

# Challenges of Sustained Petascale Computing

## SOS11 – Key West – June 12, 2007



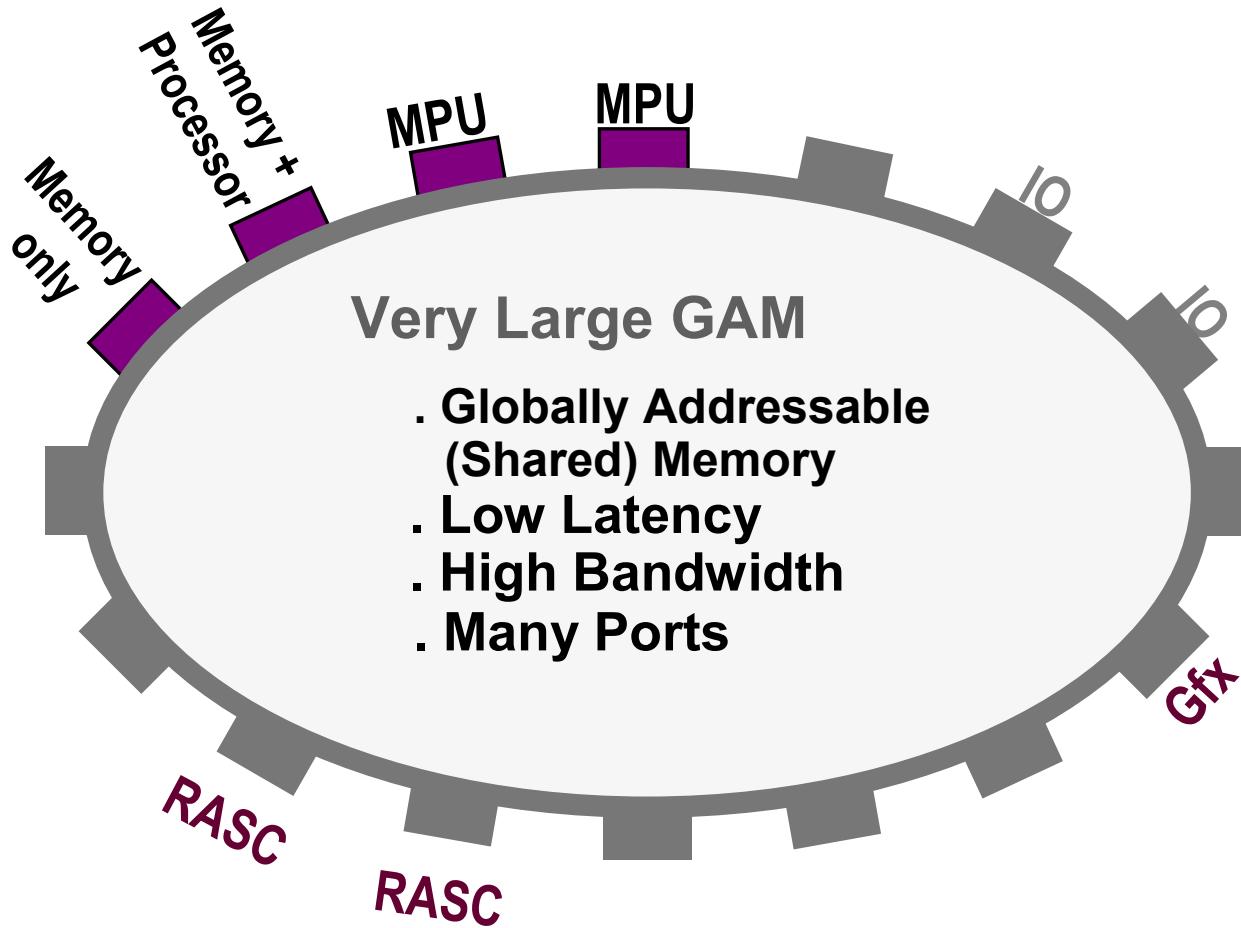
Jean-Pierre Panziera  
Chief engineer, Applications group

