

分散OS COSMOS 2における負荷分散

6P-1

李在己 計宇生 相田仁 斉藤忠夫
東京大学 工学部

1. はじめに

COSMOS 2は当研究室で実装されている統合型分散OSである。これは従来OSが持つ主な機能のプロセス管理、ファイル管理、入出力管理などに通信処理機能を加えて、これらの諸機能を論理機能モジュールとして分割し、高速なLAN上に結ばれている各ノードに分散配置する。このように分散された各論理機能モジュールは、COSMOSカーネルバックプレーン(CKB)を介して相互協同することにより、システム全体を一つのTSSのように動作させることができる。これにより、分散システムの特徴である、高価な資源の共有、機能と負荷の分散、高性能、拡張性、信頼性などが得られる。

本稿では、COSMOS 2における負荷分散の方法について述べる。

2. COSMOS 2の構成

図1は一般化したCOSMOS 2のシステム構成を示している。

COSMOS 2では、図1のように独立した論理機能モジュールを物理的な各ノードに適当に分散配置する。これによりネットワークトランザクションに柔軟に対応するノード構成となり、高度な負荷分散が可能である。

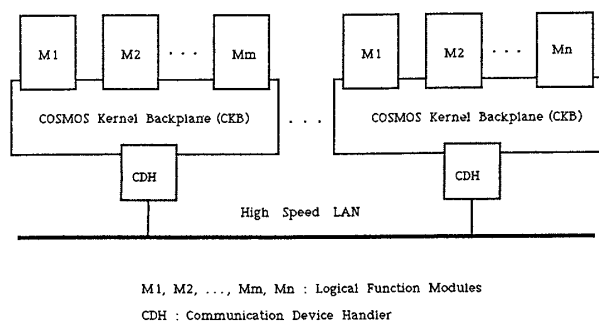


図1. COSMOS 2のシステム構成

各ノードを構成する論理機能モジュールはマネージャ、サーバとデバイスハンドラの3つに大別される。

マネージャモジュール(Process Manager, File Managerなど)はシステム内に1つのみであり、ネットワークワイドな各種資源と情報を統合管理する。各種のサーバ(Process Server, File Server, Print Serverなど)はシステム内に複数存在することができ、ユーザからの要求に対して、各々指定されたサービスを提供するモジュールである。

デバイスハンドラ(Disk Device Handler, Communication Device Handlerなど)は各種の入出力機器を直接に制御する機能を持つモジュールである。^[1]

CKBではノードの初期化、ノード内の論理機能モジュール間の通信とスケジューリング、及びLANを介したノード間の通信を司る。特に、COSMOS 2では、この2つの通信がユーザに透明に行われる。CKBに対して、各論理機能モジュールは1つのソフトウェアボードのような構成を成す。特に、COSMOS 2では、この2つの通信がユーザに透明に行われる。

通信デバイスハンドラは^[2]ノード間のメッセージ通信に使われる部分である。ここだけを変更することにより、各種のLANに対応することができる。

3. COSMOS 2の負荷分散

負荷分散とは、システム中の各ノードが均等な負荷を持つようにすることである。これを用いることにより、システム内のアイドルノード或は軽負荷のノードと過負荷のノードに均等な負荷を割り当てて、システム全体の性能を上げることができる。

COSMOS 2では、ある意味で2段階の負荷分散を行っている。第一段階はシステム構築の際に、システムマネージャが静的にノードの構成を決めることにより行われる。これは機能分散による方法である。第二段階の負荷分散はプロセスサーバ内において、プロセスの実行時、動的に行われる方法である。^[3]

ここではCOSMOS 2の第二段階負荷分散の方法と、これを実現するための関連論理機能モジュール間の動作について述べる。

3. 1 負荷の測定

COSMOS 2において、負荷分散に直接に携わる論理機能モジュールは、システム内の唯一なプロセスマネージャ (PM) と全プロセスサーバ (PS) である。

PSはUnixと同じ概念を持ち、プロセスの生成、実行、終了及びスケジューリングを行う。PMにはマネージャプロセスがあり、このプロセスが各PSとのトランザクションを行う。マネージャプロセスはPSからの要求を順次に処理する。

