

論文の内容の要旨

REPLICATION MECHANISMS OF PROCESS AND FILE SYSTEM'S METADATA FOR FAULT TOLERANCE

(耐故障のためのプロセスおよびファイルシステムのメタデータ
の複製機構)

氏名 廖 劍偉

Software-based replication is commonly employed in critical computing systems to achieve higher availability and reliability. With multiple redundant replications, when one of them fails, the other replicas are still able to work without any service interrupts. Correctness criteria for replication-based mechanisms are replica consistency and data integrity (also called constraint consistency). One particularly difficult challenge is to ensure the correctness criteria with minimal overhead. That is because the original is supposed to be hanged during replication to ensure the replica has an integral and consistent state. Moreover, maintaining replica consistency brings about the extra synchronization overhead, as well. The main focus of this thesis is to demonstrate that it is possible to build a fault-tolerant system through software-based replication mechanisms with slight extra overhead. Aiming at critical processes and the crucial metadata of parallel file systems, two replication mechanisms have been proposed, implemented and evaluated separately in this thesis:

One replication mechanism is a checkpoint approach called TIC-CKPT that replicates the running state of real-time and interactive processes to nonvolatile storage. Unlike conventional checkpoint mechanisms, which hang the checkpointed process while replicating its running state, the newly proposed mechanism allows the checkpointed process to continue running without stopping while checkpoints are set to a large extent. Through tracing TLB misses, it

blocks the accesses to target memory pages first time while dumping memory address space (the most time-consuming step when setting a checkpoint). At that time, a kernel thread, called checkpointer, copies the memory access target pages to the designated memory buffer, and then resumes the memory accesses. Finally, the memory pages copied in the designated buffer, will be used to construct an integral and consistent running state of the checkpointed process. Moreover, compared with the traditional concurrent checkpoint mechanism, since TIC-CKPT does not operate on the page table of the checkpointed process, it can reduce downtime brought by setting checkpoints much more.

Another replication mechanism that replicates part of files' metadata to the storage servers to achieve high availability and reliability metadata service for parallel file systems, which has active/standby configured metadata servers (MDSs). Since partial replication of metadata is adopted, only a small part of metadata operations trigger metadata synchronization between a metadata server and storage servers. This newly proposed mechanism has been applied in a prototype parallel file system called PARTE, which replicates and distributes a part of a file's metadata to the header sections of the corresponding stripes on the storage servers (OSTs), while the client file system (client) keeps certain sent metadata requests. If the active MDS has crashed for some reason, these backup requests will be replayed by the standby MDS to restore the lost metadata. In case one or more backup requests are lost due to network problems or dead clients, the latest metadata saved in the stripes will be used to construct consistent and up-to-date metadata on the standby MDS. Moreover, the clients and OSTs can work in both normal mode and recovery mode in the PARTE file system. This differs from conventional active/standby configured MDSs file systems, which hang all I/O requests and metadata requests during restoration of the lost metadata. In PARTE, previously connected clients can continue to perform I/O operations and relevant metadata operations, because OSTs work as temporary MDSs during that period.