# Advanced Technologies and Techniques for Debugging CUDA GPU HPC Applications

**Rogue Wave**
SOFTWARE

Innovate with Confidence

Nikolay Piskun
Director of Continuing Engineering , TotalView products
GTC  March, 2019

# Agenda

- What is debugging and why TotalView?

- Overview of TotalView

- GPU debugging

- Python debugging

- Advanced C++ and Data debugging

- TotalView resources and documentation

- Questions/Comments

# What is Debugging and
# Why do you need TotalView?

# What is Debugging?

- Debugging is the process of finding and resolving defects or problems within a computer program or a system.
  - Algorithm correctness
  - Data correctness
  - Scaling/Porting correctness

# TotalView debugger enables you to do:

- **Interactive debugging**

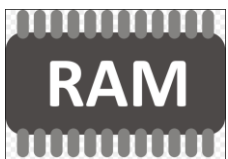  – Live control of an executing program

- **Remote debugging**

  » Debug a program running on another computer

- **Post-mortem debugging (core files and reverse debugging)**

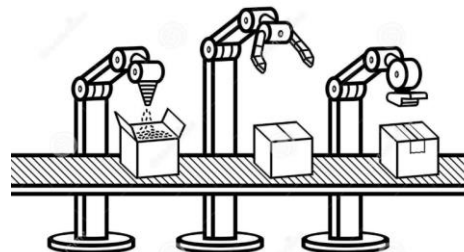  – Debugging a program after it has crashed or exited

- **Memory debugging**

  » Find memory management problems (leaks, corruption …)
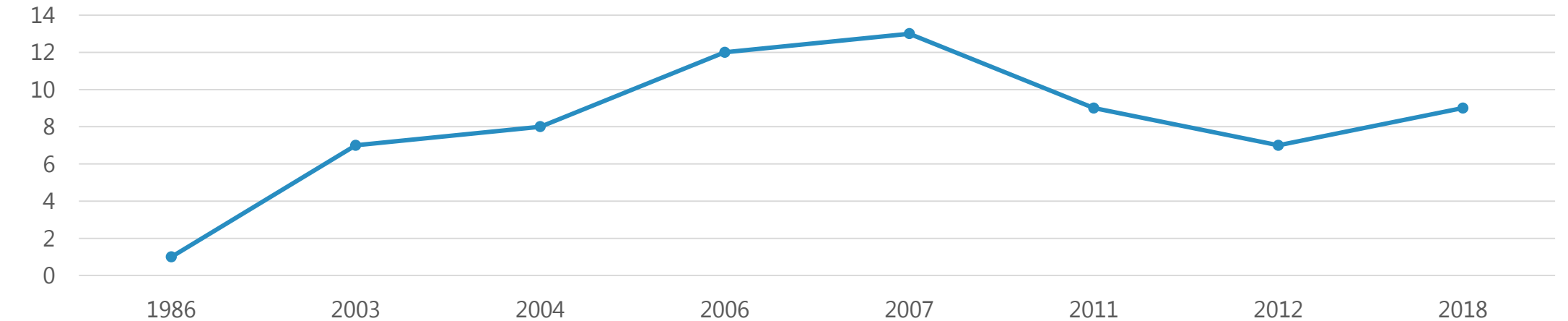
  » Comparing results between executions

