Addressing Performance and Programmability Challenges in Current and Future Supercomputers

Luiz DeRose Sr. Principal Engineer Programming Environments Director Cray Inc.

COM PUTE

STORE

ANALYZE



Outline

- Computer architecture and applications trends
- Programming Environment mission
- Tools focus for extreme scale computing
- Open issues
- Conclusions





Future Architectural Directions

Nodes are becoming more parallel

- More processors per node
- More threads per processor
- Vector lengths are getting longer
- Memory hierarchy is becoming more complex
- Scalar performance is not increasing and will start decreasing

As processors get faster, memory bandwidth cannot keep up, resulting in:

- More complex caches
- Non-uniform memory architecture (NUMA) for shared memory on node
- Operand alignment is becoming more important



Application Developers: Are you ready for the future?

- Today's petascale applications are not structured to take advantage of these architectures
 - Currently 80-90% of applications use a single level of parallelism
 - Message passing between cores of the MPP system
- Looking forward, application developers are faced with a significant task in preparing their applications for the future
 - You will have to move to a hybrid (MPI, threading, & vector) architecture
 - You must start considering how to **manage your arrays** to have them close to the computational engine when you need them
 - We are moving to a more complex memory hierarchy that will require user intervention to achieve good performance.

The Programming Environment Mission

- It is the role of the Programming Environment to close the gap between observed performance and achievable performance
- Support the application development life cycle by providing a tightly coupled environment with compilers, libraries, and tools that will hide the complexity of the system
 - Address issues of scale and complexity of HPC systems
 - Target ease of use with extended functionality and increased automation

COM PUTE



Extreme Scale Tools Focus

The current and future generation of tools should focus on

• Performance

- Help users maximize the cycles to the application
 - Address issues of scale and complexity of HPC systems
- Programmability
 - How do you get intuitive behavior and best performance with the least amount of effort
 - Provide the best environment to develop, debug, analyze, and optimize applications for production supercomputing
- Power
 - One of the biggest challenges facing extreme scale systems
 - Need to develop power-aware tools and techniques to aid the user to properly direct the system's use of power



What is needed to support the Performance, Programmability, and Power?

- Transitioning tools to assist the user in assessing and rebalancing an application for performance and power
- An environment that allows the developer to change applications a little at a time
- Provide support for user-directed tuning guidance
 for both performance and power
- Provide analysis and automation as we learn more



Creating Hybrid Codes, Why Should You Care?

- For the next decade (at least) all HPC systems will have the same basic architecture:
 - Communication between nodes
 - Application developers will continue to use MPI between nodes or sockets
 - Maybe SHMEM, UPC, Coarray
 - Shared memory programming paradigm on the node
 - MPI only will not do
 - Vectorization at the low level looping structures
 - SSE, AVX, ...
 - GPU, MIC, ...
 - etc
- Hybridizing a code gives many performance advantages
 - Resource contention (network, node memory, ...)
- While there is a potential acceptance of new languages for addressing multiple levels of parallelism, most developers cannot afford this approach until they are assured that the new language will be accepted and the generated code is within a reasonable performance range



Moving to a Hybrid Code: a Three-Task Approach

- Identify potential loops to accelerate
 - Determine where to add additional levels of parallelism
 - Assumes the MPI application is functioning correctly
 - Find top work-intensive loops
- Parallelize and vectorize identified loops
 - Split loop work among threads
 - Do parallel analysis and restructuring on targeted high level loops
- Add OpenMP (and then OpenACC directives)
 - Add parallel directives and acceleration extensions
 - Insert OpenMP directives
 - Verify application and check for performance improvements
 - Convert desired OpenMP directives to OpenACC

• We want a performance-portable application at the end

The Problem – How Do I Parallelize This Loop?

- How do I know this is a good loop to parallelize?
- What prevents me from parallelizing this loop?

