# ETNUS
## TOTALVIEW

Blue Gene/L
Workshop

Reno, NV

October, 2003

# TotalView on Blue Gene/L

John DelSignore

Chief Architect

# Outline of Today's Talk

- **About Etnus TotalView**
- **Porting TotalView to BG/L**
- **Scalability and Performance**
- **Do We Need a Paradigm Shift?**
- **Questions & Answers**

# About Etnus TotalView

- **World's leader in UNIX/Linux debugging solutions**
- **Parallel Processing, SMP and Distributed Systems**
  - Multi-process and multi-thread support
  - Distributed, clusters, MPI, pthreads, OpenMP
- **Customers include**
  - Major government research and academic labs worldwide
  - Software development companies
  - Computer hardware vendors
  - Companies in Finance, Entertainment, Telecommunications, Energy, Aerospace, Climate Modelling, Automotive
- **TotalView has been continuously developed and extended for 17 years**
- **We're proud to be the ASCI debugger of choice**

# Platforms

- **Present**
  - IBM AIX Power (RS6000, SP)
  - HP/Compaq Tru64 Alpha
  - Intel Linux X86 and IA64
  - SGI IRIX MIPS
  - Sun Solaris SPARC
- **Future**
  - AMD Linux X86_64
  - HP HPUX IA64
- **Past**
  - Compaq Linux Alpha
  - HP HP-UX PA-RISC
  - Many, many others
- **Fujitsu, NEC, Cray, Hitachi**
- **Red Storm, Blue Gene/L**

- **C/C++**
- **Fortran 77/90/95**
  - F90 types
  - Modules
  - By-descriptor arrays
- **UPC**
- **Mixed Languages**
- **Assembly**
- **Mixed Java/C/C++**
  - with the CodeRoad JNI Bridge

# Supports Leading Compilers

- IBM XL/VA
- GCC 3
- SGI MIPS Pro
- Sun ONE Studio
- Intel C/C++ for Linux
- Intel Fortran for Linux

- Intel KCC
- Intel KAI Guide
- PGI 3
- Lahey/Fujitsu
- Apogee
- HP/Compaq C/C++/F90
- Compaq UPC 2.0

# TotalView GUI & CLI



Root Window



Variable (Process Laminated)



Process WIndow

# TotalView GUI & CLI (cont)

Subset Attach

Message Queue Graph

CLI

- **Etnus is porting TotalView to BG/L**
- **Working with IBM to**
  - Iron out the details of *how* to do it
  - Collaborating on debugging interfaces
- **Remainder of this talk will outline**
  - General porting approach we plan to use
  - Outline of the functionality we plan to have
  - Scalability and performance ideas
- **Details are subject to change**
  - At participating dealers only, your mileage may vary, offer not valid in all states, some restrictions apply ☺.

# BG/L TotalView Debugger Software Stack

# Front-End Node (FEN)

**TotalView**

**ICCDP client**     **TCP**

**Automatic Process Acquisition?**

**srun**

**Linux**

**4-16 x Power4**

**Front End**

◆ **Requires a powerful FEN**
- ◆ Server-class SMP (4-16 way)
- ◆ Power4, 1GHz+
- ◆ 4GB+ memory

◆ **Runs Linux OS**

◆ **Hosts the TotalView GUI and CLI clients**

◆ **TV talks client-side ICCDP to *N* TV debugger servers over TCP/IP**

◆ **Automatic process acquisition TBD**
- ◆ May use existing approach (TV debugs mpirun)
- ◆ May need to develop a new client/server approach (requires LLNL/IBM/Etnus collaboration)

# I/O Nodes

**tvdsvr**

**ICCDP server**

**Milestone #26**

**CIOD** **Tree**

**Linux**

**2 x PPC 440**

**I/O**

- **Reasonable horsepower for the TotalView Debugger Server (tvdsvr)**
  - 700 MHz PPC 440
  - 512 MB memory
- **Runs Linux OS**
- **Hosts the BG/L tvdsvr**
- **tvdsvr talks server-side ICCDP over TCP/IP to TV**
- **1 server exerts trace-level control over 64 compute nodes**
- **Communicates with CIOD via pipe**
  - To implement low-level debugging protocol
  - CIOD debug messages sent to/from CIOD agent

# Compute Nodes

- **Hardware**
  - 2 x 700 MHz PPC 440
  - Double Hummer FPU (1 per processor)
  - 512 MB memory
- **Runs CNK OS**
- **No part of TotalView runs here**
  - So IBM knows better how this really works
- **Debugging happens in CNK exception handlers ("CIOD agent" is my name)**
  - Communicates with CIOD on Linux I/O node
  - Sends/receives debug messages over the tree
- **TotalView will view each compute node as**
  - A single process (one address space)
  - Having two threads (one thread per processor)

