis, 2010



GPU

176 ms

CPU 4,385 ms

Sp = 24.9

GPU – performance analysis

• Using NSight







GPU – performance analysis

• Using NSight

4	۰ 🕥 🖌	Timelin	ie									
rs	Y		Sections	un de la compañía de							_	
	Seconds		15.1693	119431	15.1703119431	15.17	13119431	15.1723119431	15.17	/33119431	15.1743119431	15.175
C	JDA		1	-								
	Context 0											
	Context 1 [0]										-	
	Runtime API	8		cudaMallo	cudaMallo	cudaMallo		cuda	aLaunch [436	0]	cudaLaun	cudaFree
	Driver API	¥.		cuMemAl	cuMemAll	cuMemAl		cuLaun	hGridAsync	[4360]	cuLaunc	cuMemFr
	Warnings	Y.	Ranges									
	Memory	8										
3	Compute	V.						-	synthcolKern	el	synthr	
	0.0 % [6] s	Y	1						ynthcolKern	el		
	0.0 % [6] t	Y	1									
	0.0 % [2] s	8	i			15 1710700					15 1720202	
	0.0 % [6] t	Y				(A 0.00256)					(A 0.00256)	
	0.0 % [6] s	V.	1			(12 01002 20)					synthr	





GPU – performance analysis

• Using NSight

Top Device Functions by Total Time						
Launch Summary All						
Name	Launches		Total (µs)	Min (µs)	Avg (µs)	Max (µs)
synthcolKernel		6	10,340.21	1,716.63	1,723.37	1,731.00
transcolKernel		6	6,785.76	1,122.59	1,130.96	1,137.85
setRGBdev		2	3,466.74	1,729.27	1,733.37	1,737.47
transrowKernel		6	3,331.95	551.677	555.325	559.166
synthrowKernel		6	2,458.23	404.029	409.705	420.319



synthcolKernel
transcolKernel
setRGBdev
transrowKernel

synthrowKernel





FPGA – ROCCC code

DWT kernel extraction: transrow kernel

- Modified in several ways:
- A) code
 - outer loop parallel
 - inner loop streamlined
 - innerst loop unrolled
 - bodies pipelined
- B) data
 - output "dest[][]" →
 2 arrays "desty[]"

```
(unsigned int y = 0; y < h; y++)
for (unsigned int k = 0; k < w2; k++)
  s = 0.0f;
  d = 0.0f;
    for (int m = tH.first(); m <= tH.last(); m++)</pre>
      n = 2 * k + m;
      if (n < 0) n = 0 - n;
      if (n \ge (int)w) n -= 2 * (1 + n - w);
      s += tH[m] * float(sour[y][n]);
    for (int m = tG.first(); m <= tG.last(); m++)</pre>
      n = 2 * k + m;
      if (n < 0) n = 0 - n;
      if (n \ge (int)w) n -= 2 * (1 + n - w);
      d += tG[m] * float(sour[y][n]);
    dest[y][k] = mmxround(s);
    dest[y][k+w2] = mmxroundTH(d);
```





ROCCC kernel

• DWT transrow kernel @ 12 mul, 10 add = 22 ops

• Kernel \rightarrow VHDL \rightarrow testbench





FPGA – ROCCC performance analysis

• Single stream: 194 cycles \Rightarrow 22 ops / (194 * 8 ns) = 1.42 Gop/s





FPGA – ROCCC performance analysis

• 2 streams: 107 cycles; 4 streams: 63 cycles





4 streams: 74 cy@ 74MHz = 2.2 GFlops latency = 30cy@13ns = 390 ns

				20.500 ns				94.500 ns
				and a second				and a second
N	ame	Value		20 ns	l0 ns	60 ns	80 ns	100
	Le cik	1						
	Ve rst	0					÷.	
	🔏 soury_writeen_in	0						
*	🍓 soury_channel0_in[31:0]	40800000	UUUUU		0000000000000	x00000	40800000	
*	🧋 soury_channel1_in[31:0]	40400000	UUUUU	0000000000	0000000000000	x00000	40400000	
*	📲 soury_channel2_in[31:0]	40000000	UUUUU.	0000000000	0000000000	x00000	4000000	
*	soury_channel3_in[31:0]	3£800000	UUUUU	0000000000	0000000000	xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx	3f800000	
	🕫 outputready	0						
٠	🧋 destyg_channel0_out[31:0]	3£800000		00000000		400000)/))//CO0)////b	f800\\\ c00\\\	3f800000
*	📲 destyg_channel1_out[31:0]	00000000		00000000	0000	0000 💥 0000000 🕅	f800) 00000	000
*	destyg_channel2_out[31:0]	00000000		00000000	0000	0000 💥 0000000 🕅	f800) 00000	000
*	📲 destyg_channel3_out[31:0]	bf800000		00000000	bf800000	000000)XXX()(3f8)X()(3	f800):::::::::::::::::::::::::::::::::	bf800000
*	i destyg_channel0_address_ou	22		0	3	6	15 18	22
	🕼 destyg_read_out	1						
٠	📲 destyh_channel0_out[31:0]	41a8c000		00000000	(42) (42b1)	(42a0) (426)	122b)%(\41fb)%/%	41a8c000
*	📲 destyh_channel1_out[31:0]	41a0c000		00000000	42) (42af3)	429e) (426) (426)	227 XX (41f3)	41a0c000
*	destyh_channel2_out[31:0]	4198c000		00000000	42) 42ad)	(429c0) (426)	223)%(\41e)	4198c000
*	📲 destyh_channel3_out[31:0]	41918000		00000000	42b X 42ab X	4299	21e) 👯 (41e)	41918000
*	🥂 destyh_channel0_address_ou	22		0	XXX 3 XX	6	15 18	22
	le soury_stream_num_channels	4			4			
	🕼 destyg_stream_num_channel	4			4			
	👍 destyh_stream_num_channel	4			4			
								74.000 ns
				:0 ns	20 ns	40 ns	60 ns	80 r
			21-01-02			<u></u>		

X1: 94.500 ns X2: 20.500 ns ΔX: 74.000 ns

$ROCCC \leftrightarrow Handmade \leftrightarrow GPU$

Power of $C \rightarrow VHDL$ synthesis ٠

- ROCCC image kernel: 2.2 GFlops @ 74 MHz, Virtex 4 •
- ۲
- GPU image kernel: ٠
- Manual BLAS^{*} 2 fpgas: 7.8 GFlops @ 140 MHz, StratixII 60
 - 5.4 Gflops @ 1300 MHz, GT330
- Improvements: use a pipelined multiplier \rightarrow lower cycle time ۲



[^]S. Rousseaux, et al. A High Performance FPGA-based accelerator for BLAS library Implementation, Proc. Reconfigurable Systems Summer Institute, 2009.

Conclusion

- ROCCC toolchain greatly facilitates HPC on FPGA
- Good performance, low power
- Areas for improvement:
 - Enhanced Address Generator to account for more complex index expressions^{*}
 - Host-FPGA bandwidth
- Areas of application?
 - Recognize LHF of hardware kernels
 - Optimize data locality using reuse distance analysis
 - Think spatial, parallel and pipelined







?