COMPUTE

• Can I get help building a directive?

```
subroutine sweepz
do j = 1, js
do i = 1, isz
   radius = zxc(i+mypez*isz)
   theta = zyc(j+mypey*js)
   do m = 1, npez
    do k = 1, ks
    n = k + ks*(m-1) + 6
    r(n) = recv3(1,j,k,i,m)
     p(n) = recv3(2,j,k,i,m)
    u(n) = recv3(5,j,k,i,m)
    v(n) = recv3(3,j,k,i,m)
     w(n) = recv3(4,j,k,i,m)
     f(n) = recv3(6,j,k,i,m)
    enddo
   enddo
   call ppmlr
   do k = 1, kmax
     \mathbf{n} = \mathbf{k} + \mathbf{6}
    xa(n) = zza(k)
    dx (n) = zdz (k)
     xa0(n) = zza(k)
     dx0(n) = zdz(k)
     e (n) = p(n) / (r(n) * qamm) + 0.5 \&
        *(u(n) **2+v(n) **2+w(n) **2)
   enddo
   call ppmlr
enddo
enddo
```

subroutine ppmlr

```
call boundary
call flatten
call paraset(nmin-4, nmax+5, para, dx, xa)
call parabola(nmin-4,nmax+4,para,p,dp,p6,p1,flat)
call parabola(nmin-4,nmax+4, para,r,dr,r6,r1,flat)
call parabola(nmin-4,nmax+4,para,u,du,u6,u1,flat)
call states (pl,ul,rl,p6,u6,r6,dp,du,dr,plft,ulft,&
           rlft,prgh,urgh,rrgh)
call riemann(nmin-3,nmax+4,gam,prgh,urgh,rrgh,&
            plft,ulft,rlft pmid umid)
call remap <- contains more calls
call volume (nmin, nmax, ngeom, radius, xa, dx, dvol)
call remap <- contains more calls
return
End
```

ANALYZE

Luiz DeRose © 2013

STORF

Tools Needed When Creating Hybrid Codes

- Tools to design your application to be performanceportable across a wide range of systems
 - Application developers want to develop a single code that can run efficiently on multi-core nodes with or without an accelerator

Tools to understand your application

- Where is the time spent?
- Where are the key parallel structures (loops)?
- Where are the important arrays used?
- What prohibits parallelizing these structures?

• Tools to help with scoping analysis

• Identify shared, private and ambiguous variables

Simplifying the Task with Reveal



- Navigate to relevant loops to parallelize
- Identify parallelization and scoping issues
- Get feedback on issues down the call chain (shared reductions, etc.)
- Optionally insert parallel directives into source
- Validate scoping correctness on existing directives

