# Debugging and Optimization Tools

**Richard Gerber**
NERSC User Services
**David Skinner**
NERSC Outreach, Software & Programming Group
UCB CS267
February 16, 2012

U.S. DEPARTMENT OF ENERGY | Office of Science

NeRSC
National Energy Research
Scientific Computing Center

BERKELEY LAB
Lawrence Berkeley
National Laboratory

# Outline

- **Introduction**
- **Debugging**
- **Performance / Optimization**

**Videos, presentations, and references:**

**http://www.nersc.gov/users/training/courses/CS267/**

**Also see the DOE Advanced Computational Tools:**
http://acts.nersc.gov

- **Today's Talks**

  – Strategies for parallel performance (D. Skinner)

  – Debugging and optimization tools (R. Gerber)

- **Take Aways**

  – Common problems to look out for

  – How tools work in general

  – A few specific tools you can try

  – Where to get more information

# Debugging

- **A bug is when your code**
  - crashes
  - hangs (doesn't finish)
  - gets inconsistent answers
  - produces wrong answers
  - behaves in any way you didn't want it to

- The term "bug" was popularized by Grace Hopper (motivated by the removal of an actual moth from computer in 1947)

- **"Serial"**
  - Invalid memory references
  - Array reference out of bounds
  - Divide by zero
  - Use of uninitialized variables
- **Parallel**
  - Unmatched sends/receives
  - Blocking receive before corresponding send
  - Out of order collectives
  - Race conditions

- **Find It**
  - You want to locate the part of your code that isn't doing what it's designed to do
- **Fix It**
  - Figure out how to solve it and implement a solution
- **Run It**
  - Check for proper behavior

## printf, write

- Versatile, sometimes useful
- Doesn't scale well
- Not interactive
- Fishing expedition

## Compiler / Runtime

- Bounds checking, exception handling
- Dereferencing of NULL pointers
- Function and subroutine interface checking

## Serial gdb

- GNU debugger, serial, command-line interface
- See "man gdb"

## Parallel debuggers Using X-Windows

- DDT
- Totalview

# Compiler runtime bounds checking

**Out of bounds reference in source code for program "flip"**

…

```
allocate(put_seed(random_size))
```

…

```
bad_index = random_size+1
put_seed(bad_index) = 67
```

```
ftn  -c -g -Ktrap=fp -Mbounds flip.f90
ftn  -c -g -Ktrap=fp -Mbounds printit.f90
ftn  -o flip flip.o printit.o -g


% qsub —I —qdebug —lmppwidth=48
% cd $PBS_O_WORKDIR
%
% aprun —n 48 ./flip
0: Subscript out of range for array
   put_seed (flip.f90: 50)
    subscript=35, lower bound=1, upper
   bound=34, dimension=1
0: Subscript out of range for array
   put_seed (flip.f90: 50)
    subscript=35, lower bound=1, upper
   bound=34, dimension=1
```

Lawrence Berkeley National Laboratory

# Compiler Documentation

- **For a list of compiler options, see the man pages for the individual compilers**
  - man pgcc
  - man pgCC
  - man pgf90
  - man gcc
  - man gfortran
  - Etc.

# Parallel Programming Bug

**This code hangs because both Task 0 and Task N-1 are blocking on MPI_Recv**

```
if(task_no==0) {

    ret = MPI_Recv(&herBuffer, 50, MPI_DOUBLE,
totTasks-1, 0, MPI_COMM_WORLD, &status);
    ret = MPI_Send(&myBuffer, 50, MPI_DOUBLE,
totTasks-1, 0, MPI_COMM_WORLD);

} else if (task_no==(totTasks-1)) {

    ret = MPI_Recv(&herBuffer, 50, MPI_DOUBLE, 0, 0,
MPI_COMM_WORLD, &status);
    ret = MPI_Send(&myBuffer, 50, MPI_DOUBLE, 0, 0,
MPI_COMM_WORLD);

}
```

## Compile for debugging

```
hopper% make
cc -c -g hello.c
cc -o hello -g  hello.o
```

## Set up the parallel run environment

```
hopper% qsub -I -V -lmppwidth=24
hopper% cd $PBS_O_WORKDIR
```

## Start the DDT debugger

```
hopper% ddt ./hello
```

# DDT Screen Shot

**Press Go and then Pause when code appears hung.**

**Task 0 is at line 44**

**At hang, tasks are in 3 different places.**

# DDT Screen Shot



Task 3 is at line 47

At hang, tasks are in 3 different places.

