

## 高信頼化ミドルウェア ARTEMIS の概要と チェックポイント生成方式

6 C-6

### (Advanced Reliable disTributed Environment MiddleWare System)

(<http://www2.toshiba.co.jp/ilab/artemis>)

白木原 敏雄<sup>1</sup>      平山 秀昭<sup>2</sup>      佐藤 記代子<sup>1</sup>      金井 達徳<sup>1</sup>

<sup>1</sup>(株) 東芝研究開発センター 情報・通信システム研究所

<sup>2</sup>(株) 東芝情報・通信システム技術研究所

## 1 はじめに

近年、イントラネット上でのデータベース (DB) 処理をクライアント・サーバシステム上で行う等、分散処理が広く使用されるようになってきている。このような環境では、サーバ計算機に障害が発生した場合そのサーバ計算機が行っているサービスが使用できなくなるため、信頼性が重要になってくる。高信頼化ミドルウェア ARTEMIS(Advanced Reliable disTributed Environment MiddleWare System) は分散システム全体の高信頼化を実現するものであり [1, 2]、基本メカニズムとして、チェックポイント・リスタート方式をベースにしている。本稿では、ARTEMIS の概要およびプロセスのチェックポイント (CP) 生成方法について述べる。

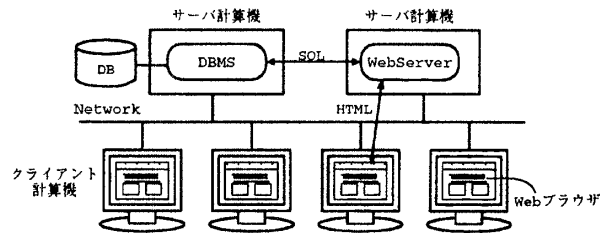


図 1: 3層クライアント・サーバシステムの例

なハードウェア故障の場合、修理が完了するまで、最悪数日、サービスが停止してしまうことになる。

## 2 ARTEMIS の概要

ARTEMISはイントラネット上の分散処理、特にデータベースサーバ、Webサーバ、Webブラウザの3層で構成される3層クライアント・サーバシステムやグループウェアをターゲットにシステム全体の高信頼性を提供するミドルウェアである。

図1は、3層クライアント・サーバシステムの構成例を示したものである。ユーザはクライアント計算機上のWebブラウザからWebサーバにDBアクセス要求を行う。Webサーバはデータベースシステムにアクセスし、結果をHTMLに変換してWebブラウザに返す。このような環境において処理途中にDBサーバのサーバ計算機に障害が発生した場合、それ以降の処理が継続できなくなる。すなわち、ハードウェアやOSのバグによるダウンの場合、OSがリブートされ、DBのジャーナルリカバリが行われるまでの数十分、致命的

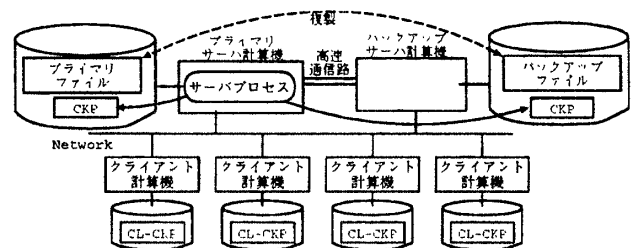


図 2: ARTEMIS の処理概要

図2は、ARTEMISの処理概要を示したものである。ARTEMISはチェックポイント・リスタート方式をベースに、以下の3つのキー技術により、上記の問題を解決する。

- 分散処理に対応した分散CP生成技術
- 2重化されたサーバ計算機で、1台の計算機に障害が発生した場合に他方で即時引継ぎを可能にする分散レプリケーション技術
- プロセス単位のCPをアプリケーションの変更なしに生成するジャケットルーチン

分散CP生成技術は分散システム上でプロセス間通信を行う複数のプロセスのCP生成を矛盾なく生成するためのプロトコルを実現している。

Overview and Checkpointing Mechanism of ARTEMIS (Advanced Reliable disTributed Environment MiddleWare System)

Toshio SHIRAKIHARA, Kiyoko SATO, Tatsunori KANAI: Communication and Information Systems Research Laboratories, Research and Development Center, TOSHIBA Corporation

Hideaki HIRAYAMA: Information & Communications Systems Laboratory, TOSHIBA Corporation

分散レプリケーション技術では、プライマリ計算機上のプロセスの CP 情報をバックアップ計算機側に転送するとともに、2 台のサーバ計算機間で複製ファイルをサポートする。これにより、プライマリ計算機障害時にバックアップ側での引継ぎを可能にする。

