

FTI - Usermanual

Author: Kai Keller

November, 2016

Contents

| | | |
|----------|---|-----------|
| 1 | Introduction | 3 |
| 2 | Multilevel Checkpointing | 4 |
| 2.1 | L1 | 4 |
| 2.2 | L2 | 4 |
| 2.3 | L3 | 4 |
| 2.4 | L4 | 4 |
| 3 | API Reference | 5 |
| 3.1 | FTI Datatypes and FTI Constants | 5 |
| 3.1.1 | FTI Datatypes | 5 |
| 3.1.2 | FTI Constants | 5 |
| 3.2 | FTI_Init() | 6 |
| 3.3 | FTI_InitType() | 7 |
| 3.4 | FTI_Protect(...) | 8 |
| 3.5 | FTI_Checkpoint(...) | 9 |
| 3.6 | FTI_Recover() | 10 |
| 3.7 | FTI_Status() | 11 |
| 3.8 | FTI_Snapshot() | 12 |
| 3.9 | FTI_Abort() | 13 |
| 3.10 | FTI_Finalize() | 14 |
| 4 | Configuration | 15 |
| 4.1 | [Basic] | 15 |
| 4.1.1 | Head | 15 |
| 4.1.2 | Node_size | 15 |
| 4.1.3 | Ckpt_dir | 15 |
| 4.1.4 | Glbl_dir | 15 |
| 4.1.5 | Meta_dir | 16 |
| 4.1.6 | Ckpt_L1 | 16 |

| | | |
|--------|----------------|----|
| 4.1.7 | Ckpt_L2 | 16 |
| 4.1.8 | Ckpt_L3 | 16 |
| 4.1.9 | Ckpt_L4 | 16 |
| 4.1.10 | Inline_L2 | 17 |
| 4.1.11 | Inline_L3 | 17 |
| 4.1.12 | Inline_L4 | 17 |
| 4.1.13 | keep_last_ckpt | 17 |
| 4.1.14 | Group_size | 18 |
| 4.1.15 | Verbosity | 18 |
| 4.2 | [Restart] | 18 |
| 4.2.1 | Failure | 18 |
| 4.2.2 | Exec_ID | 19 |
| 4.3 | [Advanced] | 19 |
| 4.3.1 | Block_size | 19 |
| 4.3.2 | Mpi_tag | 19 |
| 4.3.3 | Local_test | 20 |

1 Introduction

In high performance computing (HPC), systems are built from highly reliable components. However, the overall failure rate of supercomputers increases with component count. Nowadays, petascale machines have a mean time between failures (MTBF) measured in hours or days and fault tolerance (FT) is a well-known issue. Long running large applications rely on FT techniques to successfully finish their long executions. Checkpoint/Restart (CR) is a popular technique in which the applications save their state in stable storage, frequently a parallel file system (PFS); upon a failure, the application restarts from the last saved checkpoint. CR is a relatively inexpensive technique in comparison with the process-replication scheme that imposes over 100% of overhead.

However, when a large application is checkpointed, tens of thousands of processes will each write several GBs of data and the total checkpoint size will be in the order of several tens of TBs. Since the I/O bandwidth of supercomputers does not increase at the same speed as computational capabilities, large checkpoints can lead to an I/O bottleneck, which causes up to 25% of overhead in current petascale systems. Post-petascale systems will have a significantly larger number of components and an important amount of memory. This will have an impact on the system's reliability. With a shorter MTBF, those systems may require a higher checkpoint frequency and at the same time they will have significantly larger amounts of data to save. Although the overall failure rate of future post-petascale systems is a common factor to study when designing FT-techniques, another important point to take into account is the pattern of the failures. Indeed, when moving from 90nm to 16nm technology, the soft error rate (SER) is likely to increase significantly, as shown in a recent study from Intel. A recent study by Dong et al. explains how this provides an opportunity for local/global hybrid checkpoint using new technologies such as phase change memories (PCM). Moreover, some hard failures can be tolerated using solid-state-drives (SSD) and cross-node redundancy schemes, such as checkpoint replication or XOR encoding which allows to leverage multi-level checkpointing, as proposed by Moody et al.. Furthermore, Cheng et al. demonstrated that more complex erasure codes such as Reed-Solomon (RS) encoding can be used to further increase the percentage of hard failures tolerated without stressing the PFS.

FTI is a multi-level checkpointing interface. It provides an api which is easy to apply and offers a flexible configuration to enable the user to select the checkpointing strategy which fits best to the problem.

