Virtual Institute – High Productivity Supercomputing

VI-HPS

SOFTWARE

- ☐ 19.56 updatex
- ☐ 399.70 updateien
- ☐ 0.00 gene
- ☐ 0.00 <<iteration loop>>
- ☐ 447.52 genbc

FAST SOLUTIONS
☑ PAPI_L1_ICM
☐ PAPI_L2_DCM
☑ PAPI_L2_ICM
☐ PAPI_L1_TCM

PRODUCTIVITY

# Periscope
# Tutorial Exercise
# NPB-MPI/BT

M. Gerndt, V. Petkov, Y. Oleynik, R. Mijakovic

Technische Universität München

periscope@in.tum.de

April 2012

- Intermediate-level tutorial example

- Available in MPI, OpenMP, hybrid OpenMP/MPI variants

- Automatic performance properties search with Periscope:
  - Source code instrumentation
    - ▶ MPI calls
  - Automatic search for slow MPI communication patterns
  - Results exploration with Eclipse based GUI

- Manual instrumentation optimization

0. Loading and configuring of Periscope

   ```
   module load periscope
   ```

1. Program instrumentation: **psc_instrument**

2. Periscope analysis: **psc_frontend**

3. Performance properties exploration: Periscope GUI

- Before first use of Periscope, one has to create the configuration file `.periscope` in the home directory. Configuration could be copied from `$PERISCOPE_ROOT`:

```
% cp $PERISCOPE_ROOT/etc/periscope.sample ~/.periscope
```

- It should look like:

```
MACHINE          = curie50          //hostname
SITE             = UVSQ
REGSERVICE_HOST  = curie50          //host of registry
REGSERVICE_PORT  = 50001            //port of the registry
APPL_BASEPORT    = 51000            //first port for application
AGENT_BASEPORT   = 50002            //first port agent hierarchy
```

- The Periscope agents and the application processes register with a `registry`. The registry is started via:

```
% psc_regsrv &
```

- To enable performance measurement, the program has to be instrumented. This is done with `psc_instrument`:

```
% psc_instrument
Periscope Source-to-Source Instrumentation Wrapper
Usage: psc_instrument [-t regions] [-n] [-s sir] [-v] [-d] compiler
        [options] file [libs]
-t Types of regions to instrument separated by spaces
        (e.g. -t "user loop call")
-s Filename for the resulting SIR file (default: appl.sir)
-v Verbose output
-d Debug mode: keeps the instrumented source files
        after the compilation
-n Prints each step of the compilation instead of executing them
```

● Substitute compile/link commands in Makefile definitions (`config/make.def`) with `psc_instrument`:

```
MPIF77 = psc_instrument -s ${PROGRAM}.sir -t user,mpi mpif77

FLINK = $(MPIF77)
FFLAGS = -O

mpi-bt: $(OBJECTS)
        $(FLINK) $(FFLAGS) -o mpi-bt $(OBJECTS)
.f.o:
        $(MPIF77) $(FFLAGS) -c $<
```

- ## Return to root directory and clean-up

```
% make clean
```

- ## Re-build BT with the original command

```
% make bt-mz CLASS=B NPROCS=4
   ==========================================
   =      NAS Parallel Benchmarks 3.3      =
   =      MPI/F77/C                        =
   ==========================================
cd BT-MZ; make NPROCS=4 CLASS=B SUBTYPE= VERSION=
make[1]: Entering directory `BT-MZ'
...
psc_instrument -s bt.sir -t "user loop call" mpif77 -c  -O -g bt.f
psc_instrument -s bt.sir -t "user loop call" mpif77 -c  -O -g make_set.f
…
psc_instrument -s bt.sir -t "user loop call" mpif77 -O \
-o ../bin.periscope/bt-mz_B.4 bt.o ...
Built executable ../bin.periscope/bt-mz_B.4
make[1]: Leaving directory `BT-MZ'
```

