JÜLICH
FORSCHUNGSZENTRUM

# Scalasca components with reuse potential

July 21th 2008

Bernd Mohr, Felix Wolf

# Outline

- Overview of Scalasca
- Components with reuse potential
  - OPARI OpenMP source-code instrumenter
  - MPI tracing wrappers and wrapper generator
  - Compiler event adapters
  - Library for efficient parallel file I/O
  - Profile browser
- Ongoing and planned collaborations

# scalasca

- Started in January 2006
- Scalable performance-analysis toolset for parallel codes
  - Emphasis on detection of wait states
- Designed for large-scale systems such as IBM Blue Gene or Cray XT
- Funded through Helmholtz Impulse and Networking Funds
- Developed in cooperation with the University of Tennessee
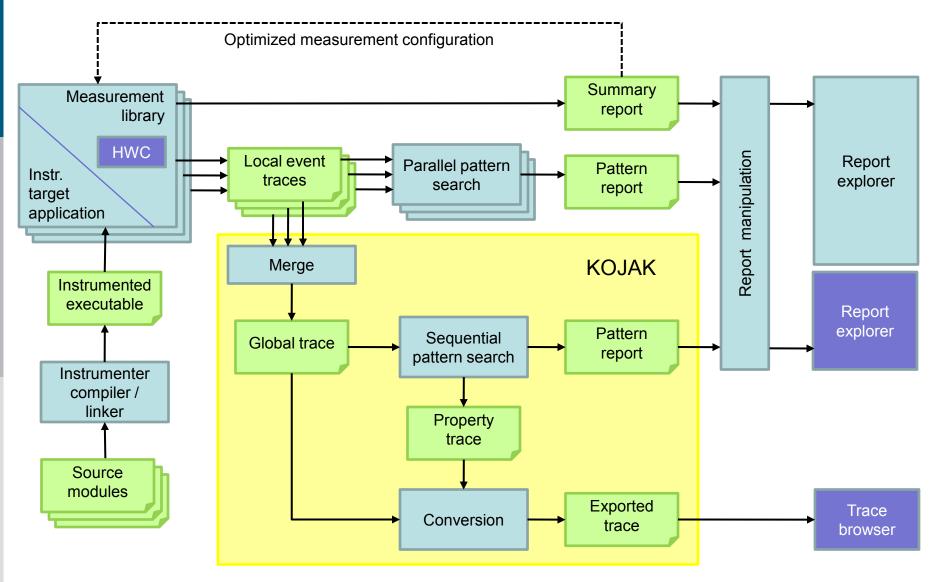- http://www.scalasca.org/

# Functionality

- Integrated performance analysis procedure
  - Runtime summaries (i.e., profiles)
    - *Overview of performance behavior*
    - *Refinement of instrumentation*
  - In-depth study of application behavior via event traces
    - *Localization and quantification of wait states*
  - Switching between both options without recompilation or re-linking
- Programming models supported
  - MPI-1
  - MPI-2 + other one-sided models (in progress)
  - OpenMP (in progress)

# Performance data flow

# OPARI OpenMP source-code instrumenter

- Instruments Fortran, C, C++ OpenMP 2.5 codes with POMP instrumentation calls

- Used by KOJAK, Scalasca, TAU, VampirTrace, ompP

- Not perfect, but works for us

- Ongoing work
  - Removal of limitations
    - *Nested and dynamic threading*
    - *Inter-compilation units dependencies*
  - Support for OpenMP 3.0 features

# MPI tracing wrappers and wrapper generator

- Complete MPI-2 tracing wrappers
  - Enter, Exit, Send, Recv, Collective, Get, Put events
  - C/C++ and Fortran support

- Basis also for Vampirtrace

- Very flexible wrapper generator

- Testsuite

# Compiler event adapters

- Many compilers have (sometimes unsupported and undocumented) options for user function instrumentation
  - GNU, Intel, PGI, Pathscale, IBM XL, Sun f90, NEC, Hitachi
- Used by KOJAK, Scalasca, Vampirtrace

- Compiler event adapter component
  - Translates compiler specific events to generic enter/exit
  - Function filtering at run-time

- Planned
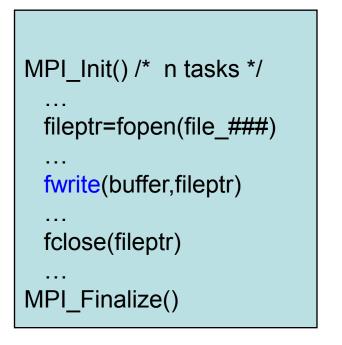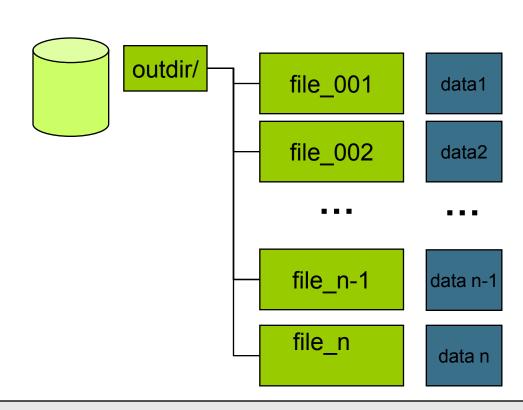  - Function filtering at compile time (GNU, IBM XL)

