

Portable Parallel I/O

SIONlib

May 26, 2014 | Wolfgang Frings, Florian Janetzko, Michael Stephan

Outline

Introduction

Motivation

SIONlib in a Nutshell

SIONlib file format

Details

Interface

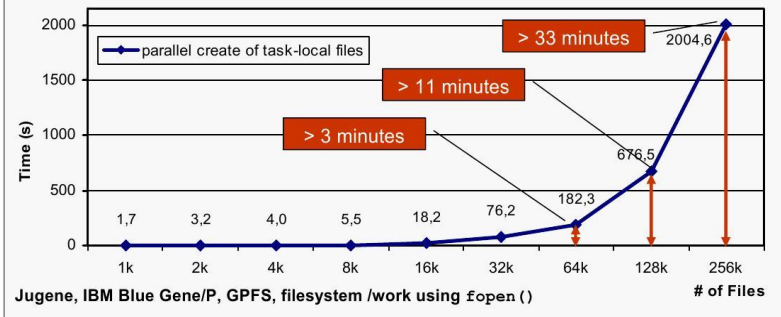
Example

Tools

Exercises

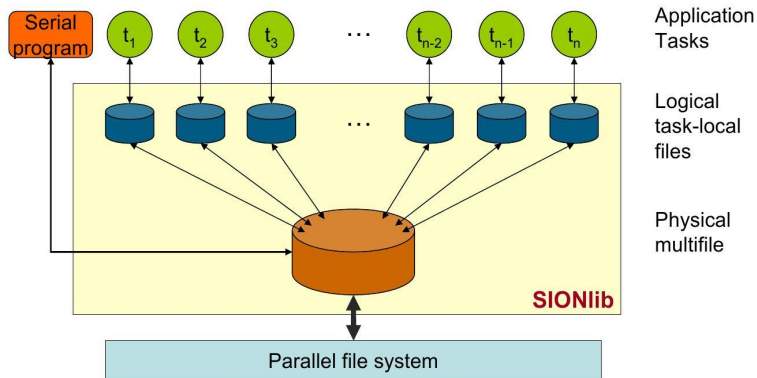
Motivation: Limitations of Task-Local I/O

Example: Creating files in parallel in the same directory



- Contention at the metadata server
- May degrade system I/O performance also for other users
- complicated file handling (e.g. archive)

Motivation: Using Shared Files



- Idea: Mapping many logical files onto one or a few physical file(s)
- → Task-local view to local data not changed

Introduction to SIONlib

- **SIONlib**: Scalable I/O library for parallel access to task-local files
- Collective I/O to binary shared files
- Logical task-local view to data
- Write and read of binary stream-data
- Metadata header/footer included in file
- Collective open/close, independent read/write
- Read/write: POSIX or ANSI-C calls
- Support for MPI, OpenMP, MPI+OpenMP
- C, C++, and Fortran-wrapper
- Optimized for large processor numbers
(e.g. 1.8M tasks on Blue Gene/Q JUQUEEN)

Parallel I/O for Large Scale Application, Types

- External Formats:
 - Exchange data with others → portability
 - Pre- and post-processing on other systems (workflow)
 - Store data without system-dependent structure (e.g. number of tasks)
 - Archive data (long-term readable and self-describing formats)
- Internal Formats:
 - Scratch files, Restart files
 - Fastest I/O preferred
 - Portability and flexibility criteria of second order
 - Read and write data “as-is” (memory dump)
- SIONlib could support I/O of internal formats

SIONlib in a Nutshell: Task-Local I/O

```
/* Open */  
sprintf(tmpfn, "%s.%06d", filename, my_nr);  
fileptr = fopen(tmpfn, "bw", ...);  
...  
/* Write */  
fwrite(bindata, 1, nbytes, fileptr);  
...  
/* Close */  
fclose(fileptr);
```

- Original ANSI C version
- No collective operation, no shared files
- Data: stream of bytes

SIONlib in a Nutshell: Add SIONlib

```
/* Collective Open */
nfiles = 1; chunksize = nbytes;
sid = sion_paropen_mpi(filename, "bw", &nfiles, &chunksize,
                      MPI_COMM_WORLD, &lcomm, &fileptr, ...);
...
/* Write */
fwrite(bindata, 1, nbytes, fileptr);
...
/* Collective Close */
sion_parclose_mpi(sid);
```

- Collective (SIONlib) open and close
- Ready to run ...
- Parallel I/O to one shared file

SIONlib in a Nutshell: Variable Data Size

```
/* Collective Open */
nfiles = 1; chunksize = nbytes;
sid = sion_paropen_mpi(filename, "bw", &nfiles, &chunksize,
                      MPI_COMM_WORLD, &lcomm, &fileptr, ...);
...
/* Write */
if(sion_ensure_free_space(sid, nbytes)) {
    fwrite(bindata, 1, nbytes, fileptr);
}
...
/* Collective Close */
sion_parclose_mpi(sid);
```

- Writing more data than defined in open call
- SIONlib moves forward to next chunk, if data too large for current block