ANALYZE

Visualize Loopmark with Performance Information



Luiz DeRose © 2013

STORE

ANALYZE

COMPUTE

Access Cray Compiler Message Information

0 0	X Reveal					
<u>F</u> ile <u>E</u> dit <u>∨</u> iew <u>H</u> elp		\varTheta 🔿 🔿 🔣 Explain				
🔻 vhone.pi 🛞		OPT_INFO: A loop starting at line %s was unrolled.				
Navigation ≪ Program View	Source - /lus/sonexion/heidi/reveal/sweepx2.f90	The compiler unrolled the loop. Unrolling creates a number of copies of the loop body. When unrolling an outer loop, the compiler attempts to fuse replicated inner loops - a transformation known as unroll-and-jam. The compiler will always employ the unroll-and-jam mode when unrolling an outer loop; iteral outer loop unrolling may occur when unrolling to satisify a user directive (pragma).				
▷ riemann.f90 ▷ remap.f90	L 32 dom = 1, npey	This message indicates that unroll-and-jam was performed with respect to the identified loop. A different message is issued when literal outer loop unrolling is performed, as this transfomation is far less likely to be				
	34 n = i + isy*(m-1) + 6 35 r(n) = recv2(1,k,i,j,m)	beneficial. For sake of illustration, the following contrasts unroll-and-jam with literal outer loop unrolling.				
▷ forces.f90 ▷ ppmlr.f90 ▷ states.f90 ▷ flatten.f90 ▷ sweepz.f90 ▷ sweepy.f90 ▷ prin.f90 ▽ sweepx2.f90 ▽ 0.53% □ Loop@28 □ Loop@32 □ Loop@33 □ Loop@58	<pre>36 p(n) = recv2(2,k,i,j,m) 37 u(n) = recv2(3,k,i,j,m) 38 v(n) = recv2(3,k,i,j,m) 39 w(n) = recv2(4,k,i,j,m) 40 f(n) = recv2(5,k,i,j,m) 40 f(n) = recv2(6,k,i,j,m) 41 enddo 42 enddo 43 v 44 do i = 1,imax 45 n = i + 6</pre>	# 426 "/ptmp/ulib/buildslaves/pdgcs-81-edition-build/tbs/build/release/pdgcs/pdgcs_ftn.msg. D0 J = 1,10 D0 I = 1,100 A(J,J) = B(J,J) + 42.0 ENDD0 D0 J = 1,102 D0 I = 1,100 A(J,J) = B(J,J) + 42.0 unroll-and-jam A(J,J+1) = B(J,J+1) + 42.0 ENDD0 D0 J = 1,100 A(J,J) = B(J,J) + 42.0 literal outer unroll ENDD0 D0 J = 1,100 A(J,J) = B(J,J) + 42.0 literal outer unroll ENDD0 D0 J = 1,100 A(J,J) = B(J,J) + 42.0 literal outer unroll ENDD0 D0 I = 1,100 A(J,J) = B(J,J+1) + 42.0 ENDD0 The literal outer unroll code performs the same sequence of memory operations as the original nest, while the unroll-and-jam transformation interleaves operations from outer loop iterations. The compiler employs literal				
vhone.pl loaded. vhone_loops.ap2	Access integrated	flow impediment, prevent fusion of the replicated inner loops. Literal outer loop unrolling is generally not desirable. It is provided to ensure expected behavior and for those rare instances where the user has determined that it is beneficial.				
	message 'explain' support by right clicking on message	Explain other message X Close				
	COMPUTE STORE	ANALYZE				

View Pseudo Code for Inlined Functions

COMPUTE



VI-HPS - SC'13

Luiz DeRose © 2013

STORE

ANALYZE

(15

Scoping Assistance – Review Scoping Results



Reveal Gives Feedback on Scoping Results

Variable from inlining – hover over 'l' to see what symbol means

			Call or I/O at line 53 of sweepx1.f90
Name	Type	Scope	Info
	-		WARN: LastPrivate of array may be very expensive.
radius I	Scalar	Private	FAIL: incompatable with 'natural' scope.
u	Array	Private	FAIL: incompatable with 'natural' scope.
			WARN: LastPrivate of array may be very expensive.
v	Array	Private	FAIL: incompatable with 'natural' scope.
			WARN: LastPrivate of array may be very expensive.
w	Array	Private	FAIL: incompatable with 'natural' scope.
			WARN: LastPrivate of array may be very expensive.
xa	Array	Private	FAIL: incompatable with 'natural' scope.
			WARN: LastPrivate of array may be very expensive.
xa0	Array	Private	FAIL: incompatable with 'natural' scope.
			WARN: LastPrivate of array may be very expensive.
gamm	Scalar	Shared	
js	Scalar	Shared	
ks	Scalar	Shared	
ngeomx	Scalar	Shared	
niettx	Scalar	Shared	
nrightx	scalar	Shared	
senai	Array	Shared	
svel KI	Scalar	Shared	Reduction Performance Performa
Enabl	le FirstPri	vate	Nene
Enabl	le LastPri	vate	None
		vate	
Find Nam	ne:		

Reduction variable from down the call stack - hover over 'RI' to see what symbol means