ジャケットルーチンは、プロセス起動時に動的にリンクされるダイナミックリンクライブラリで、リンクされたプロセスの CP 生成を行う。ジャケットルーチンは OS と同じシステムコール API を持ち、実行時に動的にリンクされるため、アプリケーションおよび OS の変更は不要 (バイナリコンパチ) である。

以降では、3 つのキー技術の 1 つであるジャケットルーチンについて説明する。

### 3 チェックポイント生成方式

ジャケットルーチンでは、プロセスが発行するシステムコールをフックし、プロセスの状態の保存 (CP 生成)・回復 (リスタート) を行う。プロセスの状態は OS に依存する部分が多いが、本稿では、UNIX<sup>1</sup> (Solaris<sup>2</sup>) の場合について述べる。

プロセス資源はアドレス空間およびレジスタセットといったプロセスから直接アクセスできる資源と、ファイル、共有メモリ、socket といった OS により提供される資源の 2 つに大別できる。アドレス空間の情報は、前回の CP 生成時から現在までに更新されたページ (ダーティページ) だけを保存する差分保存方式をとっている [3]。レジスタセットに関しては、setjmp/longjmp により保存・回復を行う [4]。OS 資源についてはプロセスが OS 呼び出しを行うシステムコールをジャケットルーチンでフックし、必要な情報を保存・回復する。

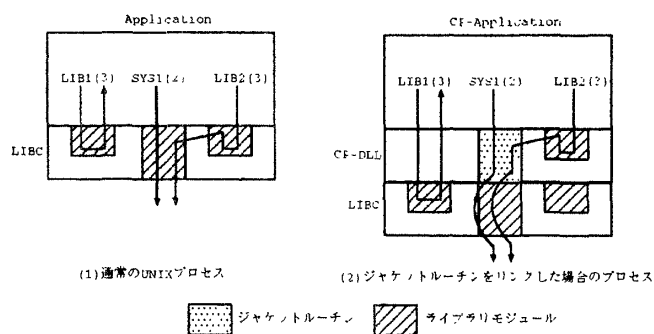


図 3: システムコールのフック方法

図 3 はシステムコールをフックするための仕組みを示した図である。通常のプロセスは図 3(1) に示すように OS をアクセスするためのシステムライブラリ (LIBC) がリンクされる。LIBC は通常、以下の 3 種類のインタフェースを提供する。

- システムコール (SYS1、ex. open)
- ライブラリコール (LIB1、ex. strcpy)
- 間接的システムコール (LIB2、ex. fopen)

これに対して ARTEMIS では、図 3(2) に示すようにジャケットルーチンが LIBC より先にリンクされる。ARTEMIS では必要なシステムコールをフックするジャケットルーチンを持っており、アプリケーションがシステムコール SYS1 を発行した場合、LIBC の SYS1 ではなく、ジャケットルーチンの SYS1 が呼ばれる。また、LIB2 のように間接的に SYS1 を呼び出すような関数に関しては、LIBC から LIB2 のオブジェクトを取り出し、ともにリンクする。この場合 LIB2 からの SYS1 の呼び出しは SYS1 を呼び出すように解決されるため、最終的にはジャケットルーチン SYS1 が呼び出される。

このように、ジャケットルーチンをアプリケーションにシステムライブラリより先にリンクすることにより、アプリケーションの変更を必要とせず、そのアプリケーションに関する OS 内の状態を保存することが可能になる。

### 4 おわりに

本稿では、高信頼化ミドルウェア ARTEMIS の概要とプロセスの CP を生成する仕組みであるジャケットルーチンについて述べた。アプリケーションに関する OS 内の状態の保存・回復方法はこの技術のポイントの 1 つであるが、その詳細については別の機会で述べたい。

### 参考文献

- [1] 佐藤他, “高信頼化ミドルウェア ARTEMIS の分散チェックポイント生成方式”, 情処第 54 回全国大会, 1997.
- [2] 平山他, “高信頼化ミドルウェア ARTEMIS の分散レプリケーション方式”, 情処第 54 回全国大会, 1997.
- [3] J. S. Plank, M. Beck, G. Kingsley and K. Li, “Libckpt: Transparent Checkpointing under UNIX”, USENIX Winter 1995 Technical Conference, 1995.
- [4] 森山他, “利用者レベルで実現したプロセス移送ライブラリ”, 情処 OS 研究会報告 91-OS-51, 1991.

<sup>1</sup>UNIX は X/Open の商標です。

<sup>2</sup>Solaris は米国 Sun Microsystems 社の商標です。