# Challenges of Petascale Computing

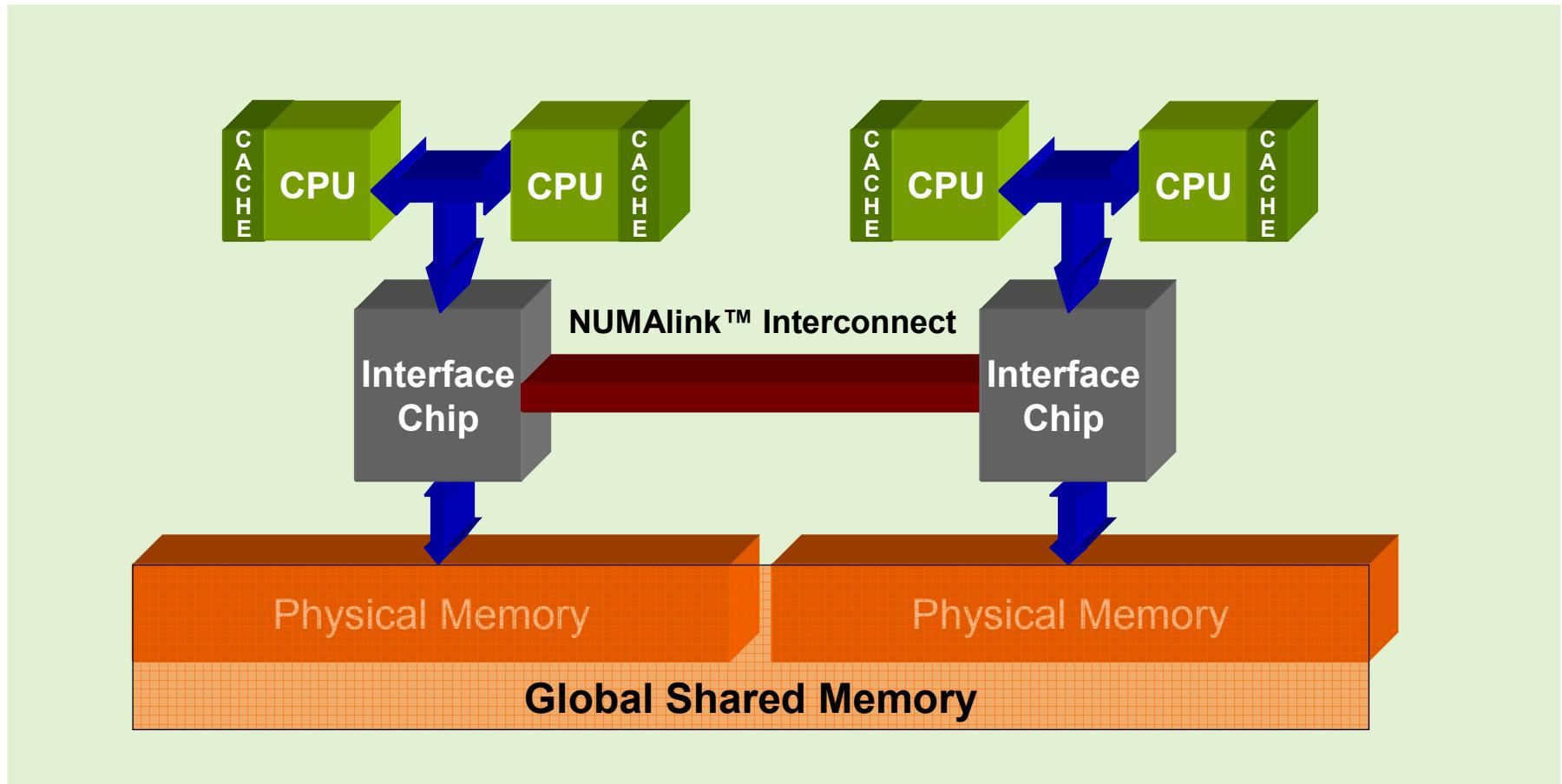
- the Globally Addressable Memory (GAM) concept
- NUMAflex architecture evolution
- a current configuration: LRZ
- Scaling HW to Petascale
- Reconfigurable Application Specific Computing (RASC)
- Petascale system software
- Petascale applications