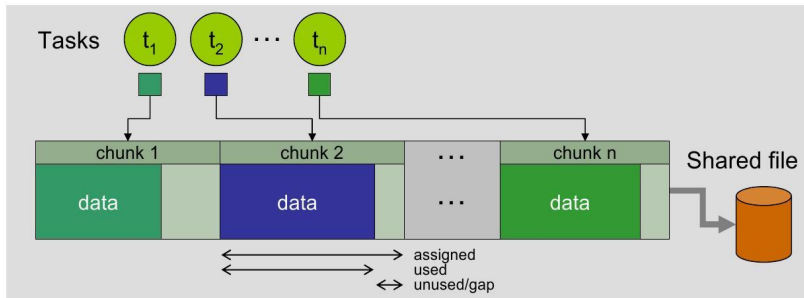
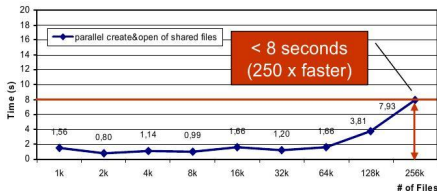
SIONlib in a Nutshell: Wrapper Function

```
/* Collective Open */
nfiles = 1; chunksize = nbytes;
sid = sion_paropen_mpi(filename, "bw", &nfiles, &chunksize,
                      MPI_COMM_WORLD, &lcomm, &fileptr, ...);
...
/* Write */
sion_fwrite(bindata, 1, nbytes, sid);
...
/* Collective Close */
sion_parclose_mpi(sid);
```

- Includes check for space in current chunk
- Parameter of fwrite: fileptr → sid

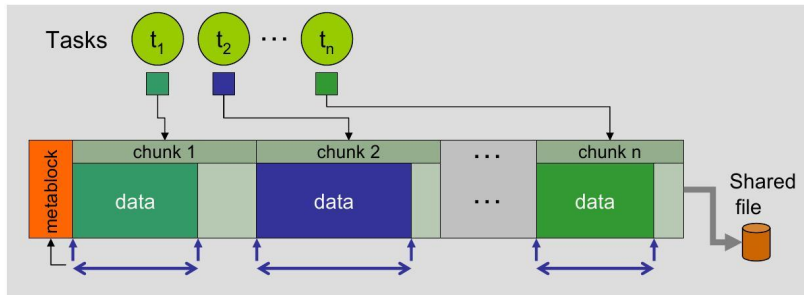
File Format (1): a Single Shared File

- create and open fast
- simplified file handling
- logical partitioning required



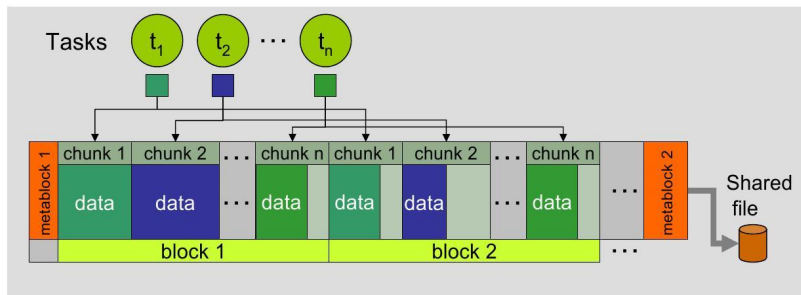
File Format (2): Metadata

- Offset and data size per task
- Tasks have to specify chunk size in advance
- Data must not exceed chunk size



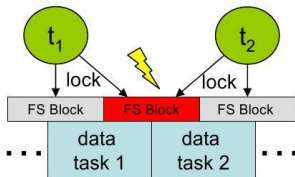
File Format (3): Multiple Blocks of Chunks

- Enhancement: define blocks of chunks
- Metadata now with variable length ($\#task * \#blocks$)
- Second metadata block at the end
- Data of one block does not exceed chunk size



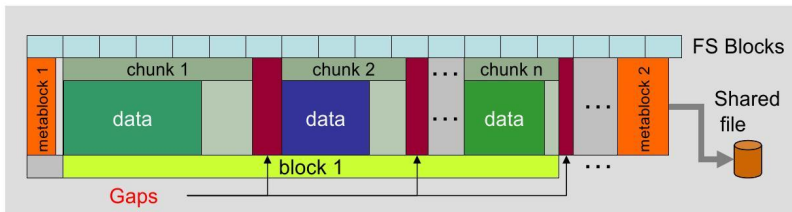
File Format (4): Alignment to Block Boundaries

- Contention: writing to same file-system block in parallel



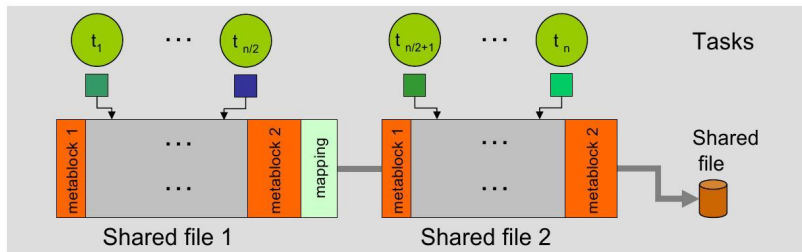
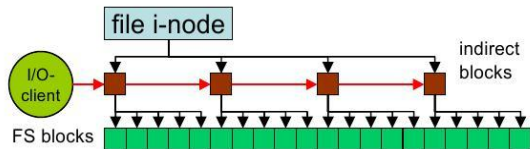
#tasks	data size	blksize	write bandwidth
32768	256 GB	aligned	3650.2 MB/s
32768	256 GB	not aligned	1863.8 MB/s

Jugene (JSC, IBM Blue Gene/P, GPFS, fs:work)



File Format (5): Multiple Physical Files

- Variable number of underlying physical files
- Bandwidth degradation of GPFS by using single shared file



Version, Download, Installation

- Version: 1.5
- Version: file format: 5
- Open-source license, registration
<http://www.juelich.de/jsc/sionlib>
- Installation
`configure; make; make test; make install`

Available Modules

- Modules on Juqueen:

```
juqueen> module avail
----- /usr/local/modulefiles/IO -----
 hdf5/1.8.11_serial(default) sionlib/1.3.6
 hdf5/1.8.9_serial           sionlib/1.3.7
 nco/4.2.5                   sionlib/1.4.3
```

