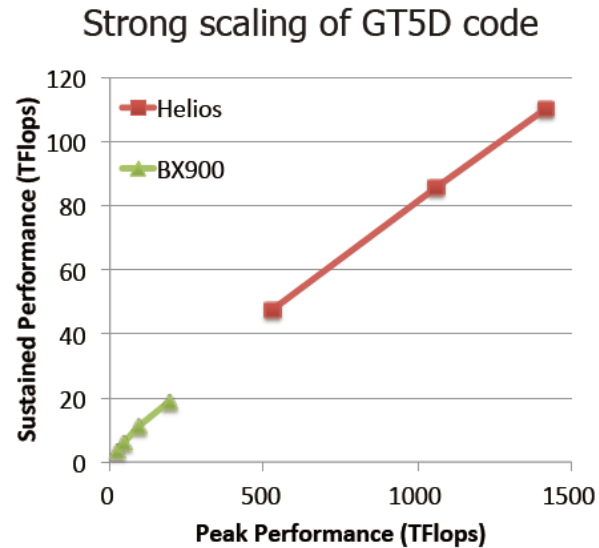# Report of JA large jobs

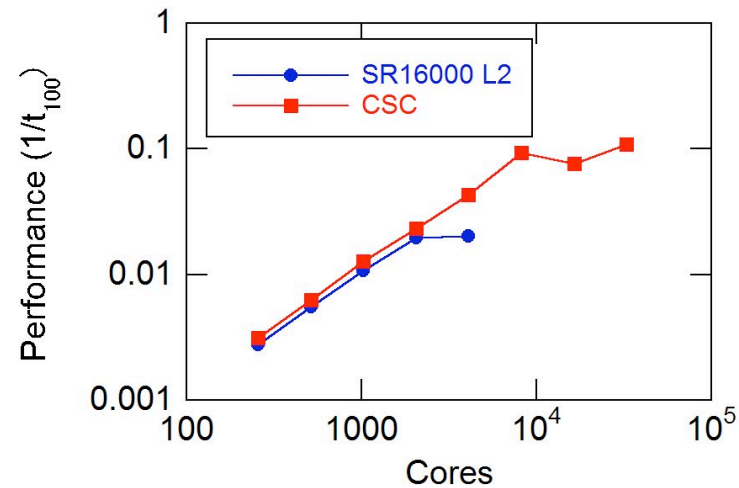IFERC PL: N.Nakajima (NIFS)

1. Background
   - Experiences in the Lighthouse projects
2. Results of the large job session
   - Large job session
   - GT5D
   - MEGA
   - MIPS
3. Summary of large jobs
4. Measures up to coming March

# 1. Background

Experiences in the Lighthouse projects

### Strong scaling of GT5D code



*Strong scaling for GT5D*



*Strong scaling for MEGA*

- 34 M (=256x512x256) grids
- 34 M particles for α and D beam, w/ FLR
- Good scaling up to 8192 cores (512 nodes)

It was very hard to implement large jobs using a large fraction of the system.
It is very important to smoothly implement large jobs for Japanese important
codes; GT5D and MEGA proposed as the benchmark codes, as well as EU
codes; GENE and ORB5.

## 2. Results of large job session

Large job session :

- To clarify problems on large jobs with almost full system and to improve situations by using a most clean configuration
  - ☐ To contribute to improvement of system, hardware, and software configurations
  - ☐ To select adequate parameters in mpi-library and so on
- From 2012/07/25, under the support of CSC and HPC support teams
- Target codes : GT5D, MEGA, and MIPS

Hereafter, the results of the large job session will be presented.

Note that the meanings of "slowness" and "scattering" of elapsed times are fairly intuitive, and their rough definitions are as follow;

- "slowness" of the elapsed times; a deviation of the elapsed times from an expected scaling low
- "scattering" of the elapsed times; a wide scattering of the elapsed times mainly depending on the timing of job submission

## 2. Results of large job session - GT5D -

Period : 2012/9-10 (tests only in this period are summarized)

- Inevitable workarounds:
  - See Helios User Manual > FAQ > Jobs > 16, and
  - See Helios User Manual > FAQ > Jobs > 17.

- mpi-library
  - ☐ Bullxmpi: all jobs successfully done. (512, 1024, 2048, 4096 nodes)
  - ☐ Intelmpi: sometimes hanged up at allreduce, failed for a long nodelist.