- **Batch debugging (tvscript, CI environments)**

  – Unattended debugging
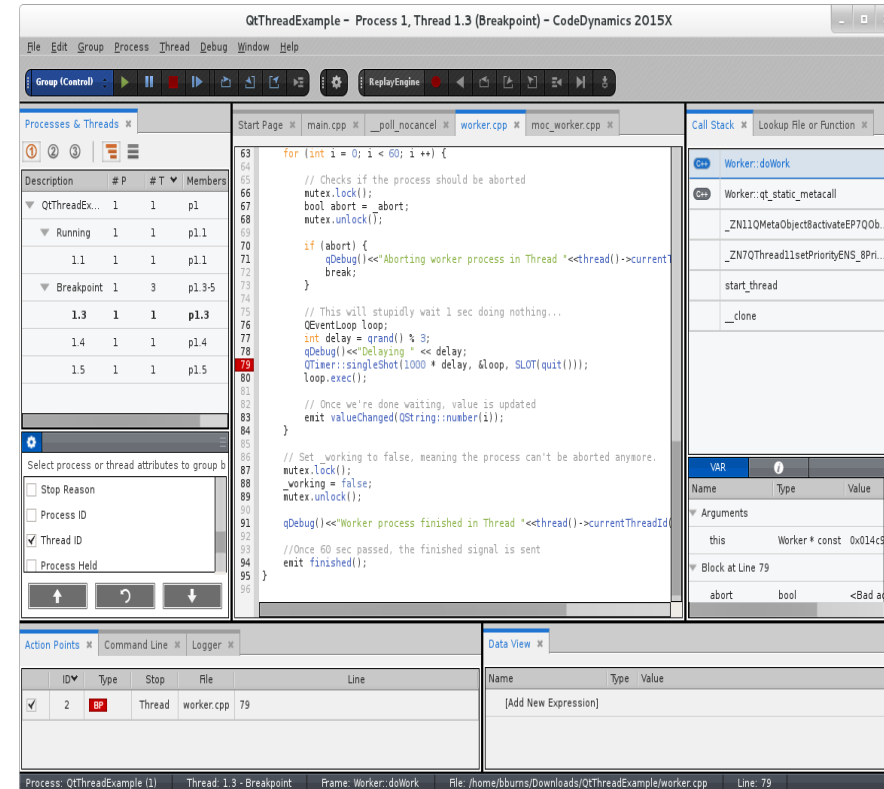
# History of TotalView

## Number of architectures supported

# TotalView for HPC and for All

- Leading debug environment for HPC users
  - Active development for 30+ years
  - Thread specific breakpoints
  - Control individual thread execution
  - View complex data types easily
  - From **MacBook** to **Top500** Supercomputers
- Track memory leaks in running applications
- Supports C/C++ and Fortran on Linux/Unix/Mac
- Integrated Reverse debugging
- Batch non-interactive debugging.

- Allowing the business to have
  - Predictable development schedules
  - Less time spent debugging

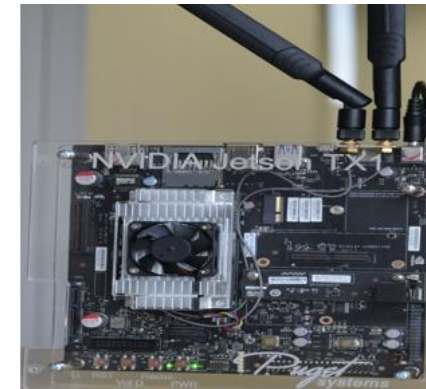**Nvidia software partner**

# High Performance Computing Applications

- HPC Applications cover many different computing areas including:
    - Healthcare and Medicine
    - Modeling and simulation
    - Security
    - Bioinformatics
    - Molecular Dynamics
    - Environment (earthquake/tsunami)/Weather
    - Machine Learning/Artificial Intelligence
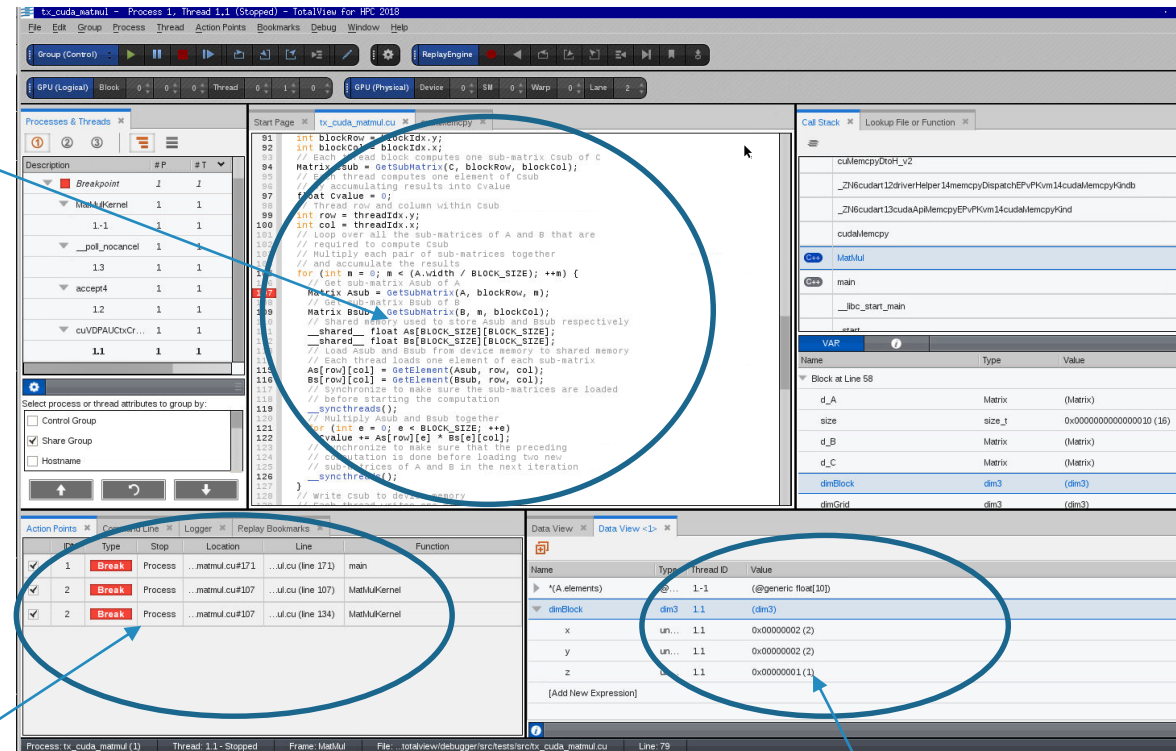
# GPU Debugging