- Modules on Juropa:

```
juropa> module avail
----- /usr/local/modulefiles/IO -----
 hdf5/1.8.5_v16              sionlib/1.2.2(default)
 hdf5/1.8.9                  sionlib/1.3.4
 hdf5/1.8.9_serial           sionlib/1.3.7
 mxml/2.6(default)           sionlib/1.3.7gnu
```

Compiling and Linking own Application

- Include file: `#include "sion.h"`
- The installation of sionlib builds (at least) two libraries:
 - `libsionxxx.a`: the parallel libraries currently supporting MPI
 - `libsionserxxx.a`: serial version of sionlib containing all functions for the serial API of sionlib
 - `xxx` could be an extensions for precision ('_32', '_64') cross compiling ('fe') or compiler ('gcc').
- Script: `sionconfig`: prints for each combination of option correct option for compiling (`-cflags`) or linking (`-libs`):

```
usage: sionconfig [--be] [--fe] [--32|--64] [--gcc] [--for]
               [--ser|--mpi] (--cflags|--libs|--path)
```

- Example: (Makefile)

```
LDFLAGS += '../..bin/sionconfig --libs --mpi -be'
CFLAGS += '../..bin/sionconfig --cflags --mpi -be'
```

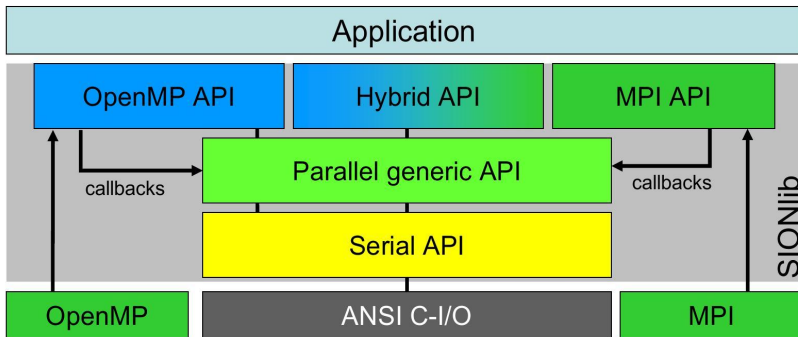
System I/O-Interfaces Used by SIONlib

- Under Unix/Linux available: ANSI-C and POSIX
- POSIX interface
 - `open()`, `close()`, `read()`, `write()`
 - Unbuffered, direct access to file
 - File descriptor: Integer
- ANSI-C
 - `fopen()`, `fclose()`, `fread()`, `fwrite()`
 - Open files and associate a stream with it
 - Typically memory buffer of file system block size
 - Buffer small consecutive reads and writes
 - File pointer: `FILE *`
- Fortran Interface: unformatted I/O
 - Uses typically internally POSIX (or ANSI-C)
 - Files opened in C cannot directly accessed from Fortran (mix languages)

SIONlib Datatypes

- Only used for parameters of SION function calls
- Data written to or read from file is a byte stream and does not need to be declared by special data types
- `sion_int32`
 - 4-byte signed integer (C)
 - INTEGER*4 (Fortran)
- `sion_int64`
 - 8-byte signed integer (C)
 - INTEGER*8 (Fortran)
 - Typically used for all parameters which could be used to compute file positions

SIONlib: Architecture



Outline

Introduction

Interface

- General Parameters
- Open/Close (Parallel)
- Open/Close (Serial)
- Read/Write
- Get Information
- Seek, Utility Functions

Example

Tools

SIONlib API Overview: Open, Close

- Parallel interface, using MPI
`sion_paropen_mpi`, `sion_parclose_mpi`
- Parallel interface, using OpenMP
`sion_paropen_omp`, `sion_parclose_omp`
- Parallel interface, using MPI+OpenMP
`sion_paropen_omp`, `sion_parclose_omp`
- Serial interface
`sion_open`, `sion_open_rank`
`sion_close`

SIONlib API Overview: Read, Write

■ Read Data

- `sion_fread` (SION, internal check of EOF)
- `fread()` (ANSI-C)
- `read()` (POSIX)
- `sion_feof` (check EOF in chunk)

■ Write Data

- `sion_fwrite` (SION, internal checks, e.g. chunk size)
- `fwrite()` (ANSI-C)
- `write()` (POSIX)
- `sion_flush` (flushes data, updates internal metadata)
- `sion_ensure_free_space` (check space in chunk)

SIONlib API Overview: Key Value

- Read Data
 - `sion_fread_key` (read value for specific key)
- Write Data
 - `sion_fwrite_key` (write value for specific key)
- Iterate
 - `sion_fread_key_iterator_next` (get next key)
 - `sion_fread_key_iterator_reset` (reset iterator)

SIONlib API Overview: Get Information I

- Get file pointer for task
 - `sion_get_fp` (ANSI-C)
 - `sion_get_fd` (POSIX)
- Byte order (big(1) or little(0) endian)
 - `sion_get_file_endianness` (endianness of File)
 - `sion_get_endianness` (endianness of current system)
- File state
 - `sion_get_bytes_written` (total number for task written)
 - `sion_get_bytes_read` (total number for task read)
 - `sion_bytes_avail_in_chunk` (rest in chunk)
 - `sion_get_position` (position in file)

SIONlib API Overview: Get Information II

- Multiple physical files
 - `sion_get_mapping` (mapping of global task to file and local task, can be used only on task 0 in parallel-mode)
 - `sion_get_number_of_files` (total number of files)
 - `sion_get_filename` (number of file for this task)
- Serial mode: get information about all tasks
 - `sion_get_locations` (returns pointer to internal chunk description arrays)
 - `sion_is_serial_opened` (indicator for open mode)
- Version
 - `sion_get_version` (returns version of library and fileformat)