- ## Change directory to bin.periscope

```
% cd bin.periscope
```

- Periscope is started via the frontend. It automatically starts application and hierarchy of analysis agents.
- Run `psc_frontend --help` for brief usage information

```
% psc_frontend --help
Usage: psc_frontend <options>
  [--help]                  (displays this help message)
  [--quiet]                 (do not display debug messages)
  [--registry=host:port]   (address of the registry service, optional)
  [--port=n]                (local port number, optional)
  [--maxfan=n]              (max. number of child agents, default=4)
  [--timeout=secs]          (timeout for startup of agent hierarchy)
  [--delay=n]                 (search delay in phase executions)
  [--appname=name]
  [--apprun=commandline]
  [--mpinumprocs=number of MPI processes]
  [--ompnumthreads=number of OpenMP threads]
…
  [--strategy=name]
  [--sir=name]
  [--phase=(FileID,RFL)]
  [--debug=level]
```

- Run Periscope analysis by executing `psc_frontend` with the following command

```
% psc_frontend --sir=bt-mz_B.4.sir --apprun=./bt-mz_B.4 --strategy=MPI
--mpinumprocs=4
[psc_frontend][DBG0:fe] Agent network UP and RUNNING. Starting search.

 NAS Parallel Benchmarks 3.3 -- BT Benchmark
 [...]
 Time step 200
 BT Benchmark Completed.

 ----------
End Periscope run! Search took 60.5 seconds (33.3 seconds for startup)
```

- Copy and change the Periscope batch script

```
#!/bin/bash
# submit from ./bin.periscope directory with "ccc_msub psc.msub"
#MSUB -r npb_btmz_psc
#MSUB -o npb_btmz_%I.oe
#MSUB -e npb_btmz_%I.oe
#MSUB -n 4       # number of MPI processes
#MSUB -c 8       # number of OpenMP threads/process
#MSUB -T 600     # max walltime in seconds
#MSUB -x         # allocate exclusive nodes
#MSUB -A tgcc0007 # project id

cd $BRIDGE_MSUB_PWD

# benchmark configuration
export OMP_NUM_THREADS=$BRIDGE_MSUB_NCORE
PROCS=$BRIDGE_MSUB_NPROC
CLASS=B
EXE=./bt-mz_$CLASS.$PROCS

# remember to check that psc_regsrv is running!
psc_frontend --apprun=$EXE --mpinumprocs=$PROCS --strategy=MPI
```

- Submit the script with `ccc_msub psc.msub`

- When running Eclipse *from* LiveDVD copy the .psc file to your local tutorial folder:

```
%scp <username>@curie.ccc.cea.fr:tutorial/NPB3.3-MZ-MPI/bin.periscope/*.psc
tutorial/NPB3.3-MZ-MPI/bin.periscope

%scp <username>@curie.ccc.cea.fr:tutorial/NPB3.3-MZ-MPI/bin.periscope/*.sir
tutorial/NPB3.3-MZ-MPI/bin.periscope
```

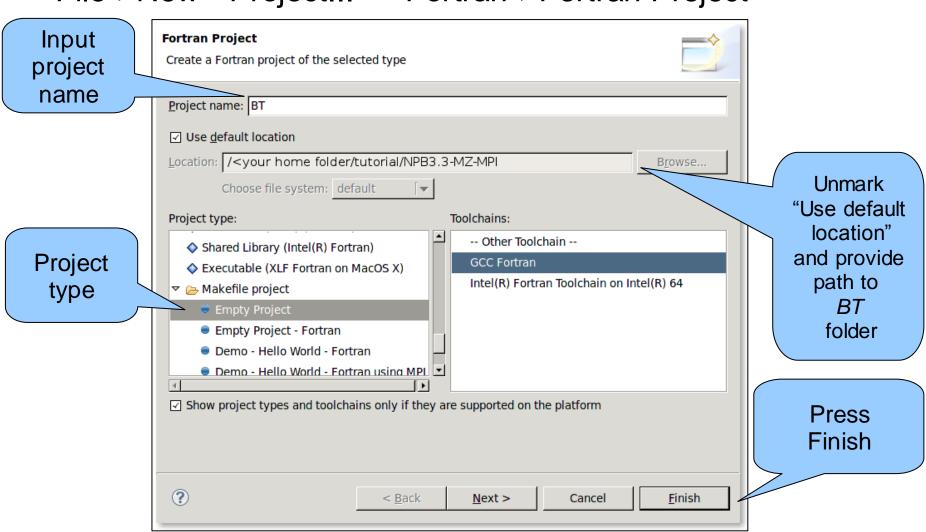- Start `Eclipse` with Periscope GUI from console

```
% eclipse
```