# Efficient parallel I/O with sionlib

- **S**calable **I**/**O** library for **n**ative parallel file access
- Efficiently reading and writing binary files from thousands of processes, e.g.,
    - Process-local scratch/restart files
    - Process-local trace files (Scalasca)
- Simplified file handling
    - Only one large file instead of thousands of small files
- Optimized I/O
    - Alignment to file system blocks
- Minimal source code changes
    - Allows use of standard file pointer (FILE* fp)
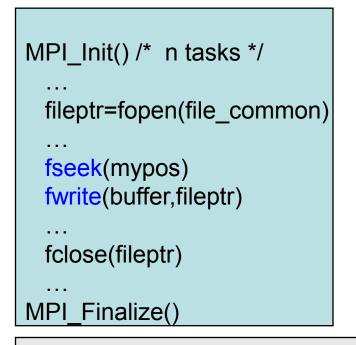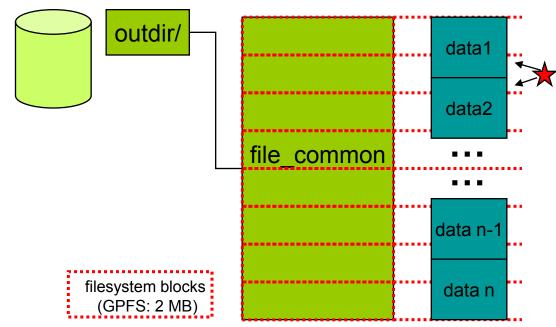
# Typical use case: parallel I/O to separate files

```
MPI_Init() /*  n tasks */
   …
   fileptr=fopen(file_###)
   …
   fwrite(buffer,fileptr)
   …
   fclose(fileptr)
   …
MPI_Finalize()
```

outdir/ — file_001 — data1
file_002 — data2
…    …
file_n-1 — data n-1
file_n — data n

**Problem 1**: file handling (backup, HSM)     ← number of files

**Problem 2**: slow create & open of files     ← Lock on outdir (serialization)

# Example: native parallel direct access

```
MPI_Init() /*  n tasks */
  …
  fileptr=fopen(file_common)
  …
  fseek(mypos)
  fwrite(buffer,fileptr)
  …
  fclose(fileptr)
  …
MPI_Finalize()
```

outdir/

file_common

data1
data2
...
...
data n-1
data n

filesystem blocks
(GPFS: 2 MB)

**Initial Problem solved:** fast open, only one file

**New Problem 1**: meta data handling, start positions and length not stored

**New Problem 2**: file system locks on blocks, overlapping parallel access to blocks

**Restriction**:       space required by each process must be known in advance

# Access with sionlib

```
MPI_Init() /*  n tasks */
  …
sid=sion_paropen_mpi(fname,
    localsize, fsblocksize,…, &fileptr)
…
sion_ensure_free_space(sid, nbytes)
fwrite(buffer,fileptr)
…
sion_ensure_free_space(sid, nbytes)
fwrite(buffer,fileptr)
…
sion_parclose(sid)
…
MPI_Finalize()
```

outdir/

file_common

. . .

metadata

data1

data2

. . .

data n-1

data n

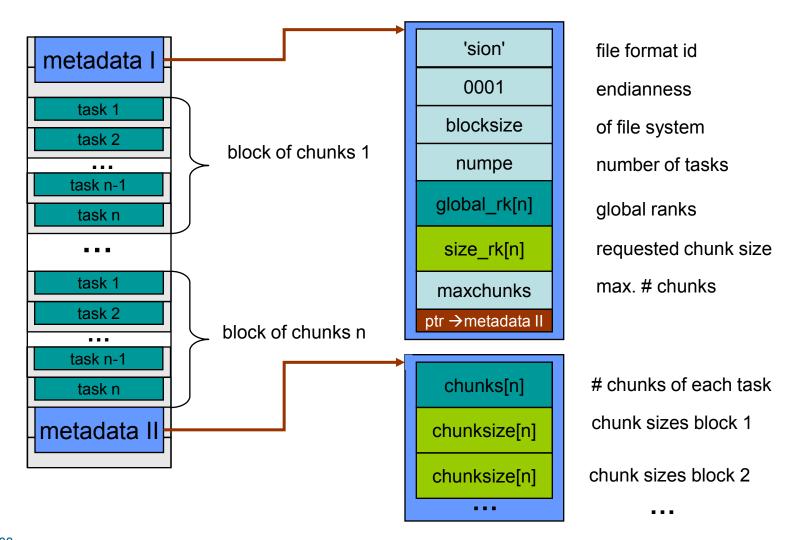filesystem blocks
(GPFS: 2 MB)