各PSは自分の負荷と、システム内の全PSが持つ負荷の平均値 (これをシステム負荷と呼ぶ) を維持する。これに対し、PMはシステム内の全PSの負荷と、システム負荷を管理する。

次に単一PSの負荷 (η_p) とシステム負荷 (η_s) の計算式を示す。

$$\eta_p = N_p / T_{ps} \quad \dots\dots\dots (1)$$

$$\eta_s = \sum \eta_p / P_n \quad \dots\dots\dots (2)$$

ここで、 N_p はPS上のユーザプロセスの実行待ち行列の長さであり、 T_{ps} は1つのノード内において、全論理機能モジュールに対するPSの実行時間の割合である。そして P_n はシステム中のPSの数である。

3. 2 負荷分散の方法

COSMOS 2での負荷分散は、PS上のプロセスの実行時、式(1)と(2)により求められた自分の負荷とシステム負荷を比較して、PMと相手のPSとのネゴシエーションにより決められる。これはリモート実行による負荷分散であり、リモート実行の必要条件是 $\eta_p > \eta_s$ である。

図2はCOSMOS 2でのPSとPM間のリモート実行の手順を示している。

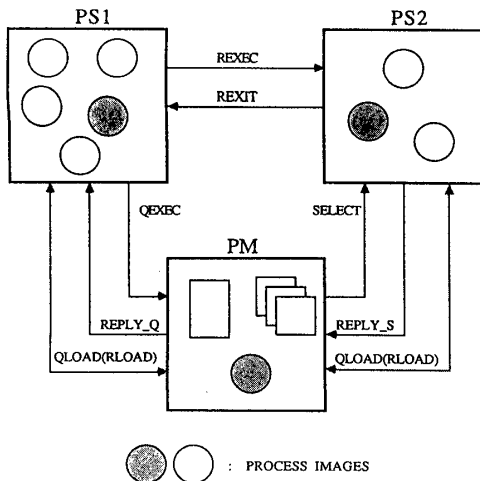


図2. COSMOS 2のリモート実行の手順

ここでPS1からリモート実行の必要が生じた場合、PS1はPMに自分の負荷とともにリモート実行の要求メッセージを送る。これに対してPMは負荷管理テーブルを参照して、PSの中で一番負荷が軽いPSを選ぶ。PMは選択されたPS (ここではPS2) に現在のシステム負荷を送り、リモート実行の可否を要求する。PS2はリモート実行の可否を決めて、PMに応答する。PS2から可の応答が来ると、PMはリモート実行を要求したPS1にリモート実行の宛先を送って応答する。その後、PS1は直接にリモート実行をPS2に依頼することができる。PS2での実行が終わると、実行の結果をPS1に送りリモート実行を終了する。もしPS2の応答が否なら、PMは他のPSを選び、同じ処理を行う。

一方、PMと各PS間の負荷の更新のやり方は、各PSが一定時間 (数秒) ごとに自分の負荷を計算して、PMに報告する。PMはこれに対する応答として、その時点のシステム負荷をPSに送る。一定時間内に、あるPSからの負荷の報告がなければ、PMはそのPSに対してポーリングを行い、システムの負荷を維持する。

4. おわりに

以上のように、COSMOS 2での負荷分散はシステム内の全PSの負荷をPMが効率的に管理することにより、リモート実行の際に動的に行われる。もしPMがダウンしても、各PSによるローカルな実行は正常に行われる。

この方法を拡張すると、システムのフォルトトレランスも容易に実現できる。

今後はこの方式の実装を進めて、なおノード間のモジュールの動的なマイグレーションについても検討するとともに、COSMOS 2の性能評価を行う予定である。

<参考文献>

[1] 計、李、鈴木、相田、斉藤、”分散OS COSMOS 2の実装(1) - OS機能モジュール”、情報処理学会第36回全国大会、4D-1、1988. 3

[2] 鈴木、李、計、相田、斉藤、”分散OS COSMOS 2の実装(2) - 通信サブシステム”、情報処理学会第36回全国大会、4D-2、1988. 3

[3] 松下、計、相田、斉藤、猪瀬、”分散オペレーティングシステムCOSMOS 2 - プロセス管理方式”、情報処理学会第34回全国大会、4B-8、1987. 3