# Debugging and Optimizing Programs Accelerated with Intel® Xeon® Phi™ Coprocessors

Chris Gottbrath, Principal Product Manager
May 7th, 2013

CUG 2013, Napa, CA

# Rogue Wave Today

**The largest independent provider of cross-platform software development tools and embedded components for the next generation of HPC applications.**

## Highlights

- Pioneers in C++/object-oriented development
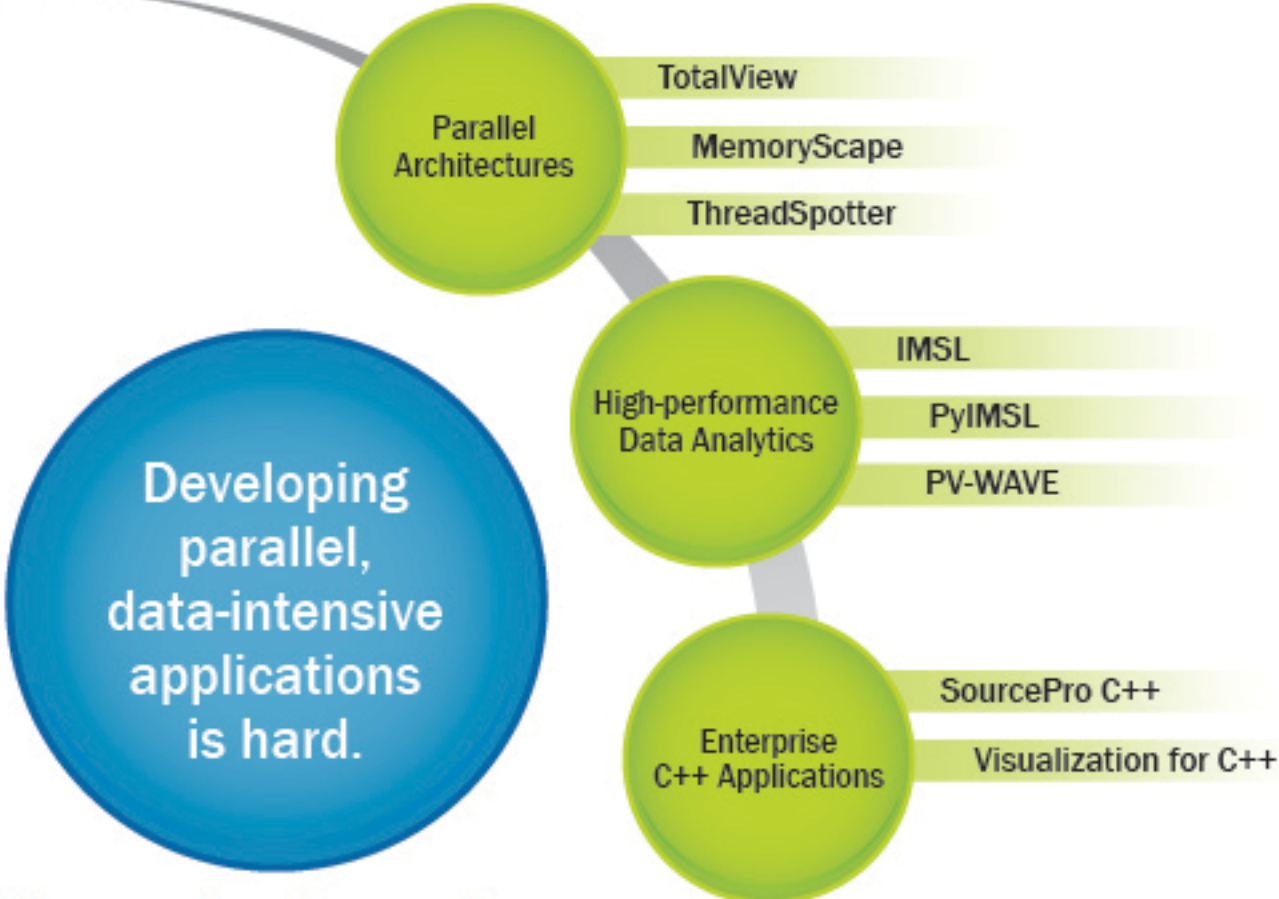- Leading the way in cross-platform, parallel development