2 Multilevel Checkpointing

2.1 L1

L1 denotes the first safety level in the multilevel checkpointing strategy of FTI. The checkpoint of each process is written on the local SSD of the respective node. This is fast but possesses the drawback, that in case of a data loss and corrupted checkpoint data even in only one node, the execution cannot successfully restarted.

2.2 L2

L2 denotes the second safety level of checkpointing. On initialisation, FTI creates a virtual ring for each group of nodes with user defined size (see [4.1.14](#)). The first step of L2 is just a L1 checkpoint. In the second step, the checkpoints are duplicated and the copies stored on the neighbouring node in the group.

That means, in case of a failure and data loss in the nodes, the execution still can be successfully restarted, as long as the data loss does not happen on two neighbouring nodes at the same time.

2.3 L3

L3 denotes the third safety level of checkpointing. In this level, the checkpoint data trunks from each node getting encoded via the Reed-Solomon (RS) erasure code. The implementation in FTI can tolerate the breakdown and data loss in half of the nodes.

In contrast to the safety level L2, in level L3 it is irrelevant which of nodes encounters the failure. The missing data can get reconstructed from the remaining RS-encoded data files.

2.4 L4

L4 denotes the fourth safety level of checkpointing. All the checkpoint files are flushed to the parallel file system (PFS).

3 API Reference

3.1 FTI Datatypes and FTI Constants

3.1.1 FTI Datatypes

FTI_CHAR : FTI data type for chars
FTI_SHRT : FTI data type for short integers.
FTI_INTG : FTI data type for integers.
FTI_LONG : FTI data type for long integers.
FTI_UCHR : FTI data type for unsigned chars.
FTI_USHT : FTI data type for unsigned short integers.
FTI_UINT : FTI data type for unsigned integers.
FTI_ULNG : FTI data type for unsigned long integers.
FTI_SFLT : FTI data type for single floating point.
FTI_DBLE : FTI data type for double floating point.
FTI_LDBE : FTI data type for long double floating point.

3.1.2 FTI Constants

FTI_BUFS : 256
FTI_DONE : 1
FTI_SCES : 0
FTI_NSCS : -1

3.2 FTI_Init()

- Reads configuration file.
- Creates checkpoint directories.
- Detects topology of the system.
- Regenerates data upon recovery.

Definition

```
int FTI_Init(char* configFile, MPI_Comm globalComm)
```

On entry

```
char* configFile  
MPI_Comm globalComm
```

configfile : Path to the config file

globalcomm : MPI communicator used for the execution

On return

```
FTI_SCES  
FTI_NSCS
```

FTI_SCES : On success

FTI_NSCS : On failure

3.3 FTI_InitType()

– initializes a data type.

Definition

```
int FTI_InitType(FTIT_type* type, int size)
```

On entry

```
FTIT_type* type  
int size
```

type : The data type to be initialized.

size : The size of the data type to be initialized.

On return

```
FTI_SCES
```

FTI_SCES : On success.

3.4 FTI_Protect(...)

- Stores metadata concerning the variable to protect.

Definition

```
int FTI_Protect(int id, void* ptr, long count, \
               FTIT_type type)
```

On entry

```
int id
void* ptr
long count
FTIT_type type
```

id : Unique ID of the variable to protect

ptr : Pointer to memory address of variable

count : Number of elements at memory address

type : FTI datatype of variable to protect

On return

```
FTI_SCES
exit(1)
```

FTI_SCES : On success

exit(1) : Number of protected variables is > FTI_BUFS

3.5 FTI_Checkpoint(...)

- Writes values of protected runtime variables to a checkpoint file of requested level.

Definition

```
int FTI_Checkpoint( int id , int level )
```

On entry

```
int id  
int level
```

id : Unique checkpoint ID

level : Checkpoint level (1=L1, 2=L2, 3=L3, 4=L4)

On return

```
FTI_DONE  
FTI_NSCS
```

FTI_DONE : On success

FTI_NSCS : On failure

3.6 FTI_Recover()

- Loads checkpoint data from the checkpoint file and initializes the runtime variables of the execution.

Definition

```
int FTI_Recover()
```

On return

```
FTI_SCES  
FTI_NSCS
```

FTI_SCES : Checkpoint data is successfully restored.

FTI_NSCS : on failure.

3.7 FTI_Status()

– Returns the current status of the recovery flag.

Definition

```
int FTI_Status()
```

On return

```
int
```

- 0** : No checkpoints taken yet.
- 1** : At least one checkpoint is taken. If execution fails, the next start will be a restart.
- 2** : The execution is a restart from checkpoint level L4 and keep_last_checkpoint was enabled during the last execution.

3.8 FTI_Snapshot()

- Loads checkpoint data and initializes runtime variables upon recovery.
- Writes multilevel checkpoints regarding their requested frequencies.

