















## CALCM Computer Architecture Lab

Carnegie Mellon

## A hard look at the performance and ops-per-Joule of today's non-von Neuman alternatives



	©CALCM	<b>Computer Archite</b>	cture Lab		and a start of the	(	Carnegie Mellon
		GPU	s, FPC	GAs an	d all t	hat	
		CPU		GPUs		FPGA	ASIC
		Intel Core i7-960	Nvidia GTX285	Nvidia GTX480	ATI R5870	Xilinx V6-LX760	Std. Cell
	Year	2009	2008	2010	2009	2009	2007
	Node	45nm	55nm	40nm	40nm	40nm	65nm
	Die area	263mm <sup>2</sup>	470mm <sup>2</sup>	529mm <sup>2</sup>	334mm <sup>2</sup>	-	-
	Clock rate	3.2GHz	1.5GHz	1.4GHz	1.5GHz	0.3GHz	-
	Single-prec. floating-point apps						
	M-M-Mult	MKL 10.2.3 multithreaded	CUBLAS 2.3	CUBLAS 3.1	CAL++	hand-coded	
	FFT	Spiral.net multithreaded	CUFFT 2.3 3.0/3.1	CUFFT 3.0	-	Spira	l.net
	Black-Scholes	PARSEC multithreaded	CUDA 2.3	-	-	hand-	coded
CM	U/ECE/CALCM/Chung&	&Hoe			Lo	s Alamos CS Symposium	n, October 2010, slide-10

()CAL	CM Computer Architecture Lab			Carnegie Mellon
	In-Core Per	formance	e and Ene	ergy
	Device	GFLOP/s actual	(GFLOP/s)/mm <sup>2</sup> normalized to 40nm	GFLOP/J normalized to 40nm
	Intel Core i7 (45nm)	96	0.50	1.14
	Nvidia GTX285 (55nm)	425	2.40	6.78
NANANA	Nvidia GTX480 (40nm)	541	1.28	3.52
IVIIVIIVI	ATI R5870 (40nm)	1491	5.95	9.87
	Xilinx V6-LX760 (40nm)	204	0.53	3.62
	same RTL std cell (65nm)	694	19.28	50.73
<ul> <li>C</li> <li>Si</li> <li>Pi</li> <li>fr</li> <li>Fi</li> </ul>	PU and GPU benchm td Cell effectively cor ower (switching+leal rom the system	arking was cor npute-bound ( kage) measure et al MICRO 2	mpute-bound; F no off-chip I/O) ments isolated 1	PGA and the core
• F(	or detail see [Chung,	et al. MICRO 2		manaium Ostahan 2010 alida 11

@CAL	.CM Computer Architecture Lab			Carnegie Mellon
	Mor	e of the	Same	
		GFLOP/s	(GFLOP/s)/mm <sup>2</sup>	GFLOP/J
	Intel Core i7 (45nm)	67	0.35	0.71
	Nvidia GTX285 (55nm)	250	1.41	4.2 4.3 - 6.5 90 Monts/1
-2 <sup>10</sup>	Nvidia GTX480 (40nm)	453	1.08	4.3
FFT	ATI R5870 (40nm)	-	-	-
	Xilinx V6-LX760 (40nm)	380	0.99	6.5
	same RTL std cell (65nm)	952	239	90
		Mopt/s	(Mopts/s)/mm <sup>2</sup>	Mopts/J
	Intel Core i7 (45nm)	487	2.52	4.88
les	Nvidia GTX285 (55nm)	10756	60.72	189
scho	Nvidia GTX480 (40nm)	-	-	-
ack-9	ATI R5870 (40nm)	-	-	-
Bla	Xilinx V6-LX760 (40nm)	7800	20.26	138
	same RTL std cell (65nm)	25532	1719	642.5
CMU/ECE/CALC	M/Chung&Hoe		Los Alamos CS Sv	mposium, October 2010, slide-1





$arphi$ and $\mu$ example values					
	МММ	Black-Scholes	FFT-2 <sup>10</sup>		
Φ	0.74	0.57	0.63		
μ	3.41	17.0	2.88		
Φ	0.77		0.47		
μ	1.83	On equa	l area basis,		
Φ	1.27	3.41x pe	rformance		
μ	8.47	at 0.74x	power		
Φ	0.31	relative a	BCE.29		
μ	0.75	5.68	2.02		
Φ	0.79	4.75	4.96		
μ	27.4	482	489		
	<b>Φ</b> μ Φ μ Φ μ Φ μ	ΜΜΜ           Φ         0.74           μ         3.41           Φ         0.77           μ         1.83           Φ         1.27           μ         8.47           Φ         0.31           μ         0.75           Φ         0.79	MIMM         Black-Scholes           Φ         0.74         0.57           μ         3.41         17.0           Φ         0.77         0           Φ         0.77         0           μ         1.83         On equal           Φ         1.27         3.41x pe           μ         8.47         at 0.74x           Φ         0.31         relative a           μ         0.75         5.68           Φ         0.79         4.75           μ         27.4         482		



Year	2011	2013	2016	2019	2022
Technology	40nm	32nm	22nm	16nm	11nm
Core die budget (mm <sup>2</sup> )	432	432	432	432	432
Normalized area (BCE)	19	37	75	149	298 <mark>(16</mark> x
Core power (W)	100	100	100	100	100
Bandwidth (GB/s)	180	198	234	234	252 <mark>(1.4x</mark>
Rel pwr per device	1X	0.75X	0.5X	0.36X	0.25X

U

U

• 432mm<sup>2</sup> populated by an optimally sized Fast Core and U-cores of choice











CALCM Computer Architecture Lab

## **Performance Scaling Summary**

- Perf-per-Watt and op-per-Joule matter
  - need to look beyond programmable processors
  - GPUs and FPGAs are viable programmable candidates
- GPU/FPGA/custom logic help performance only if
  - significant fraction amenable to acceleration, and
  - adequate bandwidth to sustain acceleration

3D-stacked memory could be very helpful!

**Carnegie Mellon** 

- Without adequate bandwidth, GPU/FPGA catches up with custom logic in achievable performance but remains programmable and flexible
- Custom logic best for maximizing performance under tight energy/power constraints