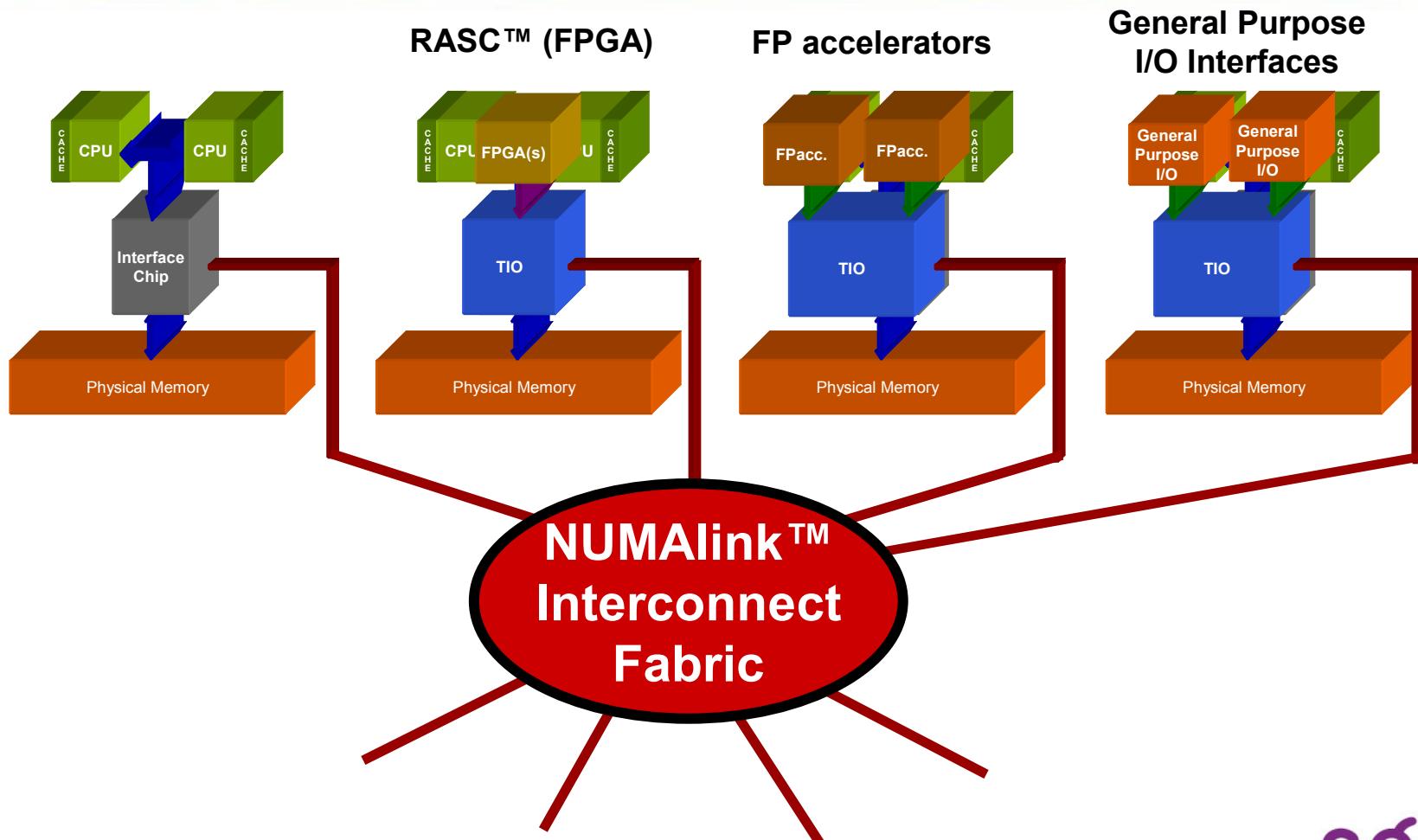
# Globally Addressable Memory (GAM) Data Centric Architecture



# ALTIX scalable system infrastructure



# ALTIX scalable system infrastructure a global shared memory space + IO and accelerator nodes



# NUMAflex architecture evolution

year	system	max SSI procs/cores	max GAM procs/cores	Memory capacity
1996	Origin2000	256 (512)	-	512 GB
2000	Origin3000	512 (1024)	-	1 TB
2003	Altix3000	512	2 048	32 TB
2006	Altix4000	4 096 (1024)	16 384	196 TB
-	UV	petascale	petascale	petascale



**sgi**<sup>®</sup>  
INNOVATION  
FOR RESULTS<sup>™</sup>

# LRZ\_2

128 racks, 4096 nodes, 9728 cores = 19 partitions x 512  
38 TB memory, 4GB/core  
300 TB disks - CXFS

2 816 routers  
14 082 links

( 2048 in backplane + 766 routers)  
( 9216 in backplane + 4866 cables)