## History

- Founded: 1989
- Acquired by Audax Group: 2012
- Acquired:
  - Visual Numerics: 2009
  - TotalView Technologies: 2009
  - Acumem: 2010
  - IBM ILOG Views C++: 2012
- 40 years of experience in HPC

## Customers

- 3,000+ customers in 36 countries
- Multiple sectors:
  - Financial services
  - Telecom
  - Oil and gas
  - Government and aerospace
  - Research and academic

**ROGUE WAVE**
SOFTWARE

# Rogue Wave Solution Portfolio



Parallel Architectures
- TotalView
- MemoryScape
- ThreadSpotter

High-performance Data Analytics
- IMSL
- PyIMSL
- PV-WAVE

Enterprise C++ Applications
- SourcePro C++
- Visualization for C++

Developing parallel, data-intensive applications is hard.
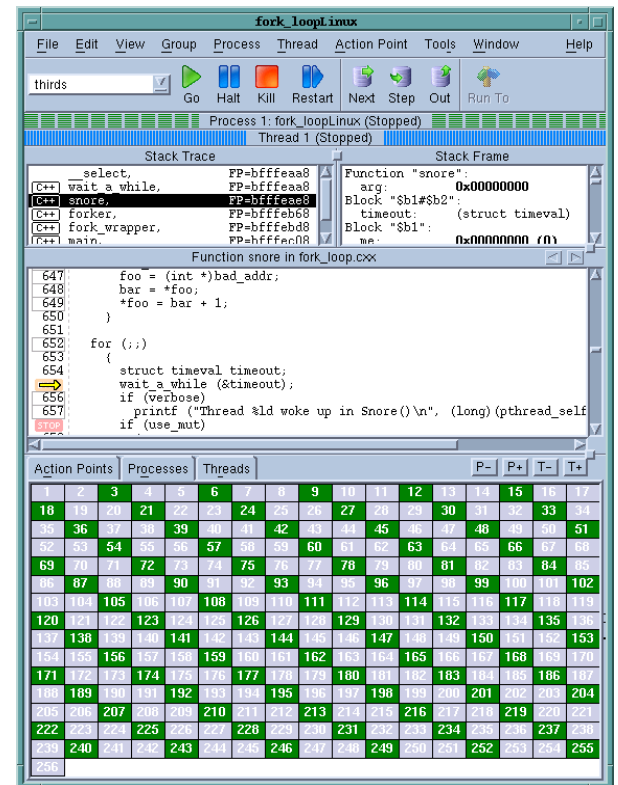
We make it easier.

ROGUE WAVE
SOFTWARE

|

# What is TotalView?

- ## Application Analysis and Debugging Tool: Code Confidently

    - Debug and Analyze C/C++ and Fortran on Linux, Unix or Mac OS X
    - Laptops to supercomputers (Cray, BG, BullX, etc..)
    - Makes developing, maintaining and supporting critical apps easier and less risky

- ## Major Features

    - Easy to learn graphical user interface with data visualization
    - Parallel Debugging
        - MPI, Pthreads, OpenMP, GA, UPC
        - CUDA and OpenACC, Xeon Phi (early access)
    - Includes a Remote Display Client freeing you to work from anywhere
    - Memory Debugging with MemoryScape
    - Deterministic Replay Capability Included on Linux/x86-64
    - Non-interactive Batch Debugging with TVScript and the CLI
    - TTF & C++View to transform user defined objects