Definition

```
int FTI_Snapshot()
```

On return

```
FTI_SCES  
FTI_DONE  
exit(1)
```

FTI_SCES : On successful call / recovery.

FTI_DONE : If checkpoint was taken successfully.

exit(1) : If recovery unsuccessful.

3.9 FTI_Abort()

- Aborts the application after cleaning the file system.

Definition

```
int FTI_Abort()
```

On return

```
exit(1)
```

exit(1) : After call.

3.10 FTI_Finalize()

- Frees the allocated memory.
- Communicates the end of the execution to dedicated threads.
- Cleans checkpoints and metadata.

Definition

```
int FTI_Finalize()
```

On return

```
FTI_SCES  
exit(0)
```

FTI_SCES : For an application process.

exit(0) : For a head process.

4 Configuration

4.1 [Basic]

4.1.1 Head

The checkpointing safety levels L2, L3 and L4 produce additional overhead due to the necessary postprocessing work on the checkpoints. FTI offers the possibility to create an mpi process, called *HEAD*, in which this postprocessing will be accomplished. This allows it for the application processes to continue the execution immediately after the checkpointing.

| Value | Meaning |
|-------|--|
| 0 | The checkpoint postprocessing work is covered by the application processes. |
| 1 | The HEAD process accomplishes the checkpoint postprocessing work (notice: In this case, the number of application processes will be $(n-1)/\text{node}$). |

(*default* = 0)

4.1.2 Node_size

Lets FTI know, how many processes will run on each node (*npp*). In most cases this will be the amount of processing units within the node (e.g. 2 CPU's/node and 8 cores/CPU \rightarrow 16 processes/node).

| Value | Meaning |
|---------------|---|
| npp (int > 0) | Number of processing units within each node (notice: The total number of processes must be a multiple of group size \times node size) |

(*default* = 2)

4.1.3 Ckpt_dir

This entry defines the path to the local hard drive on the nodes.

| Value | Meaning |
|--------|---|
| string | Path to the local hard drive on the nodes |

(*default* = `_BLANK_`)

4.1.4 Glbl_dir

This entry defines the path to the checkpoint folder on the PFS (L4 checkpoints).

| Value | Meaning |
|--------|---|
| string | Path to the checkpoint directory on the PFS (notice: The directory has to be created before execution). |

(default = /path/to/global/storage/)

4.1.5 Meta_dir

This entry defines the path to the meta files directory. The directory has to be accessible from each node. It keeps files with informations about the topology of the execution.

| Value | Meaning |
|--------|---|
| string | Path to the meta files directory. (notice: The directory has to be created before execution). |

(default = /home/username/.fti)

4.1.6 Ckpt_L1

Here, the user sets the checkpoint frequency of L1 checkpoints.

| Value | Meaning |
|--------------------------|--|
| L1 freq. (int ≥ 0) | L1 checkpointing frequency in min^{-1} . If the value is equal to 0, L1 checkpointing is disabled. |

(default = 3)

4.1.7 Ckpt_L2

Here, the user sets the checkpoint frequency of L2 checkpoints.

| Value | Meaning |
|--------------------------|--|
| L2 freq. (int ≥ 0) | L2 checkpointing frequency in min^{-1} . If the value is equal to 0, L2 checkpointing is disabled. |

(default = 5)

4.1.8 Ckpt_L3

Here, the user sets the checkpoint frequency of L3 checkpoints.

| Value | Meaning |
|--------------------------|--|
| L3 freq. (int ≥ 0) | L3 checkpointing frequency in min^{-1} . If the value is equal to 0, L3 checkpointing is disabled. |

(default = 7)

4.1.9 Ckpt_L4

Here, the user sets the checkpoint frequency of L4 checkpoints.

| Value | Meaning |
|--------------------------|---|
| L4 freq. (int ≥ 0) | L4 checkpointing frequency in min^{-1} . If the value is equal to 0, L4 checkpointing is disabled. |

(*default* = 11)

4.1.10 Inline_L2

In this entry, the user decides, whether the post-processing work on the L2 checkpoints is done by the HEAD or by the application processes.

| Value | Meaning |
|-------|--|
| 0 | The post-processing work of the L2 checkpoints is done by the HEAD (notice: This setting is only allowed if Head = 1). |
| 1 | The post-processing work of the L2 checkpoints is done by the application processes. |

(*default* = 1)

4.1.11 Inline_L3

In this entry, the user decides, whether the post-processing work on the L3 checkpoints is done by the HEAD or by the application processes.