SIONlib API Overview: Seek, Utility functions

- Change position in SION file
 - `sion_seek` (parallel mode)
 - `sion_seek_fp` (serial mode, change of file pointer possible)
- Utilities
 - `sion_swap` (change endianness of data in memory)
 - `_sion_file_stat_file` (wrapper for `stat()`, large file support)
- Experimental
 - `sion_coll_fread_mpi` (collective read)
 - `sion_coll_fwrite_mpi` (collective write)

SIONlib Parameters of Open Calls I

- `fname` (file name)
 - Character string describing path and file name
 - Will not be extended by SION-specific suffix
 - Multiple physical files are generated
 - First file: `file_name`
 - All other files: `file_name` + "." + 6-digit-number (000001 ...)
 - All commands and function calls use the base name
- `file_mode`
 - Must specify at least one of the following
 - `r,rb,br` – (read block), open existing SION file for reading
 - `w,wb,bw` – (write block), create a new SION file, open for write; overwrite if existing
 - `posix` – use internally POSIX interface for file access, otherwise ANSI-C
 - Multiple parameter: comma-separated

SIONlib Parameters of Open Calls II

- `sid`
 - `sid = sion_paropen...()` (C)
 - `FSION_PAROPEN_MPI(... , sid)` (Fortran)
 - Unique integer value, referring internally to data structure associated to SION file (internal file handle)
 - Allows multiple simultaneous opened files
 - C: return code, Fortran: last parameter of open call
 - Integer file handle for Fortran necessary
- `chunksize`
 - Pointer of type `sion_int64*` (C)
 - Size of data in bytes written by this tasks
 - May be different for each task
 - Must be set if open for writing
 - Will increased internally to the next multiple of the file system size block size

SIONlib Parameters of Open Calls III

- `fsblksize`
 - Size of file system block in bytes
 - Read-mode: file system block size at write time
 - Write-mode: will be detected by SIONlib if set to -1
- `newfname`
 - File name of physical file assigned to this task

Collective Open (MPI)

C/C++

```
int sion_paropen_mpi (char *fname, const char *file_mode,
                    int      *numFiles,
                    MPI_Comm gComm, MPI_Comm *lComm,
                    sion_int64 *chunksize,
                    sion_int32 *fsblksize,
                    int      *globalrank,
                    FILE      **fileptr,
                    char      **newfname);
```

- Open a SION file in parallel for reading or writing data
- Collective call, called by each task at the same time
- Accesses one or more physical files of a logical SION file
- Parameters are passed “by reference” to pass back information in read open mode

Collective Open (MPI)

Fortran

```
FSION_PAROPEN_MPI (FNAME, FILE_MODE, NUMFILES,  
                  GCOMM, LCOMM, CHUNKSIZE, FSBLKSIZE,  
                  GLOBALRANK, NEWFNAME, SID)  
CHARACTER*(*) FNAME, FILE_MODE, NEWFNAME  
INTEGER      NUMFILES, FSBLKSIZE, GLOBALRANK, SID  
INTEGER      GCOMM, LCOMM  
INTEGER*8    CHUNKSIZE
```

- Open a SION file in parallel for reading or writing data
- Collective call, called by each task at the same time
- Accesses one or more physical files of a logical SION file
- Parameters are passed “by reference” to pass back information in read open mode

Special Parameters of *sion_paropen_mpi*

- Communicator `gComm`
 - Call is collective over all tasks of this communicator
 - Each task gets assigned one chunk of SION file
 - Read: number of tasks must be equivalent to number tasks written to SION file
- Number of physical files
 - `numfiles`, or -1 if specified by communicator
 - `lComm`, or `MPI_Comm_null`
 - Read-mode: parameters will be set by open call
- `globalrank`
 - Rank of task in global communicator `gComm`

Collective Close (MPI)

C/C++

```
int sion_parclose_mpi (int sid);
```

- Closes a SION file in parallel on all tasks
- Collective call, called by each task of (gComm) at the same time
- Metadata will be collected from each tasks
- Metadata blocks of SION file will be written in this call
- Currently no fault tolerant handling of the metadata, call has to executed

Collective Close (MPI)

Fortran

```
FSION_PARCLOSE_MPI (SID, IERR)  
INTEGER SID, IERR
```

- Closes a SION file in parallel on all tasks
- Collective call, called by each task of (gComm) at the same time
- Metadata will be collected from each tasks
- Metadata blocks of SION file will be written in this call
- Currently no fault tolerant handling of the metadata, call has to executed

Collective Open (OpenMP)

C/C++

```
int sion_paropen_omp(char *fname,
                    const char *file_mode,
                    sion_int64 *chunksize,
                    sion_int32 *fsblksize,
                    int *globalrank,
                    FILE **fileptr,
                    char **newfname);
```

- Open a SION file in parallel for writing or reading data
- Collective call, has to be called inside a parallel region
- SION file consists of only one physical file
- Parameters are passed “by reference” to pass back information in read open mode
- Thread-number: `globalrank`

Collective Open (OpenMP)

Fortran

```
FSION_PAROPEN_OMP (FNAME, FILE_MODE,  
                  CHUNKSIZE, FSBLKSIZE,  
                  GLOBALRANK, NEWFNAME, SID)  
CHARACTER*(*) FNAME, FILE_MODE, NEWFNAME  
INTEGER      FSBLKSIZE, GLOBALRANK, SID  
INTEGER*8    CHUNKSIZE
```

- Open a SION file in parallel for writing or reading data
- Collective call, has to be called inside a parallel region
- SION file consists of only one physical file
- Parameters are passed “by reference” to pass back information in read open mode
- Thread-number: `globalrank`

Collective Close (OpenMP)