Scoping Assistance – Reveal Generates Directive



Use Reveal to Validate User Inserted Directives

00	X Reveal	
<u>File E</u> dit <u>V</u> iew <u>H</u> elp		
🔻 vhone.pi 🔞	Liser inserted	
Navigation	Source - /ufs/home/users/heidi/reveal/riemann.f90	
🔺 Program View 🗘 🏟		/
▼ boundary.f90	R 64 !\$OMP parallel do default(none)	
b dtcon t90	65 !\$OMP& private (n)	
▶ dump.f90	66 !\$OMP& shared (lmin,lmax,prgh,urgh,vrgh,plft,ulft,vlft,pmid,clft, &	
≥ evolve.f90	67 !\$OMP& crgh,l,plfti,pmold,prghi,umidl,umidr,wlft,wrgh,zlft, &	
▶ flatten.f90	68 !\$OMP& zrgh,gamfac1,gamfac2)	
⊳ forces.f90	LSg 69 do l = lmin, lmax	
▶ images.f90	r L 70 do n = 1, 12	
▶ init.f90	71 pmold(1) = pmid(1)	
▶ parabola.f90	$\frac{7}{72} \text{wlft} (\mathbf{l}) = 1.0 + \text{gamfacl*}(\text{nmid}(\mathbf{l}) - \text{nl} = 0.0 + 0.000 + 0.000 + 0.000 + 0.000 + 0.000 + 0.00000 + 0.00000 + 0.00000 + 0.0000 + 0.00000 + 0.00000 + 0.00000 + 0.00000 + 0.00000 + 0.00000 + 0.00000 + 0.0000 + 0.0000 + 0.0000 + 0.0000 + 0.0000 + 0.0000 + 0.0000 + 0.00000 + 0.00000 + 0.00000000$	
⊳ ppmlr.f90	$\frac{73}{73} \text{ wrgh (1)} = 1.0 + \operatorname{gamfacl*(pmid(1))} - \operatorname{pr}$	
⊳ prin.f90	74 wift (1) = clft(1) * scrt(wlft(1))	
▶ remap.f90	75 wreb (1) = creb (1) * sert (wreb (1))	
➡ riemann.f90	$75 \qquad \text{wigh}(\mathbf{x}) = c \log(\mathbf{x}) \cdot s \operatorname{squ}(\mathbf{x}) \operatorname{squ}(\mathbf{x})$	
▼ RIEMANN		
Loop@44	77 21 gir (1) = 4.0 * Vigir (1) * Wigir (1) * Wi Scalar Private WARN Scope does not agree with user OMP directive.	
Loop@69	$\frac{78}{10} = -2(1) + w(1)(\frac{1}{2}) + w(1)(\frac{1}{2}) + Scalar = Private$	
Loop@70	79 Zrgn (t) = Zrgn (t) * wrgn (t) (zrgn (t) cm Array Shared	
Loop@89	$unid(\mathbf{t}) = urt(\mathbf{t}) - (pnid(\mathbf{t}) - pirt(\mathbf{t}) range)$	
▶ states.190	81 $\operatorname{umin}(\mathbf{L}) = \operatorname{urgn}(\mathbf{L}) + (\operatorname{pmin}(\mathbf{L})) - \operatorname{prgn}(\mathbf{L})$ scalar sc	
▷ sweepx1.f90	82 pmid $(\mathbf{l}) = pmid(\mathbf{l}) + (umidr(\mathbf{l})) - umid gamaz Scalar Shared$	
▷ sweepx2.f90	83 pmid () = max(smallp,pmid()) imax scalar shared	
▶ sweepy.f90	84 1† (abs(pmid(l))-pmold(l))/pmid(l) < to imin Scalar Shared	
▶ sweepz.f90	H 85 enddo pin Array Shared	
vh1mods.f90	R6 enddo più Array Shared	
▶ vhone.f90	mini printa Array Stated	
▶ volume.190	clafo - Line 69	
zonemod.190	A loop starting at line 69 was not vectorized for an unspecified reason. Enable First Private	
	A loop starting at line 69 was partitioned.	0
	Search:	
vhone.pl loaded	Insert Directive Show Directive	Close