# GPU debugging with TotalView

- NVIDIA CUDA support
  - Multiple platforms : X86-64, PowerLE, ARM64 (in beta)
  - Multiple cards:  from Jetson to Volta (Turing testing)
- Features and capabilities include
  - Support for dynamic parallelism
  - Support for MPI based clusters and multi-card configurations
  - Flexible Display and Navigation on the CUDA device
    - Physical (device, SM, Warp, Lane)
    - Logical (Grid, Block) tuples
  - CUDA device window reveals what is running where
  - Support for CUDA Core debugging
  - Leverages CUDA memcheck
  - Support for OpenACC

# CUDA Debugging Model Improvements

- First in class Unified Source debugging
- Improves and streamlines debugging CUDA applications



- Set breakpoints in CPU **and** GPU kernel code before it is launched on the GPU

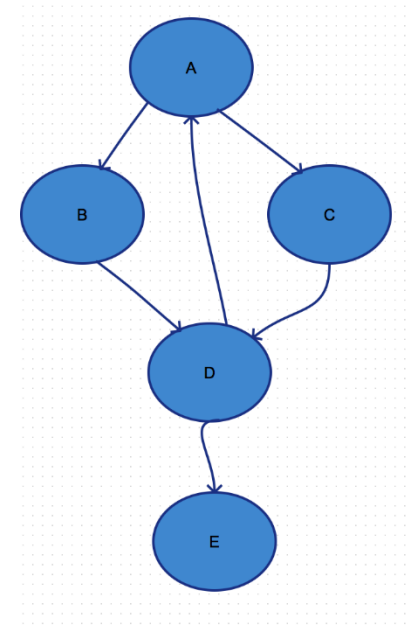- Compare variables in CPU and GPU code together

# CUDA Debugging Demo

# Python/C++ Debugging

# Debugging multiple languages

- Debugging one language is difficult enough
  - Especially with many threads/processes

- The language intersection is tougher
  - Data comparison
  - Glue code

- Issues are:
  - Type mismatches
  - Extraneous stack frames

# Why Python ?

- Use Python to build applications that call out to C++

- Provides access to

  - High-performance routines

  - Leverage existing algorithms and libraries

  - Utilize advanced multi-threaded capabilities

- Calling between languages easily enabled using technologies such as SWIG, ctypes, Cython, CFFI, et al

- Debugging mixed language applications is not easy

  - Good for debugger developers ☺

# Python debugging with TotalView

- What TotalView provides:

  - Easy Python debugging session setup

  - Fully integrated Python and C/C++ call stack

    - "Glue" layers between the languages removed

  - Easily examine and compare variables in Python and C++

  - Utilize reverse debugging and memory debugging


- What TotalView does not provide (yet):

  - Setting breakpoints and stepping within Python code

# Demo

```
#!/usr/bin/python

def callFact():
    import tv_python_example as tp
    a = 3
    b = 10
    c = a+b
    ch = "local string"
     ......
    return tp.fact(a)
if __name__ == '__main__':
    b = 2
    result = callFact()
    print result
```