62 TFlops (peak)  
43 TB/s aggregate local memory bandwidth  
409 GB/s bisection bandwidth

petascale? → x20 !!! (current generation)  
x2- 4 with next generation

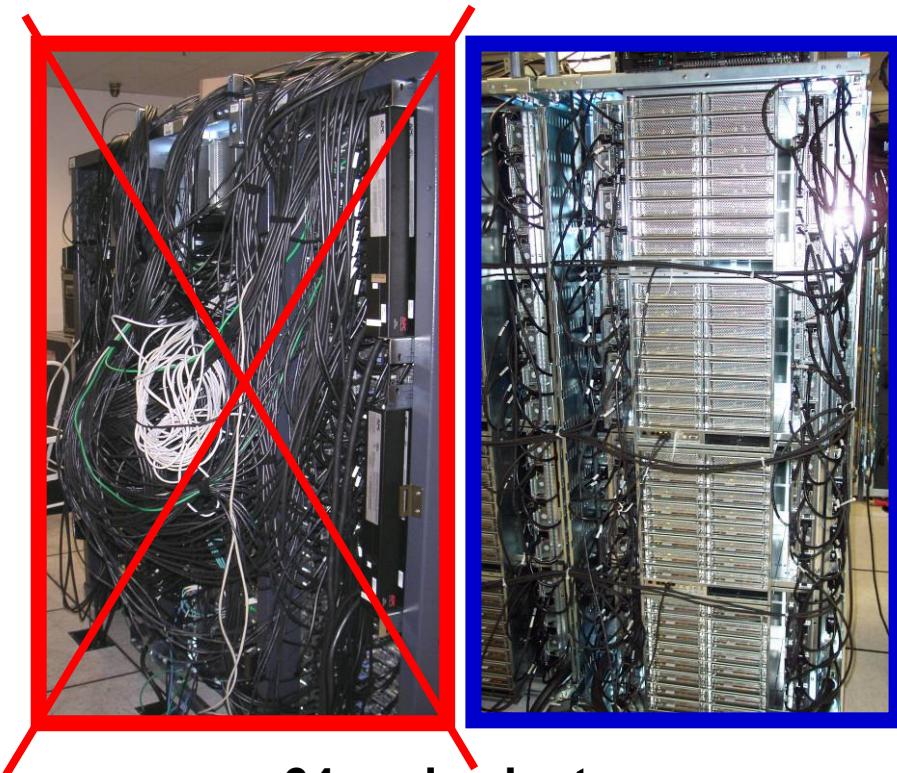
May 2007



# Hardware Challenges of Sustained Petascale Computing

- Scale to more nodes/cores/memory → NUMAflex architecture
- Scale memory bandwidth → Bytes/Flop ratio
- Scale interconnection bandwidth → NUMAlink enhancements
  
- System packaging, denser, tighter, easier to deploy
- Power management, efficient (water) Cooling,
- Integrate Accelerators (FPGA, FP accelerators, ...)
- Resilient systems, improved RAS features

# Packaging: rack optimized cabling and water cooling

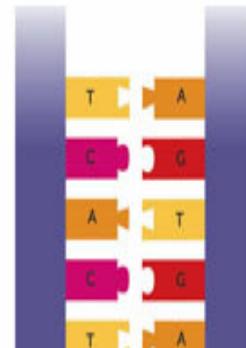
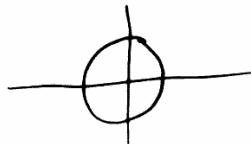


