

# DEEP Tutorial: Parallel Performance Profiling at the Source Code Level.

David McNamara  
Veridian PSR  
September 2000

# DEEP Tutorial: Topics

- Who is Veridian PSR?
- What is DEEP?
- What can DEEP do for my programs?
- DEEP program analysis demonstration.
  - Gathering the data.
  - Starting DEEP.
  - Finding Bottlenecks.
  - Improving the code.
- DEEP Debugger.
- DEEP future plans.

# Veridian -- Pacific-Sierra Research High Performance Computing Group

- PSR founded 1971, part of Veridian since 1998.
- High Performance Computing group has been in Santa Monica, CA since 1979.
- Products: VAST, DEEP
- Code optimization and porting, seminars on parallel programming, consulting.
- Past and present customers include NEC, Cray, IBM, Fujitsu, Hitachi, Convex, Alliant, Intel, DARPA, Absoft, Apogee, ...

# DDevelopment EEnvironment for PParallel Programs: DEEP

- **Motivation:** Understanding parallel program performance and behavior is **hard**.
- **Goal:** Provide an integrated, cross-platform, parallel development environment *at the user's source code level*.
- Originally developed under DARPA SBIR contract.
- **Status:** First available in 1999.

# DEEP: Target Parallel Systems

- One environment for all popular parallel programming choices.
  - Fortran and C (and mixed). C++ soon.
- **Shared Memory Parallel.**
  - Automatic parallelization
  - OpenMP Programming
- **Distributed Memory Parallel.**
  - Data Parallel Programming (HPF and Data Parallel C)
  - MPI Programming
- Scalar Performance
  - Cache misses, branch performance, ...
- Vector Performance
  - Mflops rate, vector length, ....