# DDT video

- **http://vimeo.com/19978486**

- **Or http://vimeo.com/user5729706**

- **This is out of date; I need to change the NX server from "Euclid" to "nx.nersc.gov and "hopp2" to "hopper"**

# Other Debugging Tips

- **Try different compilers**

  – Diagnostic messages and language spec compliances differ

- **Look for memory corruption**

  – Bad memory reference in one place (array out of bounds) can make code crash elsewhere

  – It might appear that you're crashing on a perfectly valid line of code

- **Check the arguments to your MPI calls**

- **Call the NERSC Consultants (800-66-NERSC or 510 486-8600)**

# Performance / Optimization

- How can we tell if a program is performing well?

- Or isn't?

- If performance is not "good," how can we pinpoint why?

- How can we identify the causes?

- What can we do about it?

- **Primary metric: application time**
  - but gives little indication of efficiency

- **Derived measures:**
  - rates (Ex.: messages per unit time, Flops per second, clocks per instruction), cache utilization

- **Indirect measures:**
  - speedup, parallel efficiency, scalability

# Optimization Strategies

- **Serial**
  - Leverage ILP on the processor
  - Feed the pipelines
  - Exploit data locality
  - Reuse data in caches
- **Parallel**
  - Minimize latency effects (aggregate messages)
  - Maximize work vs. communication
- **Both**
  - Minimize data movement (recalculate vs. send)
  - Memory locality on NUMA processors - first touch

- **Sampling**
  - Regularly interrupt the program and record where it is
  - Build up a statistical profile of time spent in various routines
  - Concentrate first on longest running sections or routines
- **Tracing**
  - Insert hooks into program to record and time program events (logging)
  - Reasonable for sequential programs
  - Unwieldy for large parallel programs (too much data!)

- **Hardware Event Counters**

  – Special registers count events on processor

  – E.g. number of floating point instructions

  – Many possible events

  – Only a few can be recorded at a time (~4 counters)

  – Can give you an idea of how efficiently you are using the processor hardware

- **(Sometimes) Modify your code with macros, API calls, timers**
- **Compile your code**
- **Transform your binary for profiling / tracing with a tool**
- **Run the transformed binary**
  - A performance data file is produced
- **Interpret the results with a tool**

# Performance Tools @ NERSC

- **Vendor Tools:**
  - CrayPat on Crays
- **Community Tools :**
  - TAU (U. Oregon via ACTS)
  - PAPI (Performance API)
  - gprof
- **IPM: Integrated Performance Monitoring**
  - A low overhead, low effort NERSC tool

- **Suite of tools that provides a wide range of performance-related information**

- **Can be used for both sampling and tracing**
  - with or without hardware or network performance counters
  - Built on PAPI

- **Supports Fortran, C, C++, UPC, MPI, Coarray Fortran, OpenMP, Pthreads, SHMEM**

- **Man pages**
  - intro_craypat(1), intro_app2(1), intro_papi(1)

# Using CrayPat

1. **Access the tools**
   - `module load perftools`
2. **Build your application; keep .o files**
   - `make clean`
   - `make`
3. **Instrument application**
   - `pat_build ... a.out`
   - Result is a new file, `a.out+pat`
4. **Run instrumented application to get top time consuming routines**
   - `aprun ... a.out+pat`
   - Result is a new file XXXXX.xf (or a directory containing .xf files)
5. **Run pat_report on that new file; view results**
   - `pat_report  XXXXX.xf  > my_profile`
   - `view my_profile`
   - Also produces a new file: XXXXX.ap2