C/C++

```
int sion_parclose_omp (int sid);
```

- Closes a SION file in parallel on all threads
- Collective call, has to be called inside a parallel region
- Metadata will be collected from each thread
- Metadata blocks of SION file will be written in this call
- Currently no fault tolerant handling of the metadata, call has to be executed

Collective Close (OpenMP)

Fortran

```
FSION_PARCLOSE_OMP (SID, IERR)  
INTEGER            SID, IERR
```

- Closes a SION file in parallel on all threads
- Collective call, has to be called inside a parallel region
- Metadata will be collected from each thread
- Metadata blocks of SION file will be written in this call
- Currently no fault tolerant handling of the metadata, call has to be executed

Collective Open (MPI+OpenMP)

C/C++

```
int sion_paropen_ompi (char *fname, const char *file_mode,
                      int          *numFiles,
                      MPI_Comm gComm, MPI_Comm *lComm,
                      sion_int64  *chunksize,
                      sion_int32  *fsblksize,
                      int          *globalrank,
                      FILE         **fileptr,
                      char         **newfname);
```

- Open a SION file in parallel for reading or writing data
- Collective call, call from each task/thread at the same time, has to be called inside a parallel region
- Parameter and further description see MPI and OpenMP functions

Collective Open (MPI+OpenMP)

Fortran

```
FSION_PAROPEN_OMPI (FNAME, FILE_MODE, NUMFILES,  
                   GCOMM, LCOMM, CHUNKSIZE, FSBLKSIZE,  
                   GLOBALRANK, NEWFNAME, SID)  
  
CHARACTER*(*) FNAME, FILE_MODE, NEWFNAME  
INTEGER      NUMFILES, FSBLKSIZE, GLOBALRANK, SID  
INTEGER      GCOMM, LCOMM  
INTEGER*8    CHUNKSIZE
```

- Open a SION file in parallel for reading or writing data
- Collective call, call from each task/thread at the same time, has to be called inside a parallel region
- Parameter and further description see MPI and OpenMP functions

Collective Close (MPI+OpenMP)

C/C++

```
int sion_parclose_ompi (int sid);
```

- Closes a SION file in parallel on all tasks/threads
- Collective call, call from each task of (gComm) at the same time, call has to be called inside a parallel region
- Metadata will be collected from each task
- Metadata blocks of SION file will be written in this call
- Currently no fault tolerant handling of the metadata, call has to be executed

Collective Close (MPI+OpenMP)

Fortran

```
FSION_PARCLOSE_OMPI (SID, IERR)  
INTEGER             SID, IERR
```

- Closes a SION file in parallel on all tasks/threads
- Collective call, call from each task of (gComm) at the same time, call has to be called inside a parallel region
- Metadata will be collected from each task
- Metadata blocks of SION file will be written in this call
- Currently no fault tolerant handling of the metadata, call has to be executed

Serial Open

C/C++

```
int sion_open(char *fname,  
              const char *file_mode,  
              int *ntasks, int *nfiles,  
              sion_int64 **chunksizes,  
              sion_int32 *fsblksize,  
              int **globalranks,  
              FILE **fileptr);
```

- Open a SION file in serial mode
- All chunks of all tasks can be selected, via `sion_seek_fp`
- Multiple physical files can be handled
- Designed for serial pre- and post-processing tools
- Reads all metadata of all tasks into memory

Serial Open

Fortran

```
FSION_OPEN (FNAME, FILE_MODE, NTASKS, NUMFILES,  
            CHUNKSIZES, FSBLKSIZE,  
            GLOBALRANKS, SID)  
CHARACTER*(*) FNAME, FILE_MODE  
INTEGER      NUMFILES, NTASKS, FSBLKSIZE, SID  
INTEGER      GLOBALRANKS(NTASKS)  
INTEGER*8    CHUNKSIZES(NTASKS)
```

- Open a SION file in serial mode
- All chunks of all tasks can be selected, via `sion_seek_fp`
- Multiple physical files can be handled
- Designed for serial pre- and post-processing tools
- Reads all metadata of all tasks into memory

Serial Open for one Rank

C/C++

```
int sion_open_rank (char      *fname,  
                   const char *file_mode,  
                   sion_int64 *chunksize,  
                   sion_int32 *fsblksize,  
                   int        *rank,  
                   FILE        **fileptr);
```

- Open SION file for one rank in serial mode
- Multiple physical files can be handled
- Designed for parallel program if collective open/close is not possible
- Reads only metadata of this task into memory

Serial Open for one Rank

Fortran

```
FSION_OPEN_RANK (FNAME, FILE_MODE,  
                 CHUNKSIZE, FSBLKSIZE,  
                 RANK, SID)  
CHARACTER*(*) FNAME, FILE_MODE  
INTEGER       FSBLKSIZE, SID, RANK  
INTEGER*8     CHUNKSIZE
```

- Open SION file for one rank in serial mode
- Multiple physical files can be handled
- Designed for parallel program if collective open/close is not possible
- Reads only metadata of this task into memory

Serial Close

C/C++

```
int sion_close (int sid);
```

- Closes a SION file in serial mode
- Metadata blocks of SION file will be written in this call
- Currently no fault tolerant handling of the metadata, call has to be executed

Serial Close

Fortran

```
FSION_CLOSE (SID, IERR)  
INTEGER      SID, IERR
```

- Closes a SION file in serial mode
- Metadata blocks of SION file will be written in this call
- Currently no fault tolerant handling of the metadata, call has to executed

Read Data

C/C++

```
size_t sion_fread (void *ptr,  
                  size_t size,  
                  size_t nmemb,  
                  int sid);
```

- Read $\text{size} * \text{nmemb}$ bytes from current position in chunk
- Internally this function reads in a while loop until all data is read from file. Reading more data than stored in one chunk is possible with this wrapper.
- Returns number of bytes read