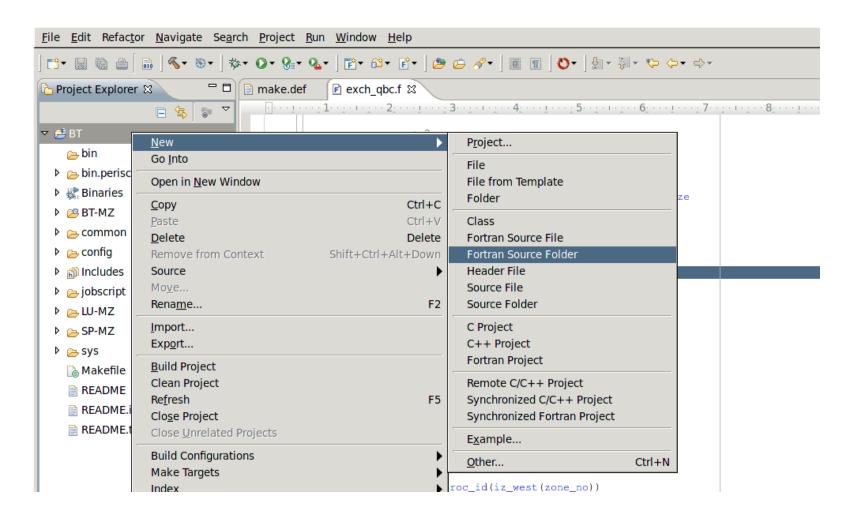- Or by double-click on Eclipse pictogram on the Desktop

- File->New->Project... → Fortran->Fortran Project

**Input project name**

**Fortran Project**
Create a Fortran project of the selected type

Project name: BT

☑ Use default location

Location: /<your home folder/tutorial/NPB3.3-MZ-MPI    Browse...

Choose file system: default ▼

Project type:

**Project type**

- ◇ Shared Library (Intel(R) Fortran)
- ◇ Executable (XLF Fortran on MacOS X)
- ▽ 📂 Makefile project
  - 🔵 Empty Project
  - 🔵 Empty Project - Fortran
  - 🔵 Demo - Hello World - Fortran
  - 🔵 Demo - Hello World - Fortran using MPI

Toolchains:

- -- Other Toolchain --
- GCC Fortran
- Intel(R) Fortran Toolchain on Intel(R) 64