Water cooling

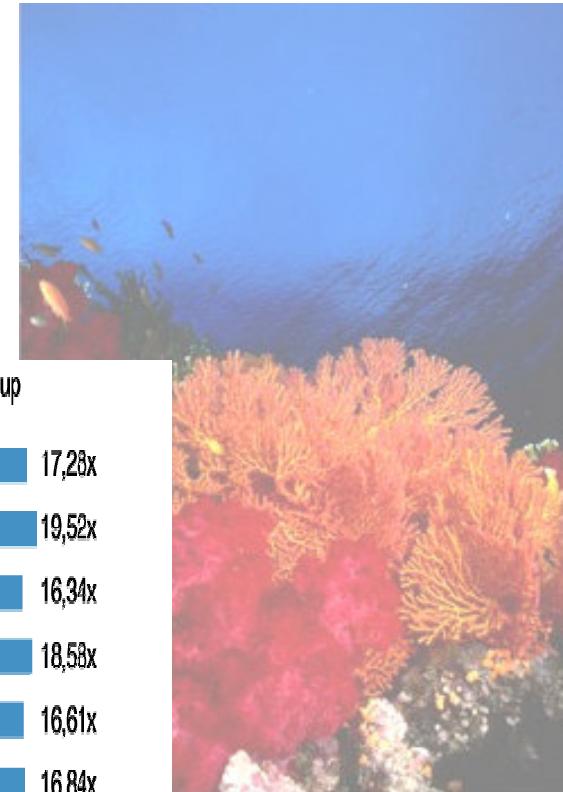
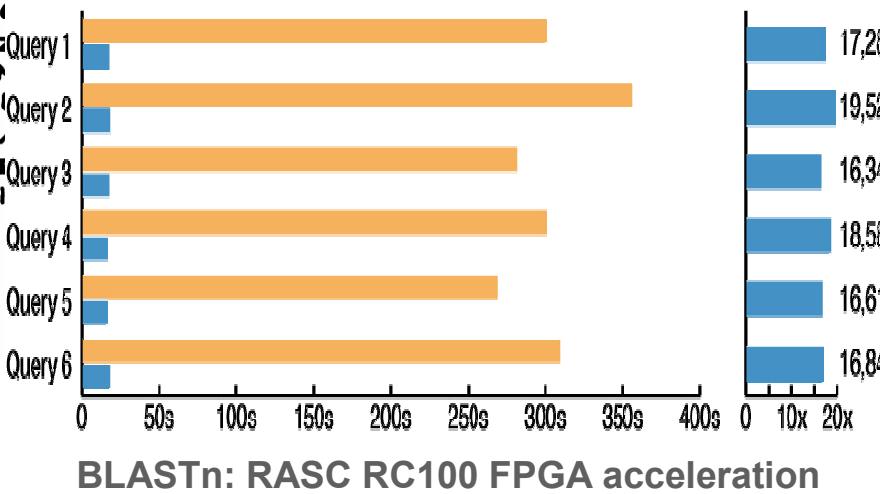
Heat Rejection  
95% water / 05% air

# RASC/FPGA accelerators encryption, BLASTn, jpeg, ...

H E R > 9 J A V P K I O L T G O D  
N 9 + B φ □ O □ D W Y . < □ K F □  
B Y □ C M + U Z G W φ □ L □ H J  
S 9 9 □ A □ L □ A □ V 9 0 + + R K 0  
□ □ M + □ T □ D 1 □ F P + P □ K /  
9 □ R □ F □ L □ O - □ O C □ F □ > □ O □  
□ □ + K □ □ I □ 0 4 □ X 6 V □ L I  
φ G □ J □ T □ O + □ N Y □ + □ L □  
0 < M + 8 + Z R □ F B □ X A 0 0 K  
- □ L □ U V + □ A J + □ O 9 A < F B Y -  
U + R / □ L E I D Y B 9 8 T M K O  
□ < □ J R J I □ □ T □ M □ 1  
□ □ A S □ □ + N I □ □ F B □  
J G F N A □ □ □ □ □ □ □ □ □ □ □ □  
Y B X □ □ I □ □ C E □ □ V U □  
I □ □ O □ □ B K □ □ 0 9 1 □ □  
R □ T + L □ □ C □ □ + F J W □  
+ + □ W c □ □ W □ □ P O S H T □  
I F K □ □ W □ □ A □ □ B □ □ Y □ □ O □ □  
> M D H N □ □ S □ □ Z □ □ A □ □