**ROGUE WAVE**
**SOFTWARE**

# TotalView for Xeon Phi

- **Support Multiple Intel Xeon Phi configurations**
  - **Native Mode**
    - **With MPI**
  - **Offload Directives**
    - **Similar to GPU**
  - **Multi-device**
  - **Multi-node**
    - **Certain configurations**
  - **CS300-AC, Future XC30**
- **User Interface**
  - **MPI Debugging Features**
    - **Process Control**
    - **View Across**
    - **Shared Breakpoints**
  - **Heterogeneous Debugging**
    - **Debug Both Xeon and Xeon-Phi Processes**



| Copyright © 2013 Rogue Wave Software | All Rights Reserved

ROGUE WAVE
SOFTWARE

# Spectrum of Execution Models

CPU-Centric
Intel® Xeon Processor

Intel® MIC-Centric
Intel® Many Integrated Core (MIC)

| Multi-core Hosted | Offload | Symmetric | Many-Core Hosted |
|---|---|---|---|

**General purpose serial and parallel computing**

**Codes with highly- parallel phases**

**Codes with balanced needs**

**Highly-parallel codes**

*Multi-core*

| Main( ) Foo( ) MPI_*() | Main( ) Foo( ) MPI_*() | Main( ) Foo( ) MPI_*() | |
|---|---|---|---|

PCIe

*Many-core*

| | Foo( ) | Main( ) Foo( ) MPI_*() | Main() Foo( ) MPI_*() |
|---|---|---|---|

**Productive Programming Models Across the Spectrum**

ROGUE WAVE
SOFTWARE

# Remote Debugging of Applications on Xeon Phi



- **Just run as**
  **totalview –r mic0 <program>**

- **Attach to running application**

- **See thread private data**

- **Investigate individual threads**

- **Kill stuck processes on MIC-coprocessor**

**ROGUE WAVE**®
**SOFTWARE**

# Debugging MPI Applications



- **Attach to subset of processes on MIC coprocessor**

- **Set breakpoints**

- **Debug "as usual" MPI**

# Debugging Applications with Offloaded Code



**One debugging session for MIC-accelerated code**

# What's New in TotalView 8.12

- **Xeon Phi Support**

- **Formal support for Cray XC**

- **AVX Instruction Support (phase 1)**

- **Cray ATP Support**

- **Mac OS X Lion and Mountain Lion support**

- **Sessions Manager**

- **STL support for set, multi-set, multi-map**

- **Improvements for specifying addresses in C++ template breakpoints**

- **Updated OS and Compiler Support**

**ROGUE WAVE**
S O F T W A R E ®

# Multi-phase R&D Projects Underway

- **Massive Scalability**
  - Collaboration with LLNL and Tri-lab partners
  - Targeting Cray, Blue Gene and Linux Clusters
- **Shiny new GUI**
  - Sleek, Modern and Fast
  - Configurable
  - Improved Usability
  - Provides aggregation capabilities for big data and scale
  - Leveraging math and stat expertise from IMSL
- **Working with customers through early access programs**
  - Customer input is key to the success of both programs

**ROGUE WAVE**
SOFTWARE

TotalView debugs 786,432 cores.
Climb with Rogue Wave towards exacale.

# Some more details on the 786,432 core test

- **The test was performed on 48 racks of Sequoia**

- **The test code**
  - Implements a Jacobi Linear Equation Solver
  - The test code is a hybrid MPI + OpenMP code
  - 16 threads per process, one process per node

- **The test operations**
  - Start up
  - Setting breakpoints / removing breakpoints
  - Single stepping all threads

- **Tests performed at a variety of scales to understand scalability**

**ROGUE WAVE**
SOFTWARE®

# Second test - Oversubscription

- ## Same framework
  - same code
  - same machine

- ## Oversubscription
  - Scheduled more than one thread per physical core
  - This is a reasonable use case since the BG/Q supports 4 logical threads per core