☑ Show project types and toolchains only if they are supported on the platform

**Unmark "Use default location" and provide path to *BT* folder**

**Press Finish**

⊘    < Back    Next >    Cancel    Finish

- Right-click -> File-> New -> Fortran Source Folder

- Choose BT-MZ as a source folder

# Loading properties



Expand BT project,
search for *.psc
and
Right click->Periscope->
Load all properties
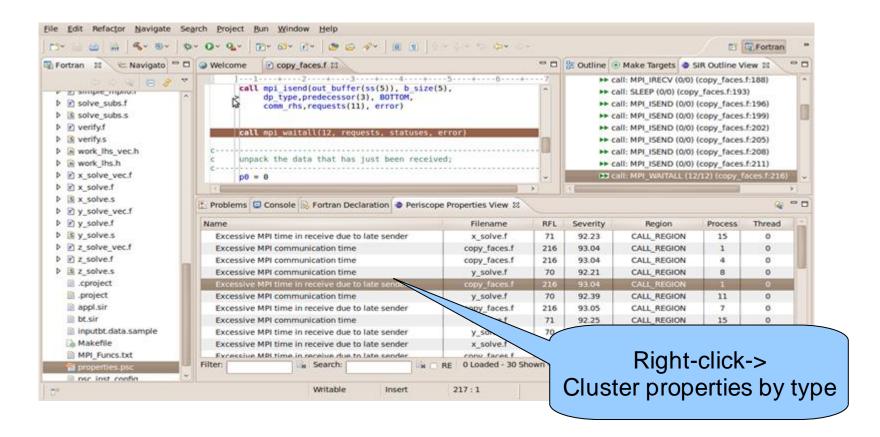
# Periscope GUI

- Multi-functional table is used in the GUI for Eclipse for the visualization of bottlenecks
  - Multiple criteria sorting algorithm
  - Complex categorization utility
  - Searching engine using Regular Expressions
  - Filtering operations
  - Direct navigation from the bottlenecks to their precise source location using the default IDE editor for that source file type (e.g. CDT/Photran editor).

- SIR outline view shows a combination of the standard intermediate representation (SIR) of the analysed application and the distribution of its bottlenecks. The main goals of this view are to assist the navigation in the source code and attract developer's attention to the most problematic code areas.
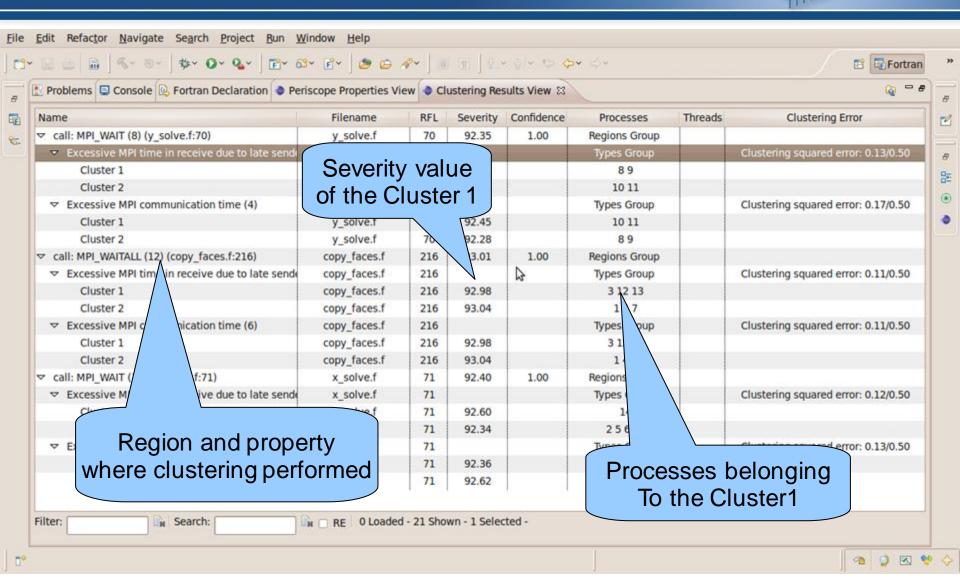
# Properties clustering