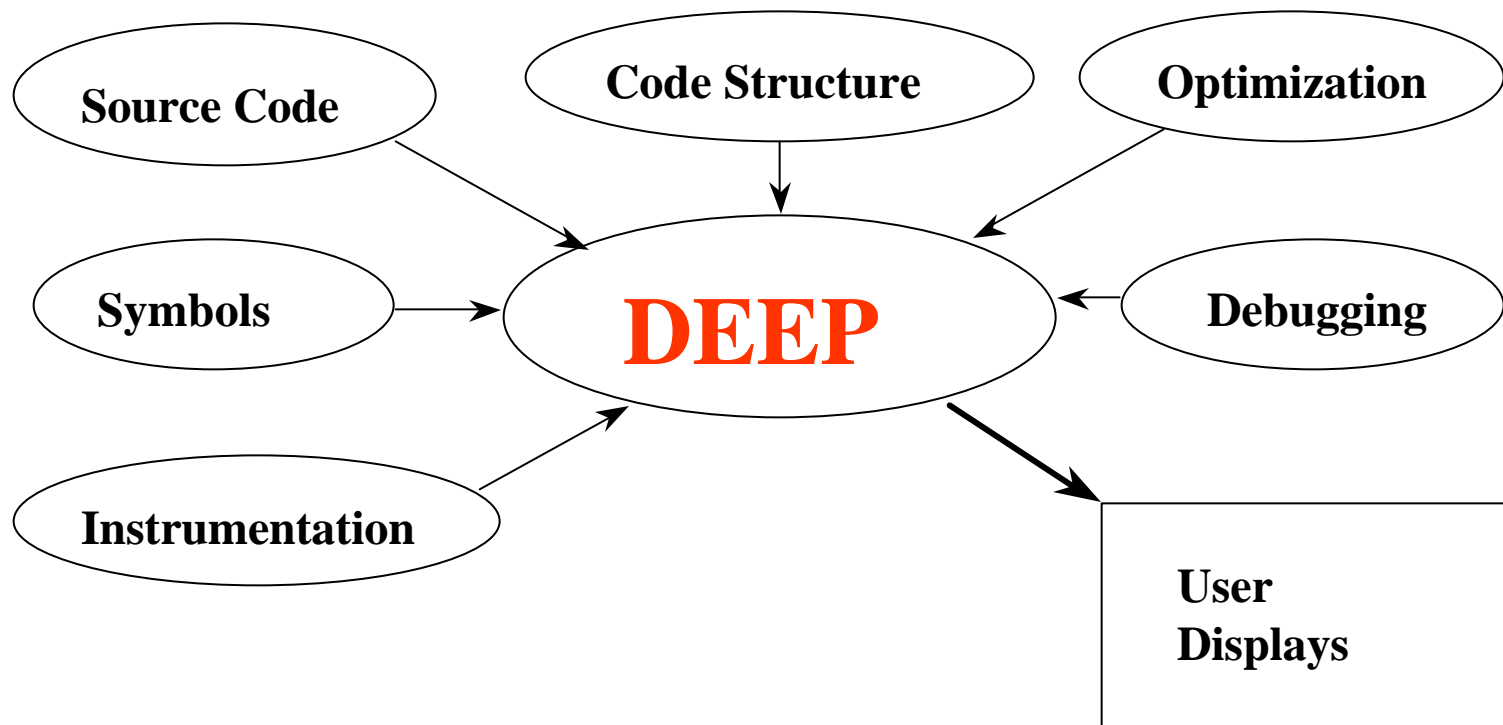
# DEEP: Overview

- Supports both Unix and Windows interfaces.
  - **Unix**/Motif (Alpha, SUN/Solaris, IBM/AIX, SGI/IRIX, ...)
  - **Linux**/Motif (x86, Alpha)
  - **Windows**/NT, 95, 98 (x86)
  - Can run DEEP on non-target system (cross analyze).
- Coordinated set of tools.
  - Move quickly between tools, point and click.
  - Single master graphical user interface.
  - Framework allows organization of lots of information
    - better than many separate windows.

# DEEP: Features

- **Profiling**
  - Where is the time spent? Which are the most important loops? How well are these loops being optimized?
  - Where are most of the messages passed? Is the processor load balanced? How much waiting is being done at synchronization points?
- **Program Structure Browsing**
  - How are the procedures connected to each other?
  - Where are global (parallel) variables referenced and set?
- **Debugging (Under development)**
  - What is the value of this distributed array?
  - What is the current status of the parallel processes?
  - What run-time data types do I have defined?