- **Optional visualization tool for Cray's perftools data**

- **Use it in a X Windows environment**

- **Uses a data file as input (`XXX.ap2`) that is prepared by `pat_report`**

  `app2 [--limit_per_pe tags] XXX.ap2`

# Apprentice Basic View

- PAPI (Performance API) provides a standard interface for use of the performance counters in major microprocessors

- Predefined actual and derived counters supported on the system
  - To see the list, run 'papi_avail' on compute node via aprun:

    ```
    qsub –I –lmppwidth=24
    module load perftools
    aprun –n 1 papi_avail
    ```

- AMD native events also provided; use 'papi_native_avail':

    ```
    aprun –n 1 papi_native_avail
    ```

- **Tuning and Analysis Utilities**
- **Fortran, C, C++, Java performance tool**
- **Procedure**
  - Insert macros
  - Run the program
  - View results with pprof
- **More info than gprof**
  - E.g. per process, per thread info; supports pthreads
- **http://acts.nersc.gov/tau/index.html**

# TAU Assignment

- **You will have a homework assignment using TAU**
  - `%module load tau`
  - Define paths in Makefile
  - Modify header file to define TAU macros
  - Add macro calls to the code
  - Compile and submit to batch queue
  - Use pprof to produce readable output
- **Good reference**
  - http://acts.nersc.gov/events/Workshop2011/Talks/TAU.pdf

- **NERSC has about 5,000 users**
  - All levels of sophistication and experience
  - We're committed to supporting both the cutting edge & production HPC computing for the masses
- **Users often ask for advice on which tools to use and we give them suggestions**
- **Our experience is that very few use programming/debugging/development tools**
- **A few users use a few tools a lot, but many try a tool only once**

# Why?

- **Extremely effective?**

- **More likely: Too confusing, difficult, didn't work, don't know how to use, don't know which to use, tied to a platform, compiler, or language**

- **It's not that we don't have tools that address specific issues**

  - TAU, PAPI, HPC Toolkit

  - Craypat, IBM HPC tools, OpenSpeedShop, Intel

  - Valgrind  (memory debugging)

  - GPU/CUDA tools & compilers

  - Vampirtrace

- **But do most users have the resources to learn how to use these tools, esp. when they don't know if there will be any benefit from any given one?**

# IPM

- **Integrated Performance Monitoring**
- **MPI profiling, hardware counter metrics, IO profiling (?)**
- **IPM requires no code modification & no instrumented binary**
  - Only a "module load ipm" before running your program on systems that support dynamic libraries
  - Else link with the IPM library
- **IPM uses hooks already in the MPI library to intercept your MPI calls and wrap them with timers and counters**

- **How it works (user perspective)**
  - `% module load IPM`*
  - Run program as normal
  - Look at results on the web
- **It's that easy!**
  - And extremely low overhead, so IPM is examining your production code

  * (As long as your system supports dynamic load libs)

# What IPM measures

- **IPM "only" gives a high-level, entire-program-centric view**
- **Still, very valuable guidance**
  - Shows whole-run info per MPI task, OpenMP thread, (CUDA under development)
  - Many pieces of data in one place
- **Reveals what many users don't know about their code**
  - High-water memory usage (per task)
  - Load balance
  - Call imbalance
  - MPI time
  - I/O time