- ## TotalView Debugged 1,048,576 threads

**ROGUE WAVE**
S O F T W A R E

# What is ThreadSpotter?

- **Runtime Cache Performance Optimization Tool: Tune into the Multi-Core Era**
  - **Realize More of the Performance Offered by Multi/Many-Core Chips**
  - **Quickly Detects and Prioritizes Issues -- and then Provides Usable Advice!**
    - **Brings Cache Performance Into Reach for Every Developer**
    - **Makes Experienced Cache Optimizers Hyper-Efficient**
- **Features**
  - **Supports Linux x86/x86-64 & Windows**
  - **Any compiled code**
  - **Runtime Analysis**
    - **Low overhead**
  - **Cache Modeling**
    - **Prioritizes Issues**
    - **Identifies Problem Lines of Code**
  - **Provides Advice**
    - **Explanations**
    - **Examples**
    - **Detailed statistics (if desired)**



|

ROGUE WAVE
SOFTWARE

# Simple modifications can make a big difference

| Program A | Program B |
|---|---|
| <pre>struct DATA<br>{<br>    int a;<br>    int b;<br>    int c;<br>    int d;<br>};<br>DATA * pMyData;<br><br>for (long i=0; i<10*1024*1024; i++)<br>{<br>    pMyData[i].a = pMyData[i].b;<br>}</pre> | <pre>struct DATA<br>{<br>    int a;<br>    int b;<br>};<br><br><br>DATA * pMyData;<br><br>for (long i=0; i<10*1024*1024; i++)<br>{<br>    pMyData[i].a = pMyData[i].b;<br>}</pre> |

## Partially Used Structures

|

**ROGUE WAVE**
SOFTWARE

# Partially Used Structures

**Defined data structure includes a,b,c,d… but only uses a & b**

**50%**

**Redefined data structure includes a,b,a,b,a,b… c,d are elsewhere.**

**100%**

**ROGUE WAVE**
S O F T W A R E

# Other opportunities for optimization include

- **Alignment Problems**

- **False Sharing**

- **Excessive communication (cache coherence) traffic**

- **Temporal locality issues**

- **Spatial locality issues**

- **Loop fusion**

# Recent improvements to ThreadSpotter

- **Improved parallel support**
  - **Support for sampling all MPI processes in an MPI job**
  - **Cray XT, XE, XK Support**
    - **ALPS, SLURM and Torque**
  - **Continued additions to the processor library**
    - **Including cross-processor analysis**



Figure A.1. MPI Sampling Principles



Figure A.2. Message Passing Toolkit, runtime system and shepherd process

**ROGUE WAVE**
SOFTWARE

# Next release: Improving ThreadSpotter MPI support

- ## Launchmon
  - Provides scalable mechanism for launching the tool in HPC clusters
  - Allows for coordination and synchronization of sampler activity
    - Will reduce "load balancing" bias that might otherwise be introduced by uncoordinated burst sampling with ThreadSpotter
  - Parallel framework can also be used for post-sampling processing
- ## Clustering Analysis
  - Some level of variability in sample results across the run
    - However the bulk of the results will be similar
  - Identify clusters of similar performance data
  - Present a small number (2-5) of reports that represent those clusters
  - Cluster analysis is done in parallel right after the sampling is completed

Rogue Wave Proprietary; plans and forward-looking statements subject to change without notice

|

**ROGUE WAVE**
S O F T W A R E

# ThreadSpotter work towards supporting the Xeon Phi

- **Xeon Phi has an interesting cache architecture**
  - **L1 & L2 caches for each core**
  - **The set of all the L2 sometimes described as "shared"**
  - **L2 caches organized around a ring-shaped bus**
    - **Duplication of data referenced by more then one**
  - **Successful cache utilization is important to achieving performance**
- **Modeling and analysis of this cache architecture**
- **Sampler**
  - **Updated for Xeon Phi vector instructions**
  - **Scaling up the sampler for many-core thread parallelism**
- **Project is still ongoing**

**ROGUE WAVE**
S O F T W A R E

# Thanks!

- **Talk to us here at CUG**

- **Contact me at: chris.gottbrath@roguewave.com**
  - Sign up for the TotalVIew 8.12 beta (Xeon Phi)
  - Learn more about ThreadSpotter
  - Feedback, suggestions, use cases

- **Learn more at: www.roguewave.com**
  - White papers
  - Product Documentation
  - Videos
  - Product evaluation

**ROGUE WAVE**
SOFTWARE