Terminal

ubuntu:~/demo_2019/PythonExamples> /usr/toolworks/totalview.2019.0.4/bin/totalview -args python2.7-dbg test_python_types.py

# Python without special debugger support

No viewing of Python data and code

**Glue code**

# Showing C code with mixed data

Glue code filtered out Python data and code available for viewing



**Shows Python & C++**

**C++ data**

**Py data**

# Stack Transformation Facility

- Hides stack frames

- Transforms stack frames

- Backbone for:
  – Python support
  – OpenMP support

- **Useful for any glue code you want to hide**
  – **Language differences**
  – **Wrapper code**

```
d1.<> dstacktransform list
Transformation Status: Enabled

Rules
 ID Transform        Operation  Filter
  1 RW_Python        modify     image('python[2-9]\.[0-9]+-dbg'),function('PyEval_EvalFrameEx')
  2 RW_Python        remove     image('python[2-9]\.[0-9]+-dbg')


Transforms
 Name               Implementation
 RW_Python          <built-in>
```

# TensorFlow basics

- Open source

- Numerical computation

- Usage in machine learning

- Written in C++
  - Called from Python

# TensorFlow



Multi-threaded application
   Glue code removed
   Added a rule for wrappers

# Advanced C++ and Data Debugging

# Advanced C++ and Data Debugging



- TotalView supports debugging the latest C++11/14 features including:
  - lambdas, transformations for smart pointers, auto types, R-Value references, range-based loops, strongly-typed enums, initializer lists, user defined literals
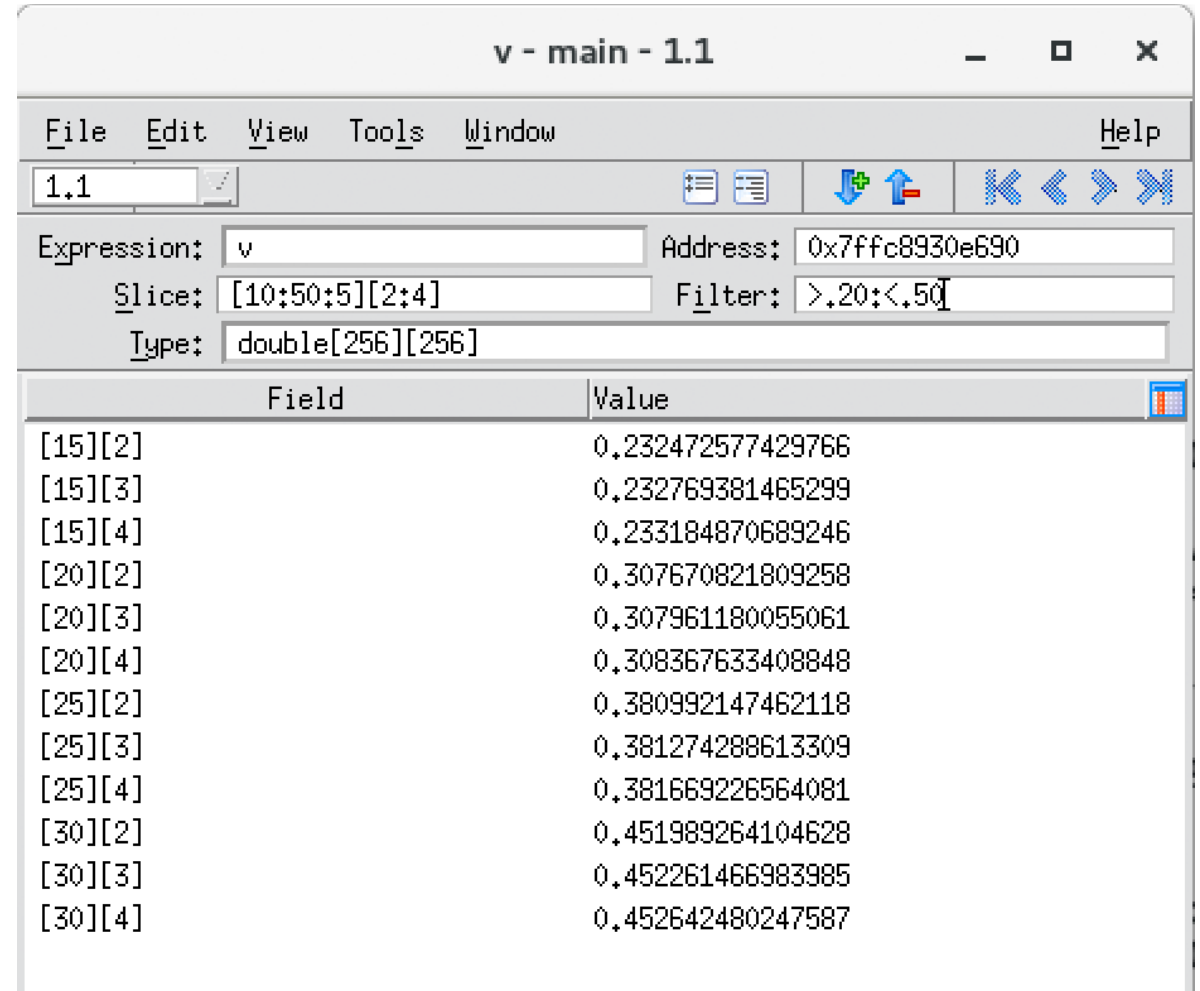
- TotalView transforms many of the C++ and STL containers such as:
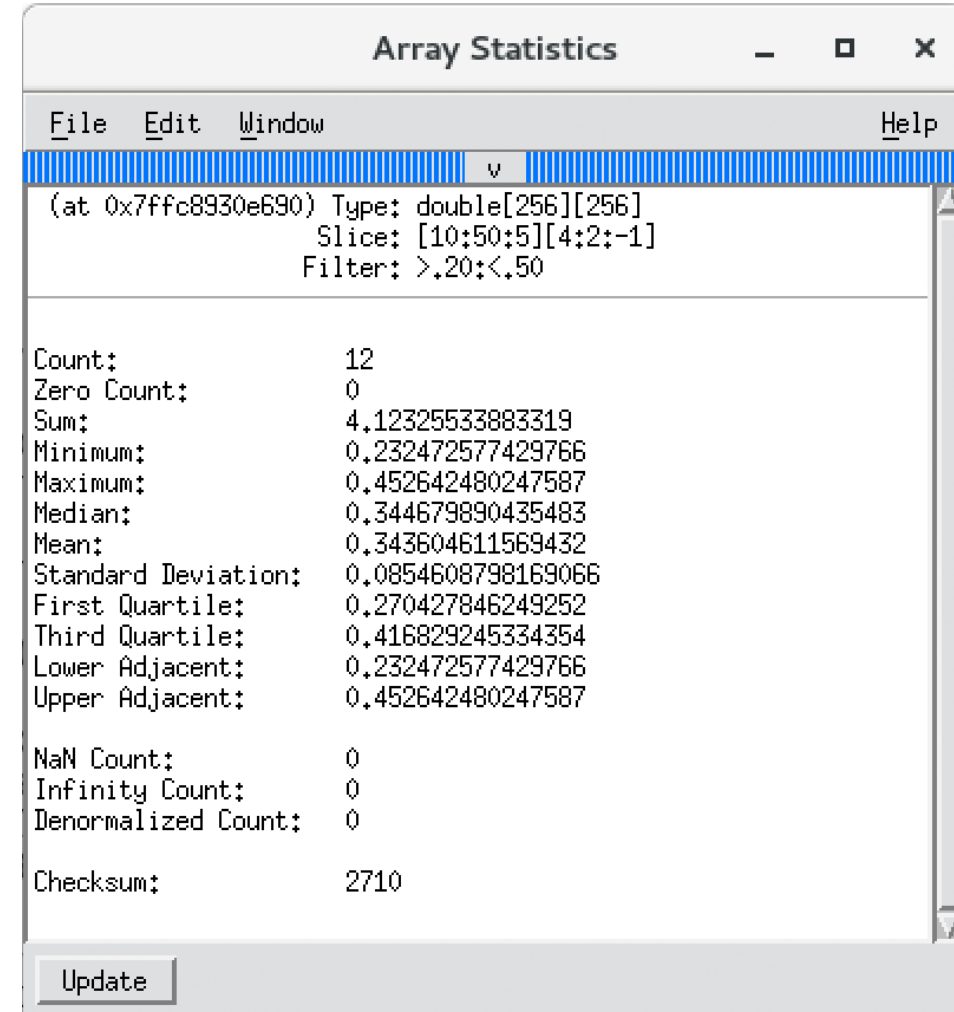  - array, forward_list, tuple, map, set, vector and others.