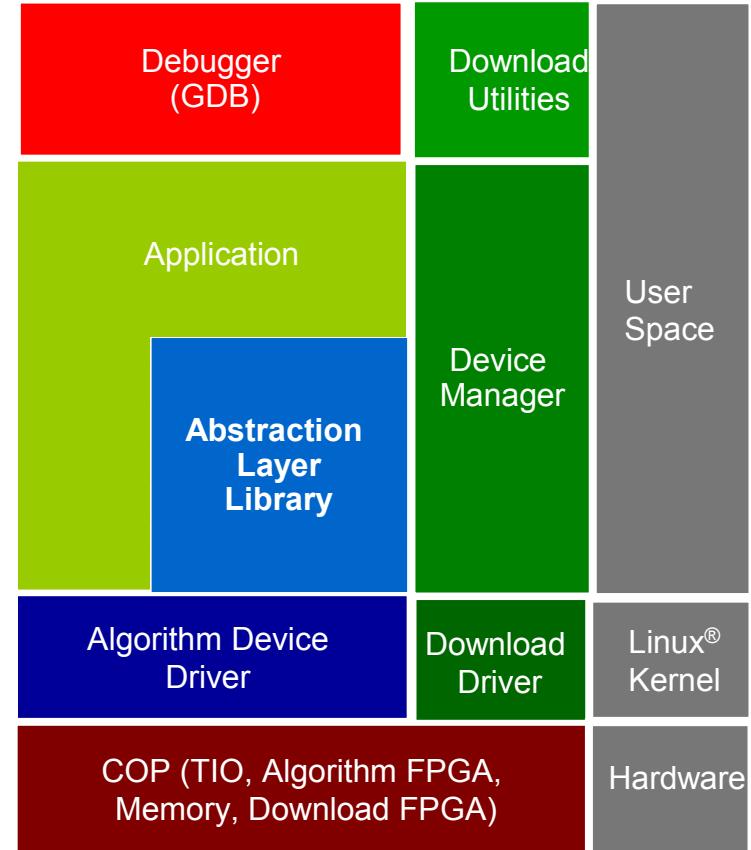


Runtime (s) CPU CPU+FPGA



# SGI® RASC™ Integrated Solution Stack Simplifies Development

Tool	Function
FPGA Aware Version of the Gnu Debugger (GDB)	Simultaneously debugs both the CPU based app and the FPGA accelerated app
RASC Abstraction Layer (RASCAL)	Enables serial or parallel FPGA scaling
RASC API and Core Services Library	Provides tools to develop reconfigurable computing elements in a multi-user, multiprocessing environment
3rd Party HLL Development Tool Support	Fully integrated tools including Celoxica Handel-C and DK Design Suite, Mitronics Mitrion C, Impulse Impulse-C (in process)



# RASC challenges

## FPGAs, FP accelerators, ...

- No standard HW interface
- No standard SW interface (IO library)
- Managing Accelerator private memory
- Low level programming (VHDL/ assembler language)
- No standard high level compiler
- Limited code size (kernel)
- FP format, FP double precision
- Data Error Protection

# Software Challenges of Sustained Petascale Computing

- Scale to more nodes/cores/memory → NUMAflex architecture
- Scale memory bandwidth → Bytes/Flop ratio
- Scale interconnection bandwidth → NUMAlink enhancements
- Integrate all types of nodes: CPU+Memory, Memory-only, IO, RASC, gfx, ...within an SSI partition
- Integrate nodes/services across partitions
- HW Optimized Middleware
- Resilient systems, improved RAS features
- Peta-scaling applications

# System software environment

- Linux SLES10 & RHELv5 out of the box
  - Linux scales up to 1024p (SLES10) ... “constellation” configurations (large partitions)
  - XFS file system
- SGI ProPack 5
  - Linux enhancements ( numatools, FFIO, cpuset, XVM)
  - Cross-partitions cluster-wide services (PCP, Array services, MPI, shmem)
- CXFS Clustered file system
- DMF, TMF, openVault

# RAS features

- Extensive components testing (DIMMs screening)
- Full System staging
- Error Prevention: memory protection, redundant data and IO paths, power supplies, cooling fans
- Error Detection: enhanced monitoring, system notification
- Error Containment: proc isolation, at worst within 1 partition only
- Maximizing serviceability: hot swappable components
- Fast checkpoint/restart at application level
- Automatic system reconfiguration

# Peta-scale Application Challenges

- Requiring new levels of parallelism 1,000s → 100,000s threads (MPI?)
- Taking advantage of multi-cores CPUs (SMP+MPI?)
- Explicit Parallel languages: UPC, co-array Fortran ?
- Integrating FPGA / FP accelerators (compilers?)
- Peta-Debugging
- Peta-Performance tools
- Minimal development effort

# Conclusion

- Sustained peta-scale computing presents formidable challenges for HW , system SW and Applications.
- SGI NUMAflex architecture provides a base for peta-scale systems
- SGI next generation UV systems will be peta-scale
- UV software environment will enable sustained peta-scale computing Applications

