### Supercomputer Scaling Life in the Fast Lane

Supercomputers for the masses: Get a bunch of cheap PCs String a lot of cable Voila! A supercomputer!





Supercomputers for the masses: Get a bunch of cheap PCs String a lot of cable Voila! A supercomputer!

### Or maybe not...





Failures multiply Always a few dead nodes Can turn into many

> Considerations: Heat removal Power conditioning Power sequencing Vibration



Communication Ethernet is high latency Infiniband

Cheap PCs often use multiple chipsets



# A Rack Cluster



# A Computer Room



# AComputer

ASCI White Lawrence Livermore National Laboratory



# A Computer Power Supply



# Another Computer



# A Computer Building



# Earth Simulator 2

1280 Nodes 131 TFlops 8 CPUs/node 4-way superscalar 256-way vector



### NEC SX-9 vector processors

# Original Teraflop



### ASCI Red: ~9K Pentium II, 3D mesh

# ASC-Q



# ASCI Purple

### ~12K Power5, 100 TFLOPS, 50 TB mem, 2 PB disk, 7.5 MW



### Red Storm

~13K quad-core Opterons, 284 TFLOPS, Cray XT-4 (PPC 440-based network)

### BlueGene L

-

596 TFLOPS, 75 TB memory, 1.9 PB disk 64K Dual PPC440, 3D Torus, Collective, Control

### Roadrunner

11 m



# Nore Blue Gene

Blue Gene/P: up to ~880K PPC450 @ 850 MHz, 3 PFLOPS Blue Gene/Q: (2011) 1.6M PPC cores, 6MW, 20 PFLOPS 

### Blue Gene Q

This one (ANL) 9PFLOPS, Sequoia (LLNL) 17PFLOPS Each chip: 18 64-bit 4-way SMT PPC cores @1.6 GHz 96 racks, 98K nodes, 3K sq. ft. 7.9MW, 2GFLOPS/W

## ORNL Titan



### 560,640 cores, plus GPU accelerators, 8.2 MW, 17.6 Petaflops

### Tianhe-2



3.12 Million cores (Intel Ivy Bridge and Phi) 17.8 MW, 33.86 Petaflops

### Blue Waters Facility 10 PFLOPS (peak) system being built at NCSA (Univ. of Illinois) 90,000 sq ft building 20,000 sq ft machine room, w/6' raised floor 32 MW power in 4 independent feeds



### Blue Waters Processors

- Based on IBM Power 7
- 8-core, 32-thread, 4 GHz processor
- 32MB level 3 cache
- 2 Memory controllers, 8 channels
- 4 chips per module 32 cores
- 1 TFlop per module
- 1.1 TB/s switch blocks



### Blue Waters Drawers

8 processor modules (TFLOPS) per drawer ■ 4 GB/core, 1 TB/drawer Water cooled ■ 39" x 72" x 2" **290 lbs** 



### Drawer Contents

### TB/s switches

### PCI-E slots

Optical fiber

500GB

RAM

### Power 7 quad-chip modules

500GB RAM

> Power Supplies

### Chilled water pipes

# Programming

- Initially MPI and Open-MP
- Migration to X10
  - Partitioned Global Address Space (PGAS)
  - Eclipse environment
  - Object-oriented, with ateach, foreach, atomic, clocks
  - Run on a virtual machine on top of HW



### Nuddied Waters

Original design cancelled IBM cannot deliver stated performance for bid price NSF can't pay more, IBM can't dump hardware 235 Cray XE6 cabinets with AMD Opteron 6200 30 future Cray XK6 cabinets with NVIDIA Tesla GPUs

# New Design

380,000 x86 cores + 3000 GPUs Solution 3D torus interconnect (Cray Gemini - sold to Intel) 1.5 PB memory, 25 PB disk 11.5 PF peak performance

# Exatlop Challenges

- 1 billion gigaflops, 1 million teraflops
- Traditionally 1 byte per flop mostly memory machine
- parallelism
- Limited power and space (25 MW costs \$15M/yr)
  - 25 picowatts per flop -- factor of 100 more efficient
- Heterogeneity, specialized processors, higher cost
- How do you checkpoint a thing like this?



### No more clock scaling -- factor of 1000 beyond current petaflop systems all via additional

# Kevin LimISCA 2008Understanding and Designing New Server Architectures for<br/>Emerging Warehouse-Computing Environments

# Benchmark Suite

52525252525252525252525					
Table 1: Summary details of the new benchmark suite to represent internet sector workloads.					
Workload	Emphasize	Description	Perf metric		
websearch	the role of unstructured data	Open source Nutch-0.9, Tomcat 6 with clustering, and Apache2. 1.3GB index corresponding to 1.3 million indexed documents, 25% of index terms cached in memory. 2GB Java heap size. QoS requires >95% queries take <0.5 seconds.	Request-per- sec (RPS) w/ QoS		
webmail	interactive internet services	Squirrelmail v1.4.9 with Apache2 and PHP4, Courier-IMAP v4.2 and Exim4.5. 1000 virtual users with 7GB of mail stored. Email/attachment sizes and usage patterns modeled after MS Exchange 2003 LoadSim heavy users. QoS requires >95% requests take <0.8 second.	RPS w/ QoS		
ytube	the use of rich media	Modified SPECweb2005 Support workload with Youtube traffic characteristics. Apache2/Tomcat6 with Rock httpd server.	RPS w/ QoS		
mapreduce	web as a platform	Hadoop v0.14 with 4 threads per CPU and 1.5GB Java heap size. We study two workloads - distributed file write (mapred-wr) and word count (mapred-wc)	Execution time		