# Array Slicing, Striding and Filtering (classic UI)

- Slicing – reduce display to a portion of the array
  - [lower_bound:upper_bound]
  - [5:10]
- Striding – Skip over elements
  - [::stride]
  - [::5], [5:10:-1]

- Filtering
  - Comparison:  ==, !=, <, <=, >, >=
  - Range of values:  [**>**] *low-value* **:** [**<**] *high-value*
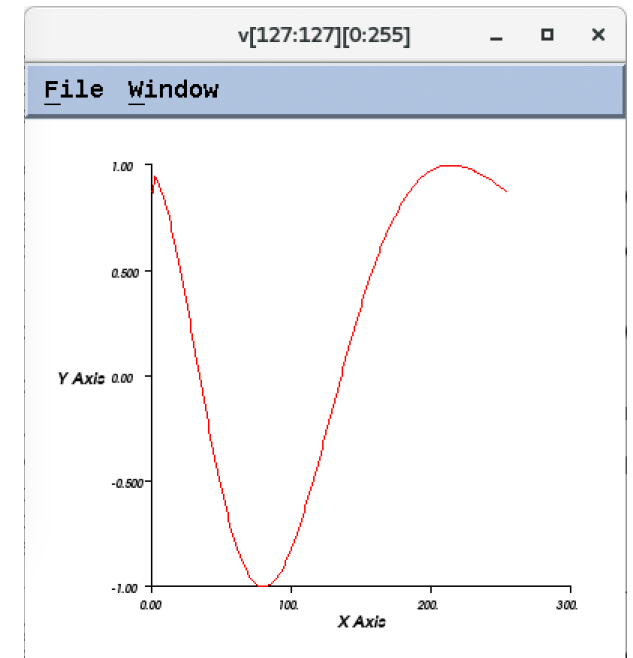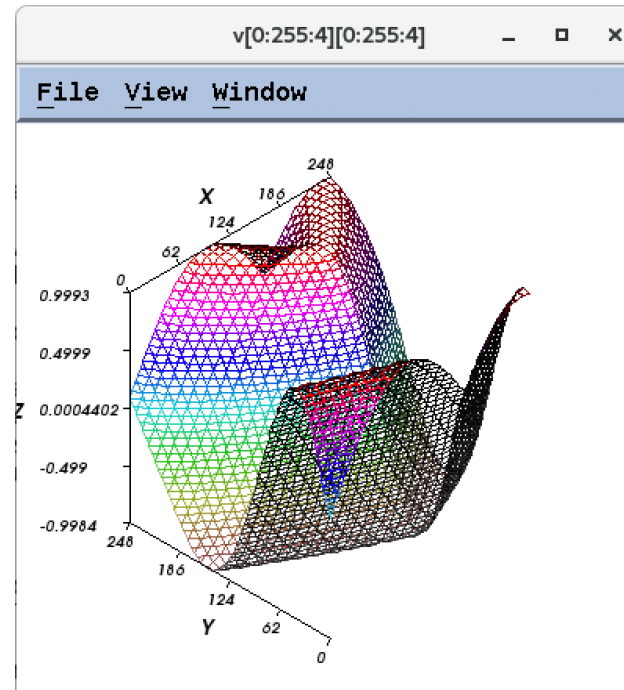  - IEEE values:  $nan, $inf, $denorm

# Array Statistics

- Easily display a set of statistics for the filtered portion of your array

# Visualizing Array Data

- Visualizer creates graphic images of your program's array data.
- Visualize one or two dimensional arrays
- View data manually through the Window > Visualize command on the Data Window
- Visualize data programmatically using the $visualize function

# Summary

- Use of modern debugger <span style="color:red">saves</span> you time.

- TotalView can help you because:
    - It's **cross-platform** (the only debugger you ever need)
    - Allow you to debug accelerators (GPU) and CPU in **one session**
    - Allow you to debug **multiple languages** (C++/Python/Fortran)

# Presenter



**Nikolay Piskun**

**Director of CE, TotalView  products**

Rogue Wave Software/Perforce Software

nikolay.piskun@roguewave.com