Read Data

Fortran

```
FSION_READ (DATA, SIZE, NMEMB, SID, IERR)  
INTEGER    SIZE, NMEMB, SID, IERR
```

- Read $\text{size} * \text{nmemb}$ bytes from current position in chunk
- Internally this function reads in a while loop until all data is read from file. Reading more data than stored in one chunk is possible with this wrapper.
- Returns number of bytes read

End of File

C/C++

```
int sion_feof (int sid);
```

- Equivalent to POSIX feof which cannot be used for shared SION files
- Internally this function flushes all buffer and checks current positions against chunk boundaries
- Moves file pointer to next chunk if end of current chunk is reached
- The function is a task-local function, which can be called independently from other MPI tasks.
- Returns 1 if pointer is behind last byte of data for this task

End of File

Fortran

```
FSION_FEOF (SID, EOF)  
INTEGER      SID, IERR
```

- Equivalent to POSIX feof which cannot be used for shared SION files
- Internally this function flushes all buffer and checks current positions against chunk boundaries
- Moves file pointer to next chunk if end of current chunk is reached
- The function is a task-local function, which can be called independently from other MPI tasks.
- Returns 1 if pointer is behind last byte of data for this task

Write Data

C/C++

```
size_t sion_fwrite (const void *data,  
                  size_t      size,  
                  size_t      nitems,  
                  int          sid);
```

- Write `size*nmemb` bytes to chunk, beginning from current position
- Internally this function uses `sion_ensure_free_space` to check whether there is enough space is available
- Returns number of bytes written

Write Data

Fortran

```
FSION_WRITE (DATA, SIZE, NMEMB, SID, IERR)  
INTEGER     SIZE, NMEMB, SID, IERR
```

- Write $\text{size} * \text{nmemb}$ bytes to chunk, beginning from current position
- Internally this function uses `sion_ensure_free_space` to check whether there is enough space is available
- Returns number of bytes written

Flush Data

C/C++

```
int sion_flush (int sid);
```

- After writing of data this function updates internal data structures to new file position
- To obtain new file position a POSIX flush will be used which could be time consuming

Flush Data

Fortran

```
FSION_FLUSH (SID, IERR)  
INTEGER     SID, IERR
```

- After writing of data this function updates internal data structures to new file position
- To obtain new file position a POSIX flush will be used which could be time consuming

Ensure Free Space in Chunk

C/C++

```
int sion_ensure_free_space (int sid, sion_int64 bytes);
```

- Ensures that there is enough space available for writing
- A new chunk will be allocated if too little space is left in the current chunk
- The function is a task-local function, which can be called independently from other MPI tasks
- The function moves the filepointer to a new position in some cases and flushes also the local filepointer
- Returns 1 if space could be ensured, there is currently no indicator if a new chunk was allocated

Ensure Free Space in Chunk

Fortran

```
FSION_ENSURE_FREE_SPACE (SID, BYTES, IERR)  
INTEGER      SID, IERR  
INTEGER*8    BYTES
```

- Ensures that there is enough space available for writing
- A new chunk will be allocated if too little space is left in the current chunk
- The function is a task-local function, which can be called independently from other MPI tasks
- The function moves the filepointer to a new position in some cases and flushes also the local filepointer
- Returns 1 if space could be ensured, there is currently no indicator if a new chunk was allocated

Get File Pointer

```
FILE * sion_get_fp (int sid); /* ANSI-C */  
int    sion_get_fd (int sid); /* POSIX */
```

- Returns ANSI-C file pointer (`_fp`) or
- Returns POSIX File descriptor (`_fd`)
- File pointer/descriptor corresponds to physical file containing data of current task
- SION file must be opened with corresponding option
- The POSIX file descriptor can be obtained from a ANSI-C file pointer: `fd = fileno(fileptr)`
- ANSI-C file pointer can be obtained from a POSIX file descriptor: `fileptr = fdopen(fd, "r")`

Get Byte Ordering (Endianness)

```
int sion_get_file_endianness (int sid);  
int sion_get_endianness ();
```

- Return endianness (1 → big endian, 0 → little endian)
- For current file (sid), or
- For current runtime environment
- Bytes have to be reordered if:
`sion_get_file_endianness() != sion_get_endianness()`
- Utility for reordering: see `sion_swap`

Seek: Change File Position

C/C++

```
int sion_seek (int      sid,  
              int      rank,  
              int      chunknr,  
              sion_int64 posinchunk );
```

- Sets the file pointer to a new position
- Seek parameters:
 - `rank`: rank number (0,...), or `SION_CURRENT_RANK`
 - `chunknum`: chunk number (0,...), or `SION_CURRENT_BLK`
 - `posinchunk`: position (0,...), or `SION_CURRENT_POS`
- In parallel write mode seeking is currently not supported
- For serial opened file please use `sion_seek_fp`, because physical file pointer could change.

Seek: Change File Position

Fortran

```
FSION_SEEK (SID, RANK, CHUNKNUM, POSINCHUNK, IERR)  
INTEGER      SID, RANK, CHUNKNUM, IERR  
INTEGER*8    POSINCHUNK
```

- Sets the file pointer to a new position
- Seek parameters:
 - `rank` : rank number (0,...), or `SION_CURRENT_RANK`
 - `chunknum` : chunk number (0,...), or `SION_CURRENT_BLK`
 - `posinchunk` : position (0,...), or `SION_CURRENT_POS`
- In parallel write mode seeking is currently not supported
- For serial opened file please use `sion_seek_fp`, because physical file pointer could change.