Moving to a Hybrid Code: a Three-Task Approach

- Identify potential loops to accelerate
 - Determine where to add additional levels of parallelism
 - Assumes the MPI application is functioning correctly
 - Find top work-intensive loops
 - Performance tools + compiler loop work estimates
- Parallelize and vectorize identified loops
 - Split loop work among threads
 - Do parallel analysis and restructuring on targeted high level loops
 - Reveal with compiler feedback and source code browsing
- Add OpenMP (and then OpenACC directives)
 - Add parallel directives and acceleration extensions
 - Insert OpenMP directives
 - Reveal with compiler scoping assistance

COMPUTE

- Verify application and check for performance improvements
 - Convert desired OpenMP directives to OpenACC

• We want a performance-portable application at the end

Luiz DeRose © 2013

STORE

ANALYZE

Debugging on Extreme Scale Systems

- Systems with thousands of threads of execution need a new debugging paradigm
- Need to build tools around traditional debuggers with innovative techniques for productivity and scalability
 - Support for traditional control-centric debugging mechanism
 - STAT Stack Trace Analysis Tool
 - MRNet based scalable generation of a single, Wisconsin merged, stack backtrace tree
 - ATP Abnormal Termination Processing
 - Scalable analysis of a sick application, delivering a STAT tree and a minimal, comprehensive, core file set.

Comparative debugging

VI-HPS - SC'13

- Collaboration with University Queensland
- A data-centric paradigm
 - Ability to see data from multiple program executions in the same instance

Luiz DeRose © 2013





ANALYZE

Stack Trace Analysis Tool (STAT)

- Stack trace sampling and analysis for large scale applications
 - Sample application stack traces
 - Scalable generation of a single, merged, stack backtrace tree
 - A comprehensible view of the entire application
 - Discover equivalent process behavior
 - Group similar processes
 - Reduce number of tasks to debug

• Merge/analyze traces:

- Facilitate scalable analysis/data presentation
- Multiple traces over space or time
- Create call graph prefix tree
 - Compressed representation
 - Scalable visualization & analysis

ATP - Automatic Termination Processing

COMPUTE

• System of light weight back-end monitor processes on compute nodes

ANALYZE

Leap into action on any application process trapping

Luiz DeRose © 2013

STORE



WISCONSIN

10798:[0,3-10799]1:[1] 1:[2] MPI Waitall do SendOrStall 10798:[0,3-10799] 1:[2] MPIDI CRAY Progress wait MPIDI CRAY SMPClus Barrier 10798:[0,3-10799] 1:[2] MPIDI_CRAY_progress MPIR Barriel 0798:[0,3-10799] \18:[133,496,502,...] 1:[2] MPIDI_CRAY_Progress_wait MPIDI CRAY ptldev progress

main

Comparative Debugger

- Helps the programmer locate errors in the program by observing the divergence in key data structures as the programs are executing
 - Allows comparison of a "suspect" program against a "reference" code using assertions
 - Simultaneous execution of both
 - Ability to assert the match of data at given points in execution
 - Focus on data not state and internal operations
 - Narrow down problem without massive thread study
 - Data comparison
 - Tolerance control nobody expect it to be perfect
 - Array subsets correlate serial to parallel bits
 - Array index permutation loops rearranged
 - Automated asserts let it run until a problem is found
 - Forcing correct values continue on with correct data