# DEEP: Data gathering





# DEEP: Profiling at All Levels

**Distributed Memory Parallel System (load balance, communication traffic).**

**Message passing library, data parallel languages. Top level loops.**

**Shared Memory Node (synchronization bottlenecks).**

**OpenMP, automatic parallelization. Outer loops.**

**Cache Memory Levels (cache miss rates).**

**L3, L2, L1. Loop nests.**

**Vector Unit (MOP rates, % utilization)**

**Hand or auto vectorization. Inner loops.**

**Scalar Unit (instruction issue rates)**

**Inner loops.**

# DEEP Program Analysis: Drill Down

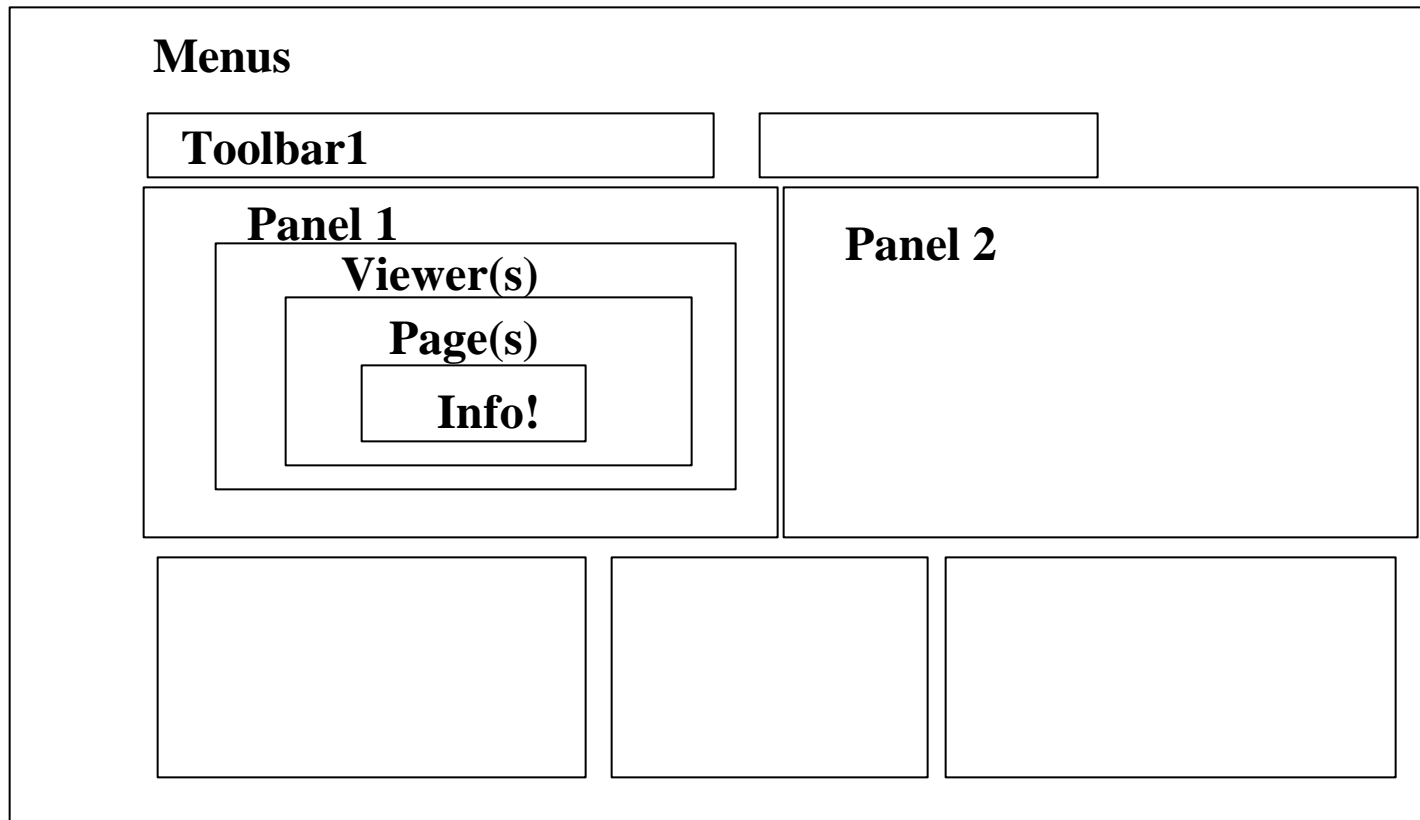
- Start with whole program data.
- Identify “bottleneck” procedures.
  - charts, tables
- Move to details for procedure.
  - code abstraction, loop performance
- Identify run-time problems.
  - “hot” loops, color-coded
- Move to source code.
  - click on important loops
- Remove bottlenecks.
- Recompile, run again and look at performance changes. Iterate if needed.

# DEEP: Drill Down

- Or, drill up?



# DEEP: Framework



# DEEP: MPI Support

- MPI message load balance is recorded and displayed.
- All MPI call sites in the application are profiled.
- Wait time in MPI calls is presented, so inefficient synchronizations can be spotted.
- MPI usage is summarized for each routine in the program.
- MPI operations are highlighted in the call tree and the code abstract, so you can move quickly to that spot in the code.

# DEEP: PerfAPI Support

- PerfAPI (PAPI)
  - PAPI is a cross-platform library that allows access to hardware performance counters.
  - Cache misses for various levels, instruction counts, branch prediction counts, instruction stalls etc.
  - over 50 events (not all supported by all platforms).
- DEEP PAPI support
  - User can select counters to be profiled.
  - Example: show L2 cache misses and branch misprediction.
  - Counts are available at loop and routine levels.
  - Automatically graphed and included in tables.

# DEEP: Gathering Data

- Run profiler on input source files.
  - This creates “static information files”.
    - Compiler optimization information, program structure, symbol table.
  - Profiler also instruments the source code.
    - Adds data gathering directly to the source code.
- Run the instrumented program.
  - This creates a “run file” for each distributed processor.
  - Run files are accumulated information for a particular processor (not a huge trace file).
- Run DEEP to view complete picture.

# DEEP: Instrumentation

- DEEP instruments the source code directly.
- This allows gathering data on loops, loop nests, call sites, and other language structures lower than the function level.
- Measured instrumentation overhead can vary greatly depending on the program (1% to 200%).
  - If loops turn out to have very small iteration counts, then overhead of instrumentation can swamp computation time.
- If instrumentation impact is high, then use sampling mode.
  - Samples first 1000 instances, then a few instances (about one out of 100) from then on.