Seek: Change File Position + FilePtr

C/C++

```
int sion_seek_fp (int      sid,  
                 int      rank,  
                 int      chunknr,  
                 sion_int64 posinchunk,  
                 FILE      **fileptr );
```

- Sets the file pointer to a new position
- Seek parameters → see `sion_seek`, in addition:
 - `fileptr`: ANSI-C pointer to file, should be used after seeking instead of `fileptr` of `open` call
- No Fortran wrapper, `fileptr` unknown in Fortran

Seek: Change File Position + FilePtr

Fortran

```
FSION_SEEK (SID, RANK, CHUNKNUM, POSINCHUNK, IERR)  
INTEGER    SID, RANK, CHUNKNUM, IERR  
INTEGER*8  POSINCHUNK
```

- Sets the file pointer to a new position
- Seek parameters → see `sion_seek`, in addition:
 - `fileptr`: ANSI-C pointer to file, should be used after seeking instead of `fileptr` of `open` call
- No Fortran wrapper, `fileptr` unknown in Fortran

Utility: Swap bytes

C/C++

```
void sion_swap(void *target, void *source,  
              int size, int n, int aflag);
```

- Perform byte-order swapping for arrays of `n` units of `size` bytes
- Bytes are swapped if and only if `aflag == 0`
- Data will be copied from `source` to `target`
- In-place swapping (`target == source`) is allowed; if `target != source`, the buffers must not overlap
- `aflag` can be initialized as follows:
`sion_get_file_endianness() == sion_get_endianness()`

Utility: Swap bytes

Fortran

```
FSION_SWAP (TARGET, SOURCE, SIZE, N, AFLAG, IERR)  
INTEGER      SIZE, N, AFLAG, IERR
```

- Perform byte-order swapping for arrays of `n` units of `size` bytes
- Bytes are swapped if and only if `aflag == 0`
- Data will copied from `source` to `target`
- In-place swapping (`target == source`) is allowed; if `target != source`, the buffers must not overlap
- `aflag` can be initialized as follow:
`sion_get_file_endianness() == sion_get_endianness()`

Outline

Introduction

Interface

Example

Tools

Exercises

Example Code: sion_par_write (Part 1)

C/C++

```
#include <sion.h>

/* SION parameters */
int      sid, numFiles, globalrank;
MPI_Comm lComm;
sion_int64 chunksize, left, bwrote;
sion_int32 fsblksize;
char     fname[256], *newfname=NULL;
FILE     *fileptr;
/* initialize MPI */
MPI_Init(&argc, &argv);
MPI_Comm_rank(MPI_COMM_WORLD, &my_rank);
MPI_Comm_size(MPI_COMM_WORLD, &num_procs);
/* open parameters */
chunksize = 10;  globalrank = my_rank;
strcpy(fname, "parfile.sion");
numFiles   = 1;  fsblksize   = -1;
```

Example Code: sion_par_write (Part 1)

Fortran

```
! SION parameters
integer*8      :: chunksize
integer        :: gComm,lComm,sid,globalrank,ierr
integer        :: fsblksize, nfiles
character(len=255) :: filename = 'test_sionfile.dat'
character(len=255) :: newfname
integer*4,dimension(:),allocatable :: buffer

! MPI initialization
call MPI_Init(ierr)
call MPI_Comm_size(MPI_COMM_WORLD,nranks,ierr)
call MPI_Comm_rank(MPI_COMM_WORLD,my_rank,ierr)

! create a new file
gcomm=MPI_COMM_WORLD
globalrank=my_rank
fsblksize=-1
chunksize=10
nfiles=1
```

Example Code sion_par_write (Part 2)

C/C++

```
/* create a new file */
sid = sion_paropen_mpi(fname, "bw", &numFiles,
                      MPI_COMM_WORLD, &lComm,
                      &chunksize, &fsblksize,
                      &globalrank,
                      &fileptr, &newfname);

/* write buffer to file */
p = (char *) fname;
sion_fwrite(p, 1, chunksize, sid);
/* close file */
sion_parclose_mpi(sid);
/* finalize MPI */
MPI_Finalize();
```


Example Code sion_par_write (Part 2)

Fortran

```
call fsion_paropen_mpi(trim(filename),'bw', &
                      nfiles, gComm, &
                      lComm, chunksize, &
                      fsblksize, globalrank, &
                      newfname, sid)

! write buffer to file
call fsion_write(buffer,4,veclen,sid,ierr)

! close file
call fsion_parclose_mpi(sid,ierr)

! MPI finalization
call MPI_Finalize(ierr)
```

Example: Blue Gene/Q I/O-node Task Mapping I

- Blue Gene/Q CPU-nodes are connected to I/O-nodes (JUQUEEN: 128 CPU-nodes : 1 I/O-Node)
- Good performance: one physical file per I/O-node
- Special MPI communicator containing all tasks connected to the same I/O-Bridge

Example: Blue Gene/Q I/O-node Task Mapping II

```
/* communicator consists of all task
   working with the same I/O-node */
MPI_Comm commSame;
MPIX_Pset_same_comm_create(&commSame);
MPI_Comm_size(commSame, &sizeSame);
MPI_Comm_rank(commSame, &rankSame);

/* create a new file */
sid = sion_paropen_mpi(fname, "bw", &numFiles,
                      MPI_COMM_WORLD, &commSame,
                      &chunksize, &fsblksize,
                      &globalrank, &fileptr, &newfname);
```

Outline

Introduction

Interface

Example

Tools

Managing SION files
Parallel Benchmark

Exercises

Tools: sionsplit

- Generates task-local files from SION-file
- Usage: `sionsplit [options] sionfile prefix`
- Options:

```
prefix          directory and/or filename-prefix for
                 task-local files
[-v]            verbose mode
[-g]            use global rank for numbering files
[-d <num>]     number of digits for filename
                 generation (default 5)
```

