# Periscope
# Tutorial Exercise
# NPB-MPI/BT

M. Gerndt, Y. Oleynik, V. Petkov

Technische Universität München

periscope@in.tum.de

September 2011

- Intermediate-level tutorial example

- Available in **MPI**, OpenMP, hybrid OpenMP/MPI variants
  - also MPI File I/O variants (collective & individual)

- Automatic performance properties search with Periscope:
  - Source code instrumentation
    - ▸ Loops, MPI & application function calls
  - Automatic search for slow MPI communication patterns
  - Results exploration with Eclipse based GUI

- Manual instrumentation optimization

0. Configuration of 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          = localhost          // hostname
SITE             = VIHPS
REGSERVICE_HOST  = cluster-beta       // host where the registry is running
REGSERVICE_PORT  = 50001              // please choose a random port!
APPL_BASEPORT    = 51000              // first port for application
AGENT_BASEPORT   = 50002              // first port agent hierarchy
```

- Install GUI into eclipse from http://www.lrr.in.tum.de/periscope/eclipse/ or use the eclipse with pre-installed GUI available with `module load periscope`

- The Periscope agents and the application processes register with a `registry`. It is started in batch 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
        -i Force Intel compilers
```

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

```
MPIF77 = psc_instrument -i -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 (B or W version)

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

   =========================================
   =      NAS Parallel Benchmarks 3.3      =
   =      MPI/F77/C                        =
   =========================================
cd BT; make NPROCS=16 CLASS=B SUBTYPE= VERSION=
make[1]: Entering directory `BT'
...
psc_instrument -i -s ... -t user,mpi mpif77 -c  -O -g bt.f
psc_instrument -i -s ... -t user,mpi mpif77 -c  -O -g make_set.f
…
psc_instrument -i -s ... -t "user loop call" mpif77 -O \
-o ../bin.periscope/bt_B.16 bt.o ...
Built executable ../bin.periscope/bt_B.16
make[1]: Leaving directory `BT'
```

- ## 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 in the batch script psc.lsf/.msub

```
% bsub < psc.lsf    or     msub psc.msub
Check the job output:
% bpeek
[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 37.57 seconds (33.09 seconds for startup)
```

- Frontend will write the detected properties into the file `properties_MPI_<PID>.psc` in the current directory. It should be copied into the BT source directory

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

- Start `Eclipse` with Periscope GUI from console

```
% eclipse &
```

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

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

# Loading properties

**VI-HPS**

Periscope - http://www.eclipse.org/photran/welcome/6.0/?version=6.0.0.201102161000&os=linux&arch=x86_64&lastVersion= - Eclipse

File   Edit   Refactor   Navigate   Search   Project   Run   Window   Help

Periscope

Welcome

Periscope SIR Outli | Project Explorer

- x_solve.f
- x_solve.o - [x86_64/le]
- y_solve_vec.f
- y_solve.f
- y_solve.o - [x86_64/le]
- z_solve_vec.f
- z_solve.f
- z_solve.o - [x86_64/le]
- bt.sir
- inputbt.data.sample
- Makefile
- properties.psc
- psc_inst_config
- test.cmi

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

Periscope Properties View | Clustering Results View | Con

| Name | Filename | Processes |
|------|----------|-----------|
| No properties clustered! | | |

Filter: [____]   Search: [____]   ☐ RE   0 Loaded - 1 Shown - 0 Selected -

# 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

# Properties clustering

# Plotting clusters

- 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=B NPROCS=16
```

- ## Change directory into location of executable

```
% cd bin.periscope
```

- Re-run Periscope analysis by submitting the script

```
% bsub < psc.lsf    or     msub psc.msub
Check the job output:
% bpeek
[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_<PID>.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_<PID>.psc` in Periscope GUI. Now found properties should have more precise severities values