# Using DEEP for Program Analysis

- Compile the program with the DEEP profiling driver.
  - MPI: `mpiprof myprog.f -or- mpiprof myprog.c`
  - PAPI: `dprof myprog.f -or- dprof myprog.c`
- Run the program.
  - Do this however you normally execute your program.
    - MPI: `mpirun -np 4 myprog < inputfile1.dat`
    - PAPI: `a.out < myinputfile.dat`
- Start DEEP.
  - `deep`
- Provide path of source code and DEEP information.
  - File dialog box: select any source file.
- Select Run (if more than one). Click on row.
- View information. Determine Bottlenecks. Repeat.

# DEEP Files

- All files (compile-time and run-time) are written to subdirectory “/deep”.
  - If you compile and run in different directories, you need to combine the “/deep”s into one “/deep” prior to starting deep.
- The DEEP files are:
  - \*.sif – static information files, one for each function.
  - \*.sym – symbol table information. This is a binary file.
  - run.\*.\*.dif – runtime profile information.
    - Written by the Instrumenter profiling library at the end of execution of the application.
    - One .dif file is written for each process for each run of the program.
    - Name: *run.run\_no.process\_no.dif*, where run\_no is the number of the run (starting with 0) and process\_no is the number of the logical process (starting with zero).

# DEEP: Instrumentation Options

- Use the `-Wv,-z` switch: `mpiprof -Wv,-zfimoz myfile.c`
- `-z [a] [d] [f] [i] [m] [o] [z]`
  - za : shortcut for -zfimo options for MPI, or `-zfiop` for PAPI.
- **-zd** enable all loop profiling (do not use with `-zo`)
- **-zf** enable external function profiling
- **-zi** enable I/O call profiling
- **-zm** enable message passing profiling for MPI versions
- **-zp** enable PAPI profiling for PAPI versions
- **-zo** enable outermost loop profiling (not with `-zd`)
- **-zz** enable sampling

# DEEP/PAPI: Selecting Counters

- Hardware can only count a few items at a time.
- By default, DEEP/PAPI selects two PAPI events
  - Floating Point Operations and Total Cycles.
- Change before run time (after compile time) by setting a series of environment variables:
  - `setenv PROF_PAPI_1 FP_INS`
  - `setenv PROF_PAPI_2 L2_TCM`
- The general form is:
  - `setenv PROF_PAPI_x FLAG`
- where 'x' is a number from 1 to the maximum number of PAPI events that can be counted
- DEEP automatically adjusts its charts and tables to reflect your choice of PAPI counters.

# DEEP/PAPI Counter Options

<b>PAPI Flag</b>	<b>PAPI Item Description</b>	<b>PAPI Flag</b>	<b>PAPI Item Description</b>
L1_DCM	Level 1 data cache misses	MEM_SCY	Cyc Stalled Waiting for Memory Access
L1_ICM	Level 1 instruction cache misses	MEM_RCY	Cyc Stalled Waiting for Memory Read
L2_DCM	Level 2 data cache misses	MEM_WCY	Cyc Stalled Waiting for Memory Write
L2_ICM	Level 2 instruction cache misses	STL_ICY	Cyc with No Instruction Issue
L3_DCM	Level 3 data cache misses	<b>FUL_ICY</b>	<b>Cyc with Maximum Instruction Issue</b>
L3_ICM	Level 3 instruction cache misses	CCY	Cyc with No Instruction Completion
L1_TCM	Level 1 total cache misses	FUL_CCY	Cyc with Maximum Instruction Completion
<b>L2_TCM</b>	<b>Level 2 total cache misses</b>	HW_INT	Hardware interrupts
L3_TCM	Level 3 total cache misses	BR_UCN	Uncond branch instructions executed
CA_SNP	Snoops	BR_CN	Cond branch instructions executed
CA_SHR	access to shared cache line (SMP)	BR_TKN	Cond branch instructions taken
CA_CLN	access to clean cache line (SMP)	BR_NTK	Cond branch instructions not taken
CA_INV	Cache Line Invalidation (SMP)	BR_MSP	Cond branch instructions mispredicted
CA_ITV	Cache Line Intervention (SMP)	BR_PRC	Cond branch instr correctly predicted
L3_LDM	Level 3 load misses	FMA_INS	FMA instructions completed
L3_STM	Level 3 store misses	TOT_IIS	Total instructions issued
BRU_IDL	Cycles branch units are idle	TOT_INS	Total instructions executed
FXU_IDL	Cycles integer units are idle	INT_INS	Integer instructions executed
FPU_IDL	Cycles floating point units are idle	<b>FP_INS</b>	<b>Floating point instructions executed</b>
LSU_IDL	Cycles load/store units are idle	LD_INS	Load instructions executed
<b>TLB_DM</b>	<b>Data translation lookaside buffer misses</b>	SR_INS	Store instructions executed
TLB_IM	Inst xlation lookaside buffer misses	BR_INS	Total branch instructions executed
TLB_TL	Total xlation lookaside buffer misses	<b>VEC_INS</b>	<b>Vector/SIMD instructions executed</b>
L1_LDM	Level 1 load misses	FLOPS	Floating Point instructions per second
L1_STM	Level 1 store misses	RES_STL	Any resource stalls
L2_LDM	Level 2 load misses	FP_STAL	FP units are stalled
L2_STM	Level 2 store misses	TOT_CYC	Total cycles
BTAC_M	BTAC miss	IPS	Instructions executed per second
PRF_DM	Prefetch data instruction caused a miss	<b>LST_INS</b>	<b>Total load/store inst. executed</b>
TLB_SD	Xlation lookaside buf shootdowns (SMP)	SYC_INS	Sync. inst. executed
CSR_FAL	Failed store conditional instructions		
CSR_SUC	Succ store conditional instructions		
CSR_TOT	Total store conditional instructions		

