



**FAST SOLUTIONS**

- PAPI\_L1\_DCM
- PAPI\_L1\_ICM
- PAPI\_L2\_DCM
- PAPI\_L2\_ICM
- PAPI\_L2\_TCM
- PAPI\_L2\_TCM

# NPB3.3-MZ-MPI/BT tutorial example OpenMP+MPI application (Curie version)

Brian Wylie  
Jülich Supercomputing Centre

b.wylie@fz-juelich.de

April 2012

- Familiarise with usage of VI-HPS tools
  - complementary tools' capabilities & interoperability
- Prepare to apply tools productively to *your* application(s)
- Exercise is based on a small portable benchmark code
  - unlikely to have significant optimization opportunities
- Optional (recommended) exercise extensions
  - analyze performance of alternative configurations
  - investigate effectiveness of system-specific compiler/MPI optimizations and/or placement/binding/affinity capabilities
  - investigate scalability and analyze scalability limiters
  - compare performance on different HPC platforms
  - ...

<b>System</b>	<b><i>curie</i></b>	
Domain	ccc.cea.fr	
Vendor	Bull	
Network	Infiniband	
	(fat nodes)	(thin nodes)
<b>Processors</b>	Intel X7560	Intel E5-2680
Frequency	2.26 GHz	2.7 GHz
<b>Compute nodes</b>	360	5040
Chips per node	4	2
Cores per chip	8	8
Threads per core	2	2
Memory per node	128 GB	64 GB

<b>System</b>	<i>curie</i>
domain	ccc.cea.fr
<b>Filesystem</b>	<i>Lustre</i>
Parallel filesys	\$WORKDIR
<b>Compiler</b>	<i>Intel</i>
OpenMP flag	-openmp
<b>MPI</b>	<i>Bullx</i>
C compiler	mpicc
C++ compiler	mpicxx
F77 compiler	mpif77
F90 compiler	mpif90
<b>Queue</b>	<i>SLURM</i>
job submit	ccc_msub job
list jobs	qstat

- Environment managed via modules

```
% module list
cc
intel/12.1.7.256          # Intel Compiler Suite 12.1.7.256
bullxmpi/1.1.14.1         # Bullx MPI
```

- other compiler and MPI library versions also available
- along with VI-HPS and other tools that can be loaded

k cachegrind/0.7.1	valgrind/3.6.0	likwid/2.3.0
periscope/default	vampir/7.3	mpc/2.3.0
scalasca/1.3.3	vampir/7.5	mpc/2.3.1
scalasca/1.4.1	vampirtrace/5.10.1	papi/4.1.4
tau/2.21.2	vampirtrace/5.12.2	

- Tutorial sources should be copied to your own directory where you can then work on them

```
% cp -r ~wylieb/tutorial/NPB3.3-MZ-MPI $WORKDIR
```

- NAS Parallel Benchmark suite (sample MZ-MPI version)
  - Available from <http://www.nas.nasa.gov/Software/NPB>
  - 3 benchmarks (all in Fortran77, using OpenMP+MPI)
  - Configurable for various sizes & classes
- Move into the NPB3.3-MZ-MPI root directory

```
% cd NPB3.3-MZ-MPI; ls
BT-MZ/  LU-MZ/  SP-MZ/
bin/    common/ config/ jobsctpt/      Makefile          README   sys/
```

- Subdirectories contain source code for each benchmark
  - plus additional configuration and common code
- The provided distribution has already been configured for the tutorial, such that it's ready to "make" benchmarks and install them into a (tool-specific) "bin" subdirectory

- Type “make” for instructions

```
% make  
=====  
= NAS Parallel Benchmarks 3.3 =  
= MPI+OpenMP Multi-Zone versions =  
=====
```

To make a NAS multi-zone benchmark type

```
make <benchmark-name> CLASS=<class> NPROCS=<number>
```

To make a set of benchmarks, create the file config/suite.def according to the instructions in config/suite.def.template and type

```
make suite
```

```
*****  
* Custom build configuration is specified in config/make.def *  
* Suggested tutorial benchmark specification: *  
*   make bt-mz CLASS=B NPROCS=4 *  
*****
```