- Clustering can effectively summarize displayed properties and identify a similar performance behaviour possibly hidden in the large amount of data
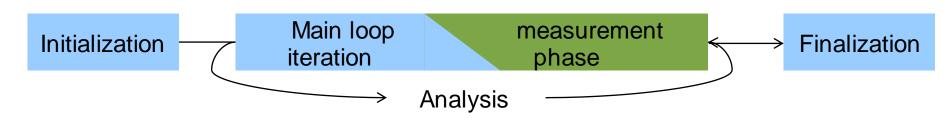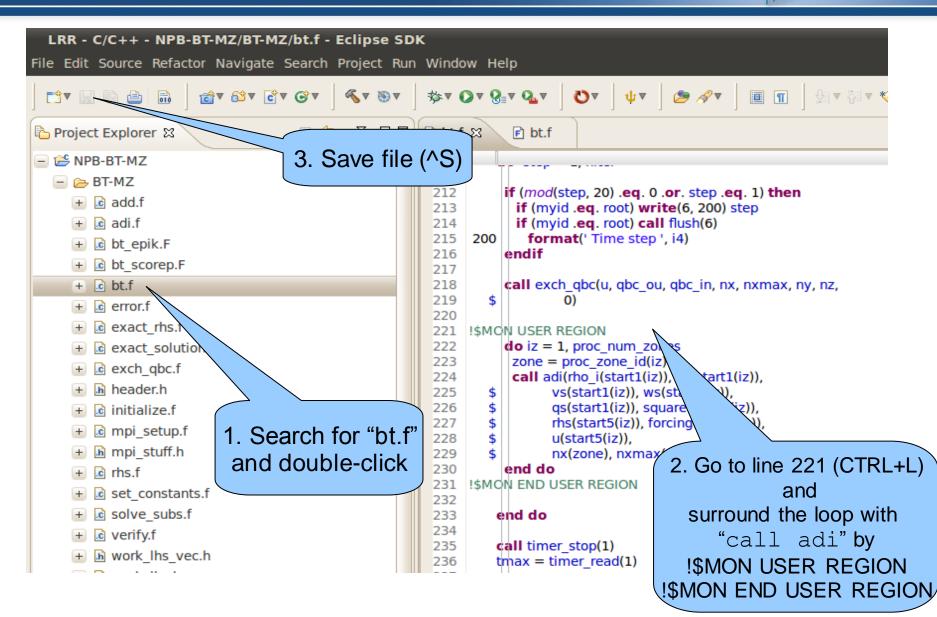


Right-click->
Cluster properties by type

# Properties clustering

- Periscope performs multiple iterative performance measurement experiments on the basis of *Phases:*
  - All measurements are performed inside phase
  - Begin and end of phase are global synchronization points
- By default phase is the whole program
  - Needs restart if multiple experiments required (single core performance analysis strategies require multiple experiments)
  - Unnecessary code parts also measured
- User specified region marked with `!$MON USER REGION` and `!$MON END USER REGION` will be used as phase:
  - Typically main loop of application → no need for restart, faster analysis
  - Unnecessary code parts are not measured → less measurements overhead
  - Severity value is normalized on the main loop iteration time → more precise performance impact estimation

| Initialization | Main loop iteration | measurement phase | Finalization |
|---|---|---|---|

Analysis

- Return to root directory and clean-up

```
% make clean
```

- Re-build BT with the original command

```
% make bt CLASS=W NPROCS=16
```

- Change directory into location of executable

```
% cd bin.periscope
```

- Re-run Periscope analysis by executing `psc_frontend`

```
% psc_frontend --sir=bt-mz_W.4.sir --apprun=./bt-mz_W.4 --strategy=MPI
--mpinumprocs=4 --force-localhost
[psc_frontend][DBG0:fe] Agent network UP and RUNNING. Starting search.


 NAS Parallel Benchmarks 3.3 -- BT Benchmark
 [...]
 Time step 1
 BT Benchmark Completed.

 ----------
End Periscope run! Search took 37.2 seconds (33.3 seconds for startup)
```

- Only 1 iteration of BT required instead of 200 previous run!
- Frontend will overwrite the properties found into the file `properties_MPI_*.psc` in the current directory, which again need to be copied into the BT source directory

```
% cp properties_MPI_*.psc ../BT
```

- Re-load `properties_MPI_*.psc` in Periscope GUI. Now found properties should have more precise severities values