# View Menu

## Program

### Charts

- Total Messages
- CPU Balance
- Message Balance
- MPI Calls
- MPI Wallclock Time
- Wallclock Time
- Custom Chart

### Views

- Calling Tree
- Loop Time
- Advisor

### Tables

- Compile-time
- MPI Call Sites
- MPI Summary
- Loops
- Inclusive Run-time
- Run-time
- Global Symbols

## Procedure

- Source Code Browser

- Code Abstract

- Loop Table

- Call Table

- MPI Call Sites

- Symbol Table

# DEEP: Screen Shot

DEEP

File Edit Mode View Options Help

a=b c=d a.int b.int a.b b.c

x\_solve.f copy\_faces.f

```

c-----
end do
call mpi_irecv(in_buffer(sr(0)), b_size
> dp_type, successor(1), WEST,
> comm_rhs, requests(0), error)
call mpi_irecv(in_buffer(sr(1)), b_size
> dp_type, predecessor(1), EAST,
> comm_rhs, requests(1), error)
call mpi_irecv(in_buffer(sr(2)), b_size
> dp_type, successor(2), SOUTH,
> comm_rhs, requests(2), error)
call mpi_irecv(in_buffer(sr(3)), b_size

```

Source Code

Pie Chart: Total messages CPU Balance Message Balance

- 16% matmul\_sub
- 16% matvec\_sub
- 10% z\_solve\_cell
- 9% x\_solve\_cell
- 9% binvcrhs
- 9% y\_solve\_cell
- 5% copy\_faces
- 3% lhsx
- 3% lhsy
- 20% other

Charts Views Tables

matmul\_sub copy\_faces

Operation	File/Line
Code: (stmts=2)	copy_faces.f/151
EndLp	copy_faces.f/153
EndLp	copy_faces.f/154
EndLp	copy_faces.f/155
Endif	copy_faces.f/156
EndLp	copy_faces.f/161
EndLp	copy_faces.f/166
MPI: MPI_IRecv	copy_faces.f/168
MPI: MPI_IRecv	copy_faces.f/171

Code Abstraction Loop Performance Call Performance

Processes Performance Advisor

DEEP Performance Advisor

matmul\_sub

The ratio of CPU time to wall clock time is low. Check for overloaded system, or excessive synchronizations or I/O. (0.415)

Loop/I x\_solve.f/557

Loop uses a large percentage of time. Make sure the loop is optimized. (34.5)

Loop has a very small average iteration count. Consider unrolling the loop entirely. (5)

Advisor Symbol Viewer

G:\DEEP\MPI\_NasPBT\copy\_faces.f 168

PACIFIC-SIERRA RESEARCH

VERTICIAN

# DEEP Toolbars

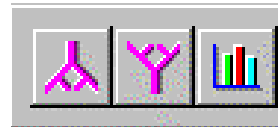
- Procedure Toolbar.

