

シミュレーションによる分散ファイルシステムの性能評価

5 G-2

金澤 裕治 村上 岳生 湯原 雅信
(株) 富士通研究所

1. はじめに

経済性・耐故障性などの理由から、計算機システムにおける資源共用の必要性が認識されつつある。中でも、長期記憶装置と、そのデータを共用することの必要性は大きく、実際、NFSなどに見るように、資源共用は長期記憶の共用から実用化されてきた。しかし、今後は、CPU資源の共有の必要性・ユーザーモビリティ・メンテナンスのしやすさなどの要請から、単にデータを共有するだけでなく、ネットワーク内の各計算機が全く同一の環境にあるように見えるような、ネットワーク透過性を有した分散ファイルシステムが必要になると考えられる。

我々は、UNIXをベースとした分散OSの研究を行なっている。分散ファイルシステム設計のために、分散ファイルシステムの性能予測が必要になるが、その性能は、アプリケーションがファイルシステムをアクセスするパターンに強く依存しているため、シミュレーション以外の方法で性能予測を行なうことは難しい。今回、そのためのシミュレーターを作成した。その1例として、分散makeのシミュレーションについて報告する。

2. シミュレータの機能

分散ファイルシステムの性能評価用シミュレーターの入力と出力として、以下のものが挙げられる。

入力:

- ネットワーク・ディスク・CPUの速度、台数、ファイルシステムのブロックサイズなどのパラメータ
- アプリケーションのファイルシステムアクセスパターン

出力:

- 評価したいファイルシステム上でそのアプリケーションを動作させたときの、実行時間・キャッシュヒット率
- (可能ならば) 誤差

また、キャッシュアルゴリズムなどの変更、シミュレートするシステムの構成の変更が容易にできるものが望ましい。

3. シミュレータの構成

分散システムでは、性質の異なる複数のプロセスが、プロセス間で通信を行なったり排他制御を行なったりしながら、CPU・ネットワークなどのハードウェア資源を順番に占有していくことにより動作が行なわれる。そのため、プロセスの論理的進行をシミュレートする論理モジュールが、物理デバイスをシミュレートする物理モジュールに作業を依頼しながらシミュレーションを進行させる手法を採用した。例えば、論理モジュールがCPU時間を消費するときは、CPUモジュールにメッセージを送り、指定したCPU時間が消費され、終了メッセージが返されてくるまで待つということを行なう。

これらのモジュールの他に、スケジューラーが1つ存在し、メッセージ・ロックなどのモジュール間の通信を、全てスケジューラーを介して渡すことにより、時間を管理している(図1)。

論理モジュールと物理モジュールに分けることには、以下のような利点が存在する。

(1) システムの構成・方式の変化に対応できる

例えば、ネットワークに接続しているホスト数を変える場合、最初に生成するCPUモジュールの数を増やすだけでよい。

また、ユーザーがCPUモジュールのプログラムを変更することにより、種々のスケジューリングのシミュレーションを行なうことも可能である。同様に、メモリを物理モジュールとして考えれば、ページングもシミュレートすることができる。

(2) プログラミングが容易になる

論理モジュールは、OS内のプロセス(スレッド)のイメージに近いので、自然にプログラムを行な

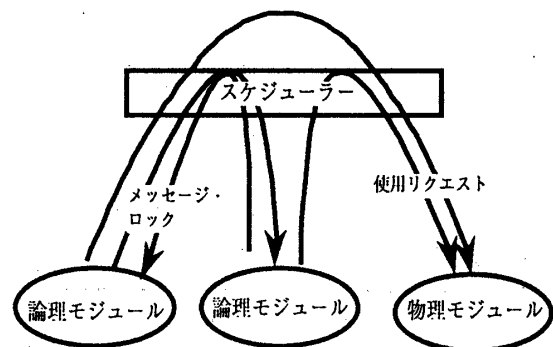


図1 シミュレータの基本構成

うことができる。

しかし、論理モジュールの動作が、実際のシステムでのプロセスの動作と完全に一致している必要はない。ハードウェアの制約により、1つの論理的流れを2つ以上のプロセスに分けている場合、実際のプロセスの流れではなく、論理的流れの方を重視することができる。

例えば、クライアントのCPUモジュール、ネットワークモジュール、サーバーのCPUモジュールと順に要求を送ることにより、rpcのシミュレートを1つの論理モジュールで行なうことも可能である。サーバー上にrpcを処理するプロセスが1つしかない場合のシミュレートを行なう場合は、サーバー上で処理される部分でロックをかけることによって可能である。

4. シミュレータの実装

SUNLWPを使って実装を行なった。

各ホストのCPU・ディスク・ネットワークをそれぞれ1つの物理モジュールとし、各モジュールとスケジューラを、1つのlwpスレッドとしてインプリメントしている。