### Cost/Performance in Server Farms

Performance is measured in sustainable metrics, such as requests per second

Total Cost of Ownership

Includes direct cost, power, cooling, depreciation

Cost may be calculated in terms of TCO, power, etc. 

# Systems Considered

### Infrastructure dollars

Table 2: Summary of systems considered.								
System	"Similar to"	System Features	Watt	Inf-\$				
Srvr1	Xeon MP, Opteron MP	2p x 4 cores, 2.6 GHz, OoO, 64K/8MB L1/L2	340	3,294				
Srvr2	Xeon, Opteron	1p x 4 cores, 2.6 GHz, OoO, 64K/8MB L1/L2	215	1,689				
Desk	Core 2, Athlon 64	1p x 2 cores, 2.2 GHz, OoO, 32K/2MB L1/L2	135	849				
Mobl	Core 2 Mobile, Turion	1p x 2 cores, 2.0 GHz, OoO, 32K/2MB L1/L2	78	989				
Emb1	PA Semi, Emb. Athlon 64	1p x 2 cores, 1.2 GHz, OoO, 32K/1MB L1/L2	52	499				
Emb2	AMD Geode, VIA Eden-N	1p x 1 cores, 600MHz, inord.,32K/128K L1/L2	35	379				

### Cost Breakdown



### Infrastructure, Power and Cooling

### Performance/Cost

	Workload	Srvr2	Desk	Mobl	Emb1	Emb2
Perf	websearch	68%	36%	34%	24%	11%
	webmail	48%	19%	17%	11%	5%
	ytube	97%	92%	95%	86%	24%
	mapred-wc	93%	78%	72%	51%	12%
	mapred-wr	72%	70%	54%	48%	16%
	HMean	71%	42%	38%	27%	10%
Perf/Inf-\$	websearch	133%	139%	112%	175%	93%
	webmail	95%	72%	55%	83%	44%
	ytube	188%	358%	315%	629%	206%
	mapred-wc	181%	302%	241%	376%	101%
	mapred-wr	141%	272%	179%	350%	140%
	Hmean	139%	162%	125%	201%	91%
Perf/W	websearch	107%	90%	147%	157%	103%
	webmail	76%	47%	73%	75%	49%
	ytube	152%	233%	413%	566%	229%
	mapred-wc	146%	197%	315%	338%	113%
	mapred-wr	114%	177%	235%	315%	157%
	Hmean	112%	105%	164%	181%	101%
Perf/TCO-\$	websearch	120%	113%	124%	167%	97%
	webmail	86%	59%	62%	80%	46%
	ytube	171%	291%	351%	600%	215%
	mapred-wc	164%	246%	268%	359%	106%
	mapred-wr	128%	221%	200%	334%	147%
	Hmean	126%	132%	140%	192%	95%
(c) Per	formance,	cost a	and po	wer e	fficienc	ies

# Proposed 2-Level Memory Architecture



### (a) Memory blade architecture Reduce amount of DRAM by sharing est. 8% to 11% better perf/\$TCO

### Integrated Systems

### N1 is current te N2 adds 2-leve



# Discussion

### Svilen Kanev ISCA 15 Profiling a warehouse-scale computer



### Some Observations



top 50 hottest binaries only cover  $\approx$ 60% of WSC cycles.



Figure 3: Individual binaries are already optimized. Example binary without hotspots, and with a very flat execution profile.



Figure 2: Workloads are getting more diverse. Fraction of cycles spent in top 50 hottest binaries is decreasing.



Figure 4: 22-27% of WSC cycles are spent in different components of "datacenter tax".

### Nore Observations



### Benchmarks like SPEC are not representative

### Caches



### SuperScalar





# Multithreading



Figure 14: SMT effects on architectural behavior. From top to bottom: (i) more ILP extracted compared to Figure 12; (ii) frontend bound cycles decrease, but (iii) instruction starvation still exists; (iv) core throughput doubles with two hyperthreads.

# Conclusions

Finding	I
workload diversity	P
flat profiles	0
datacenter tax	D
	(t
large (growing)	I-
i-cache footprints	
bimodal ILP	N
low bandwidth	Т
utilization	D
latency-bound	V
performance	

### Summary of findings and suggestions for future investigation.

### nvestigation direction

Profiling across applications.

Optimize low-level system functions.

Datacenter specific SoCs

protobuf, RPC, compression HW).

-prefetchers, i/d-cache partitioning.

Not too "wimpy" cores.

Trade off memory bandwidth for cores. The provide the second seco

Vider SMT.

# Discussion