Cray Comparative Debugger

000	0		X Cr	ay Comparativ	ve Debugg	er (CCDB)						
<u>F</u> ile	<u>V</u> iew <u>T</u> oo	ols <u>H</u> elp										
	2 💽 💽 🤇		Focus: all	all								
?		Applica	tion-0 Status				Application-1 Status					
App0{015} Stopped () driver.f:106						App1{015} Stopped () sweep.f:231						
Output	Breakpt +	driver.f		[↓] X	Output	Breakpt +	sweep f	I IX				
04	mm db	a = 1			200	- Teampt	$\frac{1}{1}$					
94:	endif	c – 1			220:		$k_0 = 1 + (k_k - 1)^{-1} k_k$ $k_1 = \min(k_0 + mk - 1 + k_1)$					
96:	if (kbc	.ne.0) then			222:		nk = k1 - k0 + 1					
97:	it_kb	c = it			223:		else					
98:	_ jt_kb	c = jt			224:		k0 = kt - (kk - 1) * mk					
99:	mm_kb	c = mm			225:		k1 = max (k0-mk+1,1)					
100:	else				226:		nk = k0 - k1 + 1					
101:	it_kb	oc = 1			227:		endif					
102:	jt_kb	oc = 1			228:							
103:	mm_kb	pc = 1			229: !	this could	be *nk* instead of *mk* if all ph	i{i,j}{b,				
104:	endif				230: !	were dimen	sioned with mmi as the second dime	nsion				
105:							nib = jt*mk*mmi					
	if (myi	d .eq. 1) then			232:		njb = it*mk*mmi					
107:	prin	it *, SWEEP3D -	Method 5 - ,	D	233:	T 1 C 1						
108:	& number	· Pipeline	a wavefront with Line	-Recurs:	234: c	1-inflows	for block (1=10 boundary)					
1109.	prin. prin	nt 100 isn isc	t mm pm it a it a kt		235. 0		if (ew row ne 0) then					
111	100 form	nat(' S'.i1.'P'	.i1.3x.'-'.1x.i2.' an	ales/oct	230.		call rev real(ew rev, philb,)	nib. ew t				
112:	& i	3.' moments'/.	, , , , ,		238:		else					
113:	، ع	global grid:	', i5, 'x', i5, 'x', i5)	∇	239:		if (i2.1t.0 .or. ibc.eq.0) the	en 🔽				
		-										
Console	<pre>App1{015} App1{015}</pre>	: Temporary break : Breakpoint 1.	kpoint 1: file sweep.f, 1 sweep at sweep.f:231	line 231.				$ \Delta $				
Output		*· ·						2				

Luiz DeRose © 2013

STORE

ANALYZE

COM PUTE

CCDB - Comparison

00	0		X CC	DB Comparison					
?	Application-0 Application-1								
	♦ Арр0{015} 🧿	sweep.f:160			5} 💽 sweep.f:160				
	Name or Expression	Туре		Results	Type Template	App-0 Decomp	App-1 Decomp	Op Eps	
	jkm	INTEGER*4	Click to se	e results				== e	
	mio	INTEGER*4	Click to se	ee results				== e	
	nk	INTEGER*4	Click to se	e results				== e	
	ti	REAL*8	Click to se	e results				== e	
	tj	REAL*8	Click to se	e results				== e	
	weta	REAL*8 (6)	Click to se	e results		None	None	== e	
	wmu	REAL*8 (6)	Click to se	e results		None	None	== e	
									7
Co	mpare Add Comparison Resu	ılt Filter: 🔷 all 🔌	🕨 fail 🔷 v	varn 🔷 pass 🔷 to	odo				Close

COMPUTE | STORE | ANALYZE

CCDB – Comparison

00	0		X CC	DB Comparison					
?	Q	Application-1							
App0{015} Sweep.f:160				App1{015} Sweep.f:160					Up Down
	Name or Expression	Туре		Results	Type Template	App-0 Decomp	App-1 Decomp	Op Eps	
	jkm	INTEGER*4	Click to se	e results				== e	
	mio	INTEGER*4	Click 1					== e	
	nk	INTEGER*4	Click 1					== e	
	ti	REAL*8	Click 1	ppu(U): Compariso	on is tru	e.		== e	
	tj	REAL*8	Click 1	ppu{1/}: Compan	ison is i	folgo		== e	
	weta	REAL*8 (6)	Click 1	ppu(ois): Compa	irison is	laise.	None	== e	
	wmu	REAL*8 (6)	Click 1		IS:		None	== e	
			A	op0{17}: 64 op0{815}: 64					V
	(
Co	mpare Add Comparison F	Result Filter: 🔷 all 🔌	🕨 fail 🔷 v	warn 🔷 pass 🔷 to	odo				Close

VI-HPS - SC'13

COM PUTE

Luiz DeRose © 2013

STORE

ANALYZE

N

Why Application Level Power Management?