| Value | Meaning |
|-------|--|
| 0 | The post-processing work of the L3 checkpoints is done by the HEAD (notice: This setting is only allowed if Head = 1). |
| 1 | The post-processing work of the L3 checkpoints is done by the application processes. |

(*default* = 1)

4.1.12 Inline_L4

In this entry, the user decides, whether the post-processing work on the L4 checkpoints is done by the HEAD or by the application processes.

| Value | Meaning |
|-------|--|
| 0 | The post-processing work of the L4 checkpoints is done by the HEAD (notice: This setting is only allowed if Head = 1). |
| 1 | The post-processing work of the L4 checkpoints is done by the application processes. |

(*default* = 1)

4.1.13 keep_last_ckpt

This setting tells FTI whether the last checkpoint taken during the execution will be kept in the case of a successful run or not.

| Value | Meaning |
|-------|--|
| 0 | After <code>FTI_Finalize()</code> , the meta files and checkpoints will be removed. No checkpoint data will be kept on the PFS or on the local hard drives of the nodes. |
| 1 | After <code>FTI_Finalize()</code> , the last checkpoint will be kept and stored on the PFS as a L4 checkpoint (notice: Additionally, the setting <i>failure</i> in the configuration file is set to 2. This will lead to a restart from the last checkpoint if the application is executed again). |

(*default* = 0)

4.1.14 Group_size

The group size entry sets, how many nodes (members) forming a group.

| Value | Meaning |
|------------------------------|---|
| int i ($2 \leq i \leq 32$) | Number of nodes contained in a group (notice: The total number of processes must be a multiple of group size \times node size). |

(*default* = 4)

4.1.15 Verbosity

Sets the verbosity level

| Value | Meaning |
|-------|---|
| 1 | Debug sensitive. Beside warnings, errors and informations, FTI debugging information will be printed. |
| 2 | Information sensitive. FTI prints warnings, errors and informations. |
| 3 | FTI prints only warnings and errors. |
| 4 | FTI prints only errors. |

(*default* = 2)

4.2 [Restart]

4.2.1 Failure

This setting should mainly set by FTI itself. The behaviour within FTI is the following:

- Within `FTI_Init()`, it remains on its initial value.
- After the first checkpoint is taken, it is set to 1.
- After `FTI_Finalize()` and `keep_last_ckpt = 0`, it is set to 0.
- After `FTI_Finalize()` and `keep_last_ckpt = 1`, it is set to 2.

| Value | Meaning |
|-------|--|
| 0 | The application starts with its initial conditions (notice: In order to force a clean start, the value may be set to 0 manually. In this case the user has to take care about removing the checkpoint data from the last execution). |
| 1 | FTI is searching for checkpoints and starts from the highest checkpoint level (notice: If no readable checkpoints are found, the execution stops) |
| 2 | FTI is searching for the last L4 checkpoint and restarts the execution from there (notice: If checkpoint is not L4 <i>or</i> checkpoint is not readable, the execution stops). |

(*default = 0*)

4.2.2 Exec_ID

This setting should mainly set by FTI itself. During `FTI_Init()` the execution ID is set if the application starts for the first time (failure = 0) *or* the execution ID is used by FTI in order to find the checkpoint files for the case of a restart (failure = 1,2).

| Value | Meaning |
|---------------------|--|
| yyyy-mm-dd_hh-mm-ss | Execution ID (notice: If variate checkpoint data is available, the execution ID may set by the user to assign the desired starting point). |

(*default = XXXX-XX-XX_XX-XX-XX*)

4.3 [Advanced]

The settings in this section, should *ONLY* be changed by advanced users.

4.3.1 Block_size

FTI temporarily copies small blocks of the L2 and L3 checkpoints to send them through MPI. The size of the data blocks can be set here.

| Value | Meaning |
|-------|---|
| int | Size in KB of the data blocks send by FTI through MPI for the checkpoint levels L2 and L3 |

(*default = 1024*)

4.3.2 Mpi_tag

FTI uses a certain tag for the MPI messages. This tag can be set here.

| Value | Meaning |
|-------|--|
| int | Tag, used for MPI messages within FTI. |

(*default = 2612*)

4.3.3 Local_test

FTI is building the topology of the execution, by determining the hostnames of the nodes on which each process runs. Depending on the settings for Group_size, Node_size and Head, FTI assigns each particular process to a group and decides which process will be Head or Application dedicated. This is meant to be a local test. In certain situations (e.g. to run FTI on a local machine) it is necessary to disable this function.

| Value | Meaning |
|-------|---|
| 0 | Local test is disabled. FTI will simulate the situation set in the configuration. |
| 1 | Local test is enabled (notice: FTI will check if the settings are correct on initialization and if necessary stop the execution). |

(*default = 1*)