# IPM

```
# host    : s05601/006035314C00_AIX     mpi_tasks : 32 on 2 nodes
# start   : 11/30/04/14:35:34           wallclock : 29.975184 sec
# stop    : 11/30/04/14:36:00           %comm     : 27.72
# gbytes  : 6.65863e-01 total           gflop/sec : 2.33478e+00 total
#                       [total]         <avg>         min           max
# wallclock            953.272         29.7897       29.6092       29.9752
# user                 837.25          26.1641       25.71         26.92
# system               60.6            1.89375       1.52          2.59
# mpi                  264.267         8.25834       7.73025       8.70985
# %comm                                27.7234       25.8873       29.3705
# gflop/sec            2.33478         0.0729619     0.072204      0.0745817
# gbytes               0.665863        0.0208082     0.0195503     0.0237541
# PM_FPU0_CMPL         2.28827e+10     7.15084e+08   7.07373e+08   7.30171e+08
# PM_FPU1_CMPL         1.70657e+10     5.33304e+08   5.28487e+08   5.42882e+08
# PM_FPU_FMA           3.00371e+10     9.3866e+08    9.27762e+08   9.62547e+08
# PM_INST_CMPL         2.78819e+11     8.71309e+09   8.20981e+09   9.21761e+09
# PM_LD_CMPL           1.25478e+11     3.92118e+09   3.74541e+09   4.11658e+09
# PM_ST_CMPL           7.45961e+10     2.33113e+09   2.21164e+09   2.46327e+09
# PM_TLB_MISS          2.45894e+08     7.68418e+06   6.98733e+06   2.05724e+07
# PM_CYC               3.0575e+11      9.55467e+09   9.36585e+09   9.62227e+09
#                       [time]         [calls]       <%mpi>        <%wall>
# MPI_Send             188.386         639616        71.29         19.76
# MPI_Wait             69.5032         639616        26.30         7.29
# MPI_Irecv            6.34936         639616        2.40          0.67
# MPI_Barrier          0.0177442       32            0.01          0.00
# MPI_Reduce           0.00540609      32            0.00          0.00
# MPI_Comm_rank        0.00465156      32            0.00          0.00
# MPI_Comm_size        0.000145341     32            0.00          0.00
```

# IPM Examples

Click on the metric you are want.

## NERSC job details

http://www.nersc.gov/REST/jobs/job_details.php?stepid=732423.sdb&timestamp=1313679078&completion=1313679081

### IPM Summary

| | | | | | |
|---|---|---|---|---|---|
| Executable | | | | | ./wrf.exe |
| Number of tasks | 512 | Aggregate GFlop/sec | 0.1482 | Average GFlop/sec/task | 0.0003 |
| Average wall secs | 8.861e+01 | Aggregate memory (GB) | 32.2626 | Average memory/task (GB) | 0.0630 |
| Average MPI secs/task | 7.898e+01 | MPI time % | 89.14 | Aggregate MPI calls made | 7.027e+07 |

### IPM Summary Statistics - 512 tasks

| Metric | Sum over all tasks | Average (per task) | Task CV (%) | Task Minimum | Task Maxium |
|---|---|---|---|---|---|
| Aggregate Floating Point Operations (Flop x 10**9) | 1.313e+01 | 2.565e-02 | 6.10 | 1.713e-02 | 2.758e-02 |
| GFlop/sec | 1.482e-01 | 2.895e-04 | 6.10 | 1.934e-04 | 3.114e-04 |
| Maximum Memory Usage (GBytes) | 3.226e+01 | 6.301e-02 | 10.12 | 5.701e-02 | 1.947e-01 |
| Time Spent in MPI Routines (sec) | 4.044e+04 | 7.898e+01 | 4.05 | 9.801e+00 | 8.359e+01 |
| Wallclock Time (sec) | 4.537e+04 | 8.861e+01 | 0.10 | 8.848e+01 | 8.895e+01 |

Memory in units of gigabytes; time in seconds.

### Hardware counter statistics - 512 tasks

| Counter Name | Sum over all tasks | Average (per task) | Task CV (%) | Task Minimum | Task Maxium |
|---|---|---|---|---|---|
| PAPI_FP_OPS | 1.161799e+12 | 2.269139e+09 | 6.09 | 1.515023e+09 | 2.439529e+09 |