Tools: sioncat

- Extract all or selected data from a SION file
- Usage: `sioncat [options] sionfile`
- Options:

```
[-v]                verbose mode
[-t <tasknum>]     write only data of task <tasknum>
[-b <blknum>]      write only data of block <blknum>
[-o <outfile>]     file data will be written to, if
                    not specified stdout will used
```

Tools: siondump

- Print metadata information of a SION file
- Usage: `siondump [options] sionfile`
- Options:

```
[-a]          print all information about all blocks  
[-m]          print all mapping information  
[-l]          print all sizes in number of bytes  
[-v]          verbose mode
```

Tools: siondefrag

- Generates new SION file from existing one
- Changes the underlying file system block size
- Can be used to remove empty gaps, caused by
 - Not completely filled chunks
 - Alignment to file system blocks
- Usage: `siondefrag [options] sionfile new_sionfile`
- Options:

```
[-Q <fsblksize>]  filesystem blocksize for new sion file  
                   in MB           (default from input file)  
[-q <fsblksize>]  filesystem blocksize for new sion file  
                   in bytes
```


Note on Sparse Files

In general, SION files can be *sparse*, i.e. the file contains unused space. This usually is handled efficiently by the underlying file system using sparse files. Transferring files using `rsync` with the `-S` or `--sparse` option preserves the sparsity, without this option the resulting files are dense.

Tools: partest, Parallel I/O benchmark I

- Usage: `partest [options]`
- Options (file settings):

```
[-f filename]      (--filename[=]) filename of direct access file
[-n <numfiles>]   (--numfiles[=]) number of files file
[-r <chunksize>]  (--chunksize[=]) sion chunk size (*)
[-q <fsblksize>]  (--fsblksize[=]) size of filesystem blocks (*)
```

- Options (test configuration):

```
[-T <type>]       (--testtype[=]) testtype (0:SION, collective)
                                   (1:SION, independent read)
                                   (2:MPI IO) (3:Task-Local)
[-b <bufsize>]    (--bufsize[=])  blocksize written by ONE fwrite (*)
[-g <totalsize>]  (--totalsize[=]) global total size of data written(*)
[-s <localsize>] (--localsize[=]) size of local data for each task(*)
[-F <factor>]    (--factor[=])    random factor (0.0 to 1.0, def: 0.0)
[-R (0|1)]       (--read[=])      switch read off/on
[-W (0|1|2)]     (--write[=])     switch write off, on, or 2x write
```

(*) Size Formats: `<d>`[g,G,Gb,GB, m,M,Mb,MB, k,K,Kb,KB, GiB, MiB, KiB]

Tools: partest, Parallel I/O benchmark II

- Options (test specific configuration):

```

[-v]          (--verbose[=](0|1))    verbose print info for each task
[-C]          (--nochecksum[=](0|1))  suppress checksum
[-d]          (--debugtask[=](0|1))    debug task 0
[-D]          (--Debugtask[=](0|1))    debug task n
[-L]          (--posix[=](0|1))        use POSIX calls instead of ANSI calls
[-M]          (--collwrite[=](0|1))    use collective write if possible
[-m]          (--collread[=](0|1))     use collective read if possible
[-Z <offset>] (--taskoffset[=])       shift tasks numbering for reading by
                                         offset to ommit data caching of
                                         file-system (0)

[-O <bytes>]  (--byteoffset[=])       start offset, write <bytes> first
                                         before using blksize (0)
[-j <#tasks>] (--serialized[=])       serialize I/O, only I/O of #tasks
                                         are running in parallel
                                         (-1 -> all tasks in parallel,
                                         -2 -> use transactions, def: -1)

[-X]          (--unlinkfiles[=](0|1))  remove files after test
  
```

Tools: partest, Parallel I/O benchmark III

- Options (Blue Gene/L, Blue Gene/P, Blue Gene/Q):

```
[-P]                (--bgionode[=](0|1))  order tasks by BG I/O-node  
[-p <numtasks>]   (--bgtaskperionode[=])  number of tasks  
                                                           per BG I/O-node (default: all)
```

- Options (MPI-I/O Hints):

```
[-w]                (--hintlargeblock[=](0|1))  IBM, Large Block IO  
[-Q <size>]        (--hintiobufsize[=])        IBM, IO bufsize in KB  
[-x]                (--hintsparseaccess[=](0|1))  IBM, sparse access
```

Outline

Introduction

Interface

Example

Tools

Exercises

Parallel Write/Read

Exercise: Parallel Write/Read I

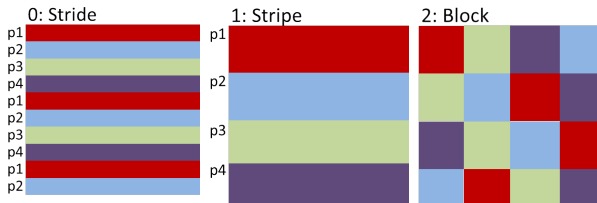
- Create two parallel applications (C, Fortran):
 - 1 Write
 - Creating a SION data set, where each task writes local data to a the corresponding chunk of the SION-file
 - A local vector of 10000 integers should be allocated and initialized with the task number
 - Each task should write the vector to the SION data set
 - 2 Read
 - Read the data of the corresponding chunk into memory
 - and check whether the data is consistent (task number)
- Run `siondump` on the SION-file to check the metadata
- Use `siondefrag` to create a dense version of the SION file, and check again the metadata of the new file
- Extract the chunks of the SION file into task-local files

Exercise: Parallel Write/Read II

- Check whether data is written with same endianness, swap if necessary

Exercise: Mandelbrot I

- Complete the code examples `mandel_sion_to_ppm_skel.c` and `output_sion_skel.c` with the missing SIONlib calls
 - The logic of the different calculation methods can be visualized in the following figures



Exercise: Mandelbrot II

- The structure of the data is shown in the following figure

File Format