- Source code.
- Code abstract.
- Symbol table.
- Loop performance.
- Call performance.



- Program Toolbar.

- Calling tree.
- Called tree.
- Charts.



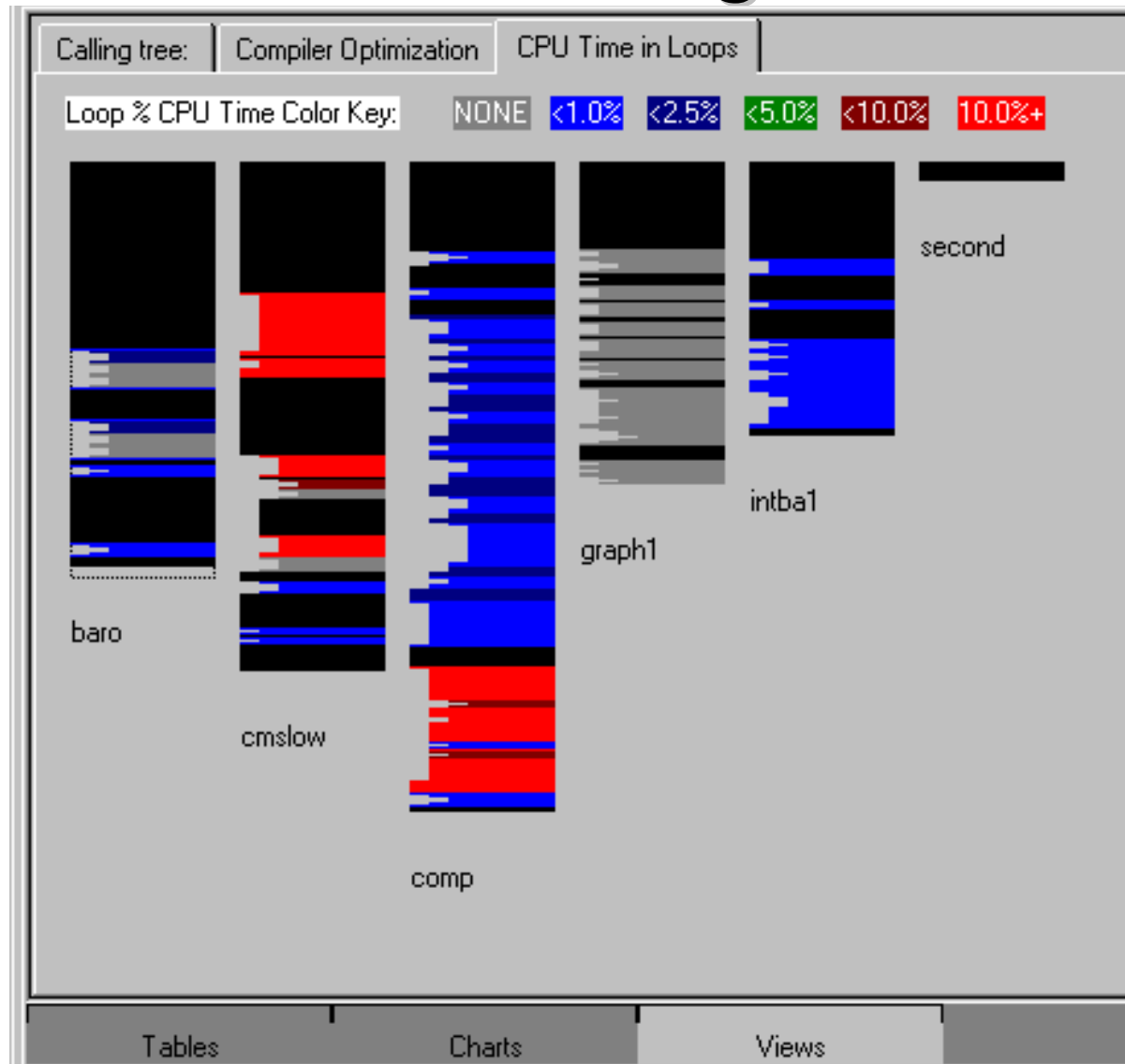
- List Manipulation Toolbar.

- Sort.
- Prune.
- Highlight.