### MPI Time Statistics - 512 tasks

| Call | Sum over all tasks | Average (per task) | Task CV (%) | Task Minimum | Task Maximum | % of MPI | % of wall |
|---|---|---|---|---|---|---|---|
| MPI_Bcast | 3.517e+04 | 6.869e+01 | 4.48 | 4.342e-01 | 7.269e+01 | 86.969 | 77.520 |
| MPI_Scatterv | 2.589e+03 | 5.057e+00 | 5.79 | 1.059e+00 | 5.540e+00 | 6.403 | 5.707 |
| MPI_Wait | 2.176e+03 | 4.249e+00 | 17.82 | 1.250e+00 | 4.968e+00 | 5.380 | 4.795 |
| MPI_Gatherv | 4.312e+02 | 8.422e-01 | 36.44 | 3.552e-03 | 2.271e+00 | 1.066 | 0.950 |
| MPI_Isend | 5.250e+01 | 1.025e-01 | 11.96 | 7.182e-02 | 1.259e-01 | 0.130 | 0.116 |
| MPI_Irecv | 1.033e+01 | 2.017e-02 | 10.21 | 1.217e-02 | 2.613e-02 | 0.026 | 0.023 |
| MPI_Gather | 1.021e+01 | 1.995e-02 | 502.07 | 1.391e-03 | 1.434e+00 | 0.025 | 0.023 |
| MPI_Comm_rank | 4.563e-01 | 8.913e-04 | 4.74 | 7.799e-04 | 1.404e-03 | 0.001 | 0.001 |
| MPI_Comm_size | 9.629e-02 | 1.881e-04 | 10.65 | 1.462e-04 | 4.859e-04 | 0.000 | 0.000 |
| MPI_Init | 0.000e+00 | 0.000e+00 | 0.000e+00 | 0.000e+00 | 0.000 | 0.000 | |
| MPI_Finalize | 0.000e+00 | 0.000e+00 | 0.000e+00 | 0.000e+00 | 0.000 | 0.000 | |

**Average MPI Time per Task**

- MPI_Bcast
- MPI_Scatterv
- MPI_Wait
- MPI_Gatherv

# IPM Examples



NERSC job details

http://www.nersc.gov/REST/jobs/ipm_summary.php?stepid=619349.sdb&name=gflops&timestamp=1310766809

Task distribution of IPM summary statistics for JobID 619349.sdb

| Metric | Sum | Mean | Std. Dev. | CV (%) | Minimum | Maximum |
|---|---|---|---|---|---|---|
| Aggregate Floating Point Operations (Flop x 10**9) | 3.011e+02 | 1.470e-01 | 4.946e-03 | 3.36e+00 | 1.395e-01 | 2.161e-01 |
| GFlop/sec | 6.147e-01 | 3.002e-04 | 1.008e-05 | 3.36e+00 | 2.847e-04 | 4.411e-04 |
| Maximum Memory Usage (GBytes) | 4.101e+02 | 2.002e-01 | 9.606e-03 | 4.80e+00 | 1.781e-01 | 2.448e-01 |
| Time Spent in MPI Routines (sec) | 1.228e+06 | 5.995e+02 | 4.984e+01 | 8.31e+00 | 5.177e+02 | 6.801e+02 |
| Wallclock Time (sec) | 1.003e+06 | 4.898e+02 | 6.428e-02 | 1.31e-02 | 4.898e+02 | 4.927e+02 |

CV = Coefficient of Variance = (Standard Deviation / Mean)

**Task distribution of** *Aggregate Floating Point Operations (Flop x 10**9)* **- as a percentage of maximum**

The MPI rank is the sum of the column and row indices in the table.

*gflops* vs. MPI Rank

# IPM Examples

# IPM Examples

# IPM Examples

- **What tools do you use?**
- **What tools do you want?**
- **What would you like centers to support?**
- **Can you get to exascale without tools?**

- **Users are asking for tools because HPC systems and programming models are changing**

- **More and more components to worry about**

  – CPU (caches, FPUs, pipelining, …)

  – Data movement to main memory, GPU memory, levels of cache

  – I/O

  – Network (message passing)

  – CPU Threads (OpenMP)

  – GPU performance

# What I Want in a Tool

- Let the users help themselves
- Work for everyone all (most of?) the time
- Easy to use
- Useful
- Easy to interpret the results
- Affordable ($$ or manpower support costs)
- Simple, supplement existing complex tools
  - Point the way for a "deeper dive" in problem areas

U.S. DEPARTMENT OF ENERGY | Office of Science

BERKELEY LAB

**National Energy Research Scientific Computing Center**