- Critical issues for both libraries
  - ☐ slowness of elapsed times
  - ☐ scattering of elapsed times

- Future measures
  - ☐ To compare source code with running results for slowness
  - ☐ To direct monitor executing jobs in the large job session by HPC team for scattering of elapsed times

# 2. Results of large job session - MEGA -

Period : 2012/8-9 (tests only in this period are summarized)

- Inevitable workarounds:

  See Helios User Manual > FAQ > Jobs > 16, and
  See Helios User Manual > FAQ > Jobs > 17.
  [Jobs with 4096 nodes become possible under those workaround.]

- mpi-library
  - ☐ Bullxmpi: all jobs successfully done. (1024, 2048, 4096 nodes)
  - ☐ Intelmpi: sometimes failed with "unexpected disconnect completion event"

- Critical issues for both libraries
  - ☐ slowness of elapsed times (independent of existence of I/O)
  - ☐ scattering of elapsed times

- Future measures
  - ☐ To compare the source code with running results for slowness
  - ☐ To direct monitor executing jobs in the large job session by HPC team for scattering of elapsed times

## 2. Results of large job session - MIPS -

Period : 2012/5-11

- Inevitable workarounds:
  - See Helios User Manual > FAQ > Jobs > 16, and
  - See Helios User Manual > FAQ > Jobs > 17.
- mpi-library
  - ❑ Bullxmpi: Jobs with full cores could run after 2012/11.
    - slowness for 2048 nodes disappears after 2012/10.
  - ❑ Intelmpi: Jobs with full cores could run after 2012/11. However, sometimes failed with "SOCKOPT ERR Connection refused"
- Tuning of source
  - ❑ Read from all process    read by rank0 + broadcast (10 times faster)
- Critical issues for both libraries
  - ❑ slowness of elapsed times
  - ❑ scattering of elapsed times
- Future measures
  - ❑ To compare the source code with running results for slowness
  - ❑ To direct monitor executing jobs in the large job session by HPC team for scattering of elapsed times

# 3. Summary of Large jobs

Large job session :

- To clarify problems on large jobs with almost full system and to improve situations by using most clean configuration
- From 2012/07/25, under the support of CSC and HPC support teams
- Target codes : GT5D, MEGA, and MIPS

Improvements have been obtained especially after 2012/10 - 11:

- Improvement of system [Workarounds in shell script (see FAQ)]
    - See Helios User Manual > FAQ > Jobs > 16, and
    - See Helios User Manual > FAQ > Jobs > 17.
- Continuous improvement of configurations
- Several patches of Luster file system

Very recent situations :

- bullxmpi (OK), intelmpi (sometimes failed, intrinsic problems?)
- Critical issues: slowness and scattering of elapse times for both libraries

# 4. Measures up to coming March

Report meeting of CSC activities will be held on 3/20 (just before the PC-12 to be held on 3/21-22).

- Results of the LHP will be reported.
- Situation of large jobs and scaling of typical codes, GENE, ORB5, GT5D and MEGA will be reported. Implementation for GT5D and MEGA (GENE and ORB5) will be done by CSC and HPC (HLST) teams.

Expected near future improvements

- Application of memory alignment to bullxmpi will be default on 11/27.
- Some hardware will be replaced and software of Lustre file system will be updated in 2012/12 or 2013/1.

For source codes

- MIPS has already been given to HPC team.
- GT5D and MEGA have been sent to CSC support team.
- GT5D has been sent to HPC team very recently, and MEGA will be sent to HPC team with some notices very soon.

# 4. Measures up to coming March

Issues to be solved up to coming March
- scattering of elapsed times (for getting stable operation)

    (Effect of memory alignment to bullxmpi on issues is not clear now.)
- slowness of elapsed times (for improving scaling properties)


Measures to solve issues reflecting users opinions
- to directly monitor job processing in the large job session by HPC team, and to compare running results with source codes,
- to submit in the large job session a set of the same jobs in every month, to monitor the change in the situation of the system continuously,
- to exchange information on large jobs among CSC and HPC teams, and HLST,
- to make more close communication with developing team of Bull in France,
- to re-consider the collaboration method with intel for improvement of intelmpi library (how to treat the relation with BA intellectual properties),
- to use some budget for improvement, if possible and necessary.