CIOD agent

MPI ranks

CNK

2 x PPC 440

Compute

# BG/L TotalView Features

- **TotalView version 6.4 or later (mainline)**
- **TotalView GUI and CLI**
  - Breakpoints, single-step, etc.
- **Documentation**
  - Electronic form, HTML help
  - Blue Gene/L specific addenda
- **Languages**
  - C, C++, Fortran
  - Assembler
  - Mixed languages

- **Compilers**
  - GCC 3
  - IBM XL / Visual Age
- **MPI**
  - Automatic process pickup
  - Message queue display
- **Double Hummer FPU**
- **Data watchpoints?**
  - If supported by CNK
- **STL types display?**
  - vector<>, list<>, map<>

# Unsupported Features

- **No data Visualizer**
- **No OpenMP**
- **No SHMEM, PVM**
- **No pthreads**
- **No shared libraries**
- **No compiled expressions**
  - Interpreted expressions only
- **No checkpoint restart**
- **No core files**

# Scalability and Performance

- TotalView's History of Scaling
- Defense Mechanisms for Scaling
- Scalability and Performance Philosophy
- Plans for Scaling to BG/L
- Do We Need a Paradigm Shift?

- **TotalView was designed to be a parallel debugger (BBN Butterfly)**
- **10 years ago, scaled to about 100 processes comfortably**
- **ASCI Path Forward projects**
  - Developed Subset Attach feature
  - Scalability was increased to
    - Handle about 1,000 processes comfortably
    - Handle about 2,000 processes less comfortably
  - But problems start at about 3,000 processes

# Defense Mechanisms for Scaling

- **Debug a smaller job**
  - Typically debugging 1 process!
  - Debugging 4 to 32 processes is very common
  - Some routinely debug 100 processes
  - Rarely 1,000 processes or more
- **Developed Subset Attach feature**
  - Debug a subset of processes in a large parallel job
  - On job launch or attach
  - Fan out attach during a session
    - Based on MPI process communication state
    - Data values in a scalar array

# Scalability and Performance Philosophy

- **Must scale in three dimensions**
  - Resource consumption: Do we fit within system limits?
  - Runtime performance: Are we responsive to the user?
  - Presentation of information: Can we present information in an easily digested format?
- **Must actively work on scalability and performance**
  - Machines and programs are getting bigger
  - Feature additions tend to slow things down
  - Continuously test and measure, using a hands-on approach, which requires access to the
    - End-user's machine resources
    - End-user's application
    - End-user's usage scenarios

- **Special effort will be required for BG/L scale machines**
- **Performance and scalability approaches under consideration**
  - Restructured finite-state machine
  - Push more computation into the tvdsvr
  - Restructure messaging: Async TV⇔tvdsvr, "psychic", aggregated, optimistic
  - Multi-threading TV and/or tvdsvr
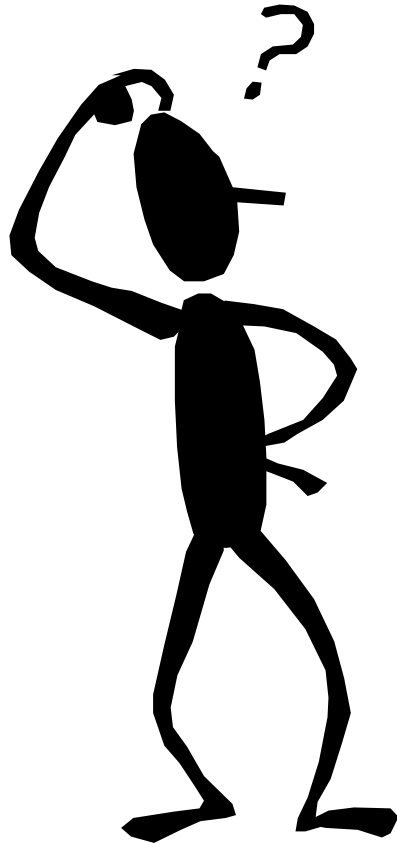- **Haven't thought too hard about how to scale the GUI yet**

# Do We Need a Paradigm Shift?

ETNUS
TotalView

- **Above concerned with scaling TV, but I have a few questions for you …**
- **Do users really want to debug 64,000 processes?**
- **What goes wrong at 64,000 vs. 1,000?**
  - Correctness? Performance? Something else?
- **Do we need additional facilities?**
  - Lightweight "watchdog" debugger
  - High-scale lightweight event tracer
  - Hardware performance counters & tools

# Questions and Answers

- **www.etnus.com**
- **info@etnus.com**