- Specify the benchmark configuration
  - benchmark name: **bt-mz**, lu-mz, sp-mz
  - the number of MPI processes: NPROCS=4
  - the benchmark class (S, W, A, B, C, D, E, F): CLASS=B

```
% make bt-mz CLASS=B NPROCS=4
cd BT-MZ; make CLASS=B NPROCS=4 VERSION=
gmake: Entering directory 'BT-MZ'
cd .../sys; cc -o setparams setparams.c -lm
..../sys/setparams bt-mz 4 B
mpif77 -c -O3 -[f]openmp bt.f
...
mpif77 -c -O3 -[f]openmp setup_mpi.f
cd .../common; mpif77 -c -O3 -[f]openmp print_results.f
cd .../common; mpif77 -c -O3 -[f]openmp timers.f
mpif77 -O3 -[f]openmp -o ../bin/bt-mz_B.4 \
    bt.o make_set.o initialize.o exact_solution.o exact_rhs.o \
    set_constants.o adi.o define.o copy_faces.o rhs.o solve_subs.o \
    x_solve.o y_solve.o z_solve.o add.o error.o verify.o setup_mpi.o \
    ../common/print_results.o ../common/timers.o
Built executable ../bin/bt-mz_B.4
gmake: Leaving directory 'BT-MZ'
```

- What does it do?
  - Solves a discretized version of unsteady, compressible Navier-Stokes equations in three spatial dimensions
  - Performs 200 time-steps on a regular 3-dimensional grid using ADI and verifies solution error within acceptable limit
  - Intra-zone computation with OpenMP, inter-zone with MPI
- Implemented in 20 or so Fortran77 source modules
- Runs with any number of MPI processes & OpenMP threads
  - bt-mz\_B.4 x8 is reasonable (vary as appropriate)
    - ▶ excess processes idle when run with more than compiled
  - bt-mz\_B.4 x8 should run in around 10 seconds
    - ▶ typically runs more efficiently with more processes than threads
  - CLASS=C does much more work and takes much longer!

- Set OMP\_NUM\_THREADS and launch as an MPI application

```
% cd bin;  OMP_NUM_THREADS=8  mpiexec -np 4 ./bt-mz_B.4
NAS Parallel Benchmarks (NPB3.3-MZ-MPI) - BT-MZ MPI+OpenMP Benchmark
Number of zones:    8 x    8
Iterations:   200      dt:    0.000300
Number of active processes:        4
```

```
Time step    1
Time step    20
Time step    40
Time step    60
Time step    80
Time step   100
Time step   120
Time step   140
Time step   160
Time step   180
Time step   200
```

Verification Successful

BT-MZ Benchmark Completed.  
**Time in seconds = 10.2**

Hint: if running interactively  
save the benchmark output  
(or note the run time)  
to be able to refer to it later

Hint: copy/edit example batch  
scripts from jobscript directory:  
**% ccc\_msub ..../jobscript/run.msub**

- The tutorial steps are similar and repeated for each tool
- Use the provided NPB3.3-MZ-MPI tutorial directory

```
% cd NPB3.3-MZ-MPI; ls  
BT-MZ/ LU-MZ/ SP-MZ/  
bin/ common/ config/ jobsript/ Makefile README sys/
```

- Edit **config/make.def** to adjust build configuration
  - Modify specification of compiler/linker: **MPIF77**
- Make clean and build new tool-specific executable

```
% make clean  
% make bt-mz CLASS=B NPROCS=4  
Built executable ./bin.$(TOOL)/bt-mz_B.4
```

- Change to the directory containing the new executable before running it with the desired tool configuration

```
% cd bin.$(TOOL)  
% export OMP_NUM_THREADS=8 ...  
% mpiexec -np 4 ./bt-mz_B.4
```

Hint: check available scripts:  
% ccc\_msub .../jobsript/\$TOOL.msub

- config/make.def

```
#           SITE- AND/OR PLATFORM-SPECIFIC DEFINITIONS
#-----
# Items in this file may need to be changed for each platform.
...
#OPENMP = -fopenmp  # GCC
OPENMP  = -openmp   # Intel
#-----
# The Fortran compiler used for hybrid MPI programs
#
MPIF77 = mpif77                                     Set flag according to compiler
# Alternative variants to perform instrumentation
#MPIF77 = psc_instrument -t user,mpi  mpif77
#MPIF77 = scalasca -instrument  mpif77
#MPIF77 = tau_f90.sh
#MPIF77 = vtf77 -vt:hyb -vt:f77  mpif77
#MPIF77 = scorep --user  mpif77
# PREP is a generic preposition macro for instrumentation preparation
#MPIF77 = $(PREP)  mpif77
...
```

Default (no instrumentation)

Hint: uncomment one of these alternative compiler wrappers to perform instrumentation ...

... or this for generic variant

# NPB-MZ-MPI/BT analyses with VI-HPS tools

VI-HPS