# DEEP: Whole Program View



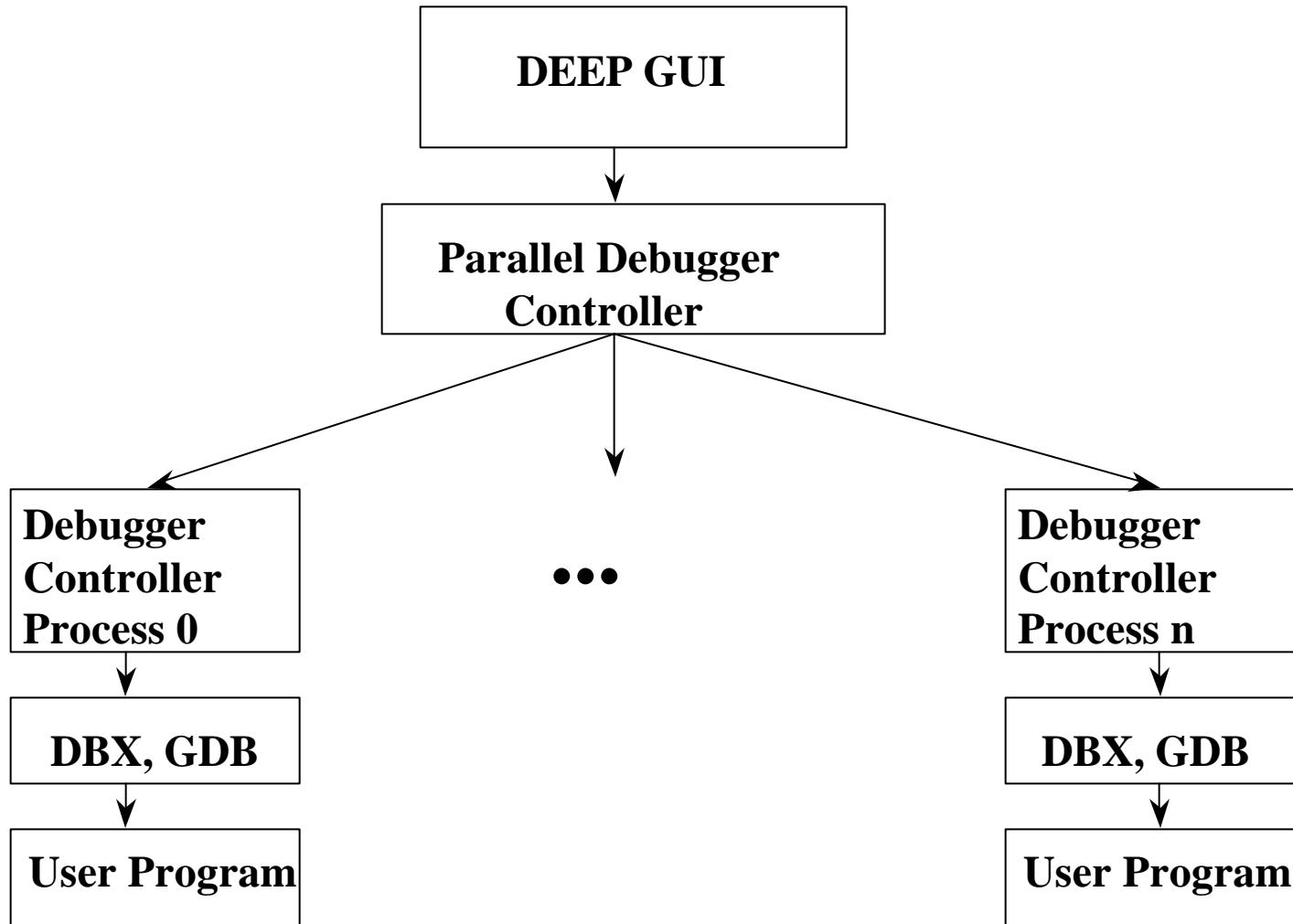
# DEEP: The big picture



# DEEP: Debugger

- Parallel Source Code Debugger currently under development to complement the existing program analysis features.
- Control existing scalar debuggers running on each parallel process.
  - rely on dbx, gdb, etc. for scalar debugger.
- Sophisticated GUI.
  - multiple viewers available.
- Examine values of parallel variables at breakpoints or through “watches”.
- Keep track of current process status for all processes.