- Power consumption is one of the biggest challenges facing extreme scale systems
 - Scaling up from today's requirements for a petaflop computer is not a viable solution
 - Sites are already increasingly constrained by power & cooling limitations as well as cost of system power and cooling

• Hardware could attempt to manage power automatically

 This strategy would most likely be reactive and certainly be suboptimal

• Power management must be proactive

- Software can help by monitoring and proactively managing how power is being used in the system
 - The programmer must understand the tradeoffs between power and performance

• More research and development is needed for

- Power measurement and monitoring tools
- Power analysis

VI-HPS - SC'13

• Management of power consumption

Luiz DeRose © 2013

ANALYZE

Power Analysis

Ultimate goal is to develop power-aware tools and techniques

- Provide suggestions of possible transformations or implementations of the application with different performance/power profiles
- Aid the user to properly direct the system's use of power
 - Using transparent automation where possible
 - Transparently adopting environmental settings to achieve a desired power reduction and energy savings with a predetermined impact to performance
- Can leverage the work done in automatic performance analysis

Challenge

 Ability to identify and correctly attribute the power weights to the individual "operation" components and map back to the application

Power Measurement Wish List

Per H/W core or memory component

- How much does each H/W attribute (L1 cache, L2, cache, TLB, etc) draw?
- Get relative power usage of components
- Get information on joules per operator or per access
 - Example: how many joules does a DCACHE hit or miss use per reference?
- Provide core utilization metrics

• Overall program power measurement

- Similar to wallclock (total watts or joules)
- Multiple power meters
 - 1 for core
 - 1 for memory system

• Correlating HWPCs to power usage in performance tools

- Can be directly or indirectly
 - Roughly measured with existing H/W performance counters

Adaptive Power Management

Statically managed

- Compiler support
 - Compile-time options to direct optimization for maximum or minimum
 - e.g., -P 1, 2, 3 like current -O 1, 2, 3
 - Compiler directives
 - Put this data there and keep it there until I tell you to release it
- Application-selected power modes/algorithms

Runtime managed

- Based on compiler hints
 - Which components are about to be utilized more or less heavily
 - Integer algorithm coming up
 - Intensive floating point region coming up
 - Poor locality coming up, turn off cache use for this next segment
- Auto-tuning
 - Auto-tuning for performance and for power (and a combination of both)
- User Guided
 - Runtime assertions/options/hints

Outline

- Computer architecture and applications trends
- Programming Environment mission
- Tools focus for extreme scale computing
- Open Issues
- Conclusions

How About Resilience?

- On extreme scale systems the probability of a node failure is substantial
- Extreme-scale applications will have built-in resilience support
 - There will be support for detection of failure and information about which images failed
 - e.g., error return for synchronization and collective operations can indicated that an image has failed
- Will tools be ready for resilient applications?
- Will tools be resilient themselves?

Other Issues to Think About

- Running on an extreme-scale environment
 - Target scalability issues in all areas of tool development

Volume of data

- Need to rethink the performance measurement approach
 - Can we focus only on things that are different?
 - What kind of runtime aggregation will be possible?
 - How about the "Heisenbug" effect?
- What can you do to prep the tools?
 - Availability of extreme-scale systems will be limited!
- Cost/benefit estimates (porting apps)
 - Can you provide an estimate to the user of the improvements that certain transformations will bring?
- Ease of use
 - Automatic program instrumentation
 - Automatic analysis



Access light version of performance tools software

> module load perftools-lite

Build program

> make



a.out (instrumented program)

ANALYZE

Run program (no modification to batch script)

COMPUTE



Luiz DeRose © 2013

STORE

Conclusions

- Extreme scale systems will be more parallel, with more processors per node, more threads per processor, longer vector lengths, more complex memory hierarchies, and ...
- Application developers will need sophisticated tools and adaptive runtime systems to help them address issues of scale and complexity of these systems
- Extreme scale tools should address issues of
 - Performance
 - Help users maximize the cycles to the application
 - Programmability
 - Ease developing, porting, and tuning efforts
 - Power
 - Help the user understand the tradeoffs between power and performance, and tune for both