**Problems solved**: simple file handling, fast open and fast I/O (fs block alignment)

**Restriction**: space required by each process must be known in advance
→ new allocation at the end of the file if writing more data than
initially allocated

# sionlib: internal file format

| Left column (file layout) | Right column (metadata I) | Description |
|---|---|---|
| metadata I | 'sion' | file format id |
| task 1 | 0001 | endianness |
| task 2 | blocksize | of file system |
| ... | numpe | number of tasks |
| task n-1 | global_rk[n] | global ranks |
| task n | size_rk[n] | requested chunk size |
| block of chunks 1 | maxchunks | max. # chunks |
| ... | ptr →metadata II | |
| task 1 | | |
| task 2 | chunks[n] | # chunks of each task |
| ... | chunksize[n] | chunk sizes block 1 |
| task n-1 | chunksize[n] | chunk sizes block 2 |
| task n | ... | ... |
| block of chunks n | | |
| metadata II | | |

# sionlib: comand line tools

**siondump** [-a] <sionfile>

- prints on stdout all information from the first meta data block , with -a also all chunk sizes from the second meta data block

**sionsplit** [-d digits] <sionfile>  <prefix>

- extracts task related files from a sion file
- a file will be generated for each task with a filename starting with <prefix>
- the task number will be appended to the <prefix>

**siondefrag** [-q blksize] [-s chunksize] <sionfile> <new_sionfile>

- generates a new sion file from an existing sion file
- the new file will have only one chunk per task which contains the data of all chunks of this task in the old sion file
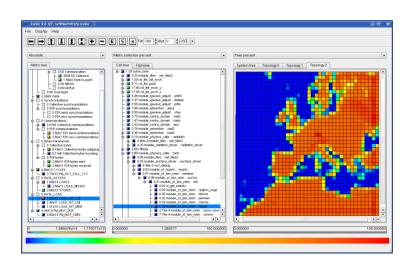- generates with "–q 1" a compact sion file without gaps

# Measurement on 16 rack Blue Gene/P

- BG/P connected to file server with 128 x 10 GiE GPFS file system bandwidth: ~ 6GB/s

- Parallel test: (file server in production)
  - Writing and reading 2 TB data, 32 MB from each task
  - 65536 MPI-tasks, 128 I/O-nodes
  - Parallel open of one SION file &rarr; ~ 1s
  - Overall write bandwidth &rarr; 3.7 GB/s
    550s for writing 2 TB
  - Overall read bandwidth &rarr; 5.4 GB/s
    380s for reading 2 TB
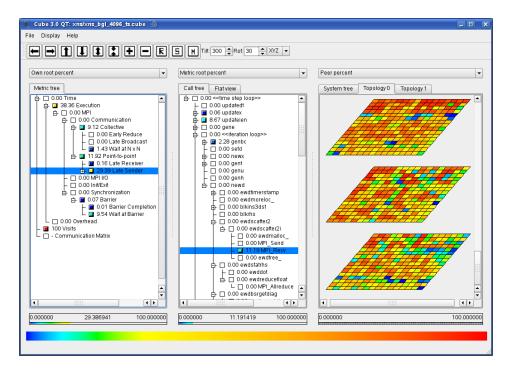
# CUBE - Call-path profile browser

- Browser based on tree widgets & topological display
- Data model and format to store call-path profiles
- Utilities to manipulate & analyze instances
  - Difference, mean, merge, cut, rank
- New version based on Qt
- Current applications
  - Scalasca trace analysis results & runtime summaries
  - TAU call-path profiles
  - MARMOT runtime errors

# Improvements of new version

- More configuration options
  - Order of trees
  - Color spectrum
  - Format and precision of numbers
  - Fonts
- Optimized to handle large data sets
  - Fast parser
  - No 3rd-party XML library
  - Dynamic loading of individual metrics
  - Faster aggregation algorithms
- More flexible and user-friendly topology widget
  - E.g., rotation of topology via mouse

# Ongoing and planned collaborations

- Vampir & Scalasca
  - Unified parallel read interface for OTF & EPILOG traces
  - Unified tracing library (planned)
- TAU & Scalasca
  - Unified instrumentation facilities
  - Unified profiling runtime (planned)

# Thank you!



## scalasca

[www.scalasca.org](http://www.scalasca.org)

Please download and try
Version 1.0