# DEEP: Debugger structure



# DEEP Debugger Screen

The screenshot displays the DEEP debugger interface with the following components:

- Source Code:** Shows the main function of a program. The current execution point is at line 881, which is a `switch( CLASS )` statement. The code includes MPI initialization and array verification logic.
- Program Output:** Displays a log of messages from the debugger, including commands like `print comm_size` and `continue`, and results for different IPIDs (0-5). It also shows breakpoint information for `main()` at line 881.
- Watches:** A table showing the current values of `my_rank` and `comm_size`.
- Processes:** A table listing the state of various MPI processes.

**Source Code:**

```

main( argc, argv )
  int argc;
  char **argv;
{
  int          i, iteration, itemp;
  double       timecounter, maxtime;

  /* Initialize MPI */
  MPI_Init( &argc, &argv );
  MPI_Comm_rank( MPI_COMM_WORLD, &my_rank );
  MPI_Comm_size( MPI_COMM_WORLD, &comm_size );

  /* Initialize the verification arrays if a valid class */
  for( i=0; i<TEST_ARRAY_SIZE; i++ )
    switch( CLASS )
    {
      case 'S':
        test_index_array[i] = S_test_index_array[i];
        test_rank_array[i] = S_test_rank_array[i];
        break;
      case 'A':
        test_index_array[i] = A_test_index_array[i];
        test_rank_array[i] = A_test_rank_array[i];
        break;
      case 'N':
    }
}
    
```

**Program Output:**

```

From Debugger
Result = 3
i00?
13 more messages expected
Command : [ipid : 3] print comm_size
Results : [ipid : 3]
Result = 6
i00?
12 more messages expected
Command : [ipid : 2] continue
Results : [ipid : 2] Continue program: is.exe
Breakpoint #2: main() at line 881 in D:\CPP\Deep\Tests\IS\is.c.
881          switch( CLASS )
i00?
11 more messages expected
Command : [ipid : 4] print my_rank
Results : [ipid : 4]
Result = 4
i00?
10 more messages expected
Command : [ipid : 1] continue
Results : [ipid : 1] Continue program: is.exe
Breakpoint #2: main() at line 881 in D:\CPP\Deep\Tests\IS\is.c.
881          switch( CLASS )
i00?
9 more messages expected
Command : [ipid : 5] print comm_size
Results : [ipid : 5]
    
```

**Watches:**

Expression	Value
my_rank	<0> 0. <1> 1. <2> 2. <3> 3. <4> 4. <5> 5
comm_size	6

**Processes:**

Logical ID	Machine ID	Status	Groups
0	Appgeo/25427	Paused at line 881	MPI_COMM_WORLD
1	Appgeo/25428	Paused at line 881	MPI_COMM_WORLD
2	Appgeo/25429	Paused at line 881	MPI_COMM_WORLD
3	Appgeo/25430	Paused at line 881	MPI_COMM_WORLD
4	Appgeo/25431	Paused at line 881	MPI_COMM_WORLD
5	Appgeo/25432	Paused at line 881	MPI_COMM_WORLD

**Status Bar:** From Debugger 0 Paused

# DEEP: Future Directions

- “Project” feature: full IDE capability
  - Specify source files, DEEP launches compiles, links, etc.
- Improved support for large numbers (<100) of processors.
  - Search tools to look for outlying cases...
  - Matrix rather than column displays...
- Allow inserting of Library calls from menu.
  - Provide correct number and name of arguments.
  - Allow insertion of arbitrary set of routines (from file).
  - Check syntax and semantics of calls. MPI, BLAS, VSIPL, etc.
- “Real-time” profiling, through debugger interface.

# DEEP: More Directions

- Distributed Objects
  - Corba/DCOM...
- Enhanced structure debugging for parallel structures.
- Event Trace Tool
  - Map events to original source code line, zoom in/out.
- ccNUMA performance tuning
  - relate local/global memory perf. back to source code.
- Support for counter multiplexing with PAPI.

# Advantages of DEEP

- Can use DEEP for all major parallel programming paradigms.
- Quickly zoom in on the bottlenecks and problems. Saves time, especially on larger applications.
- More fun than putting in your own print statements. Look at interesting displays.
- All tools in one place, don't have to learn several interfaces. Tools all work together.
- Can be customized for special requirements.



# Summary

- For further information:  
David McNamara  
Veridian -- Pacific-Sierra Research  
2901 28th Street  
Santa Monica CA 90405  
(310)-314-2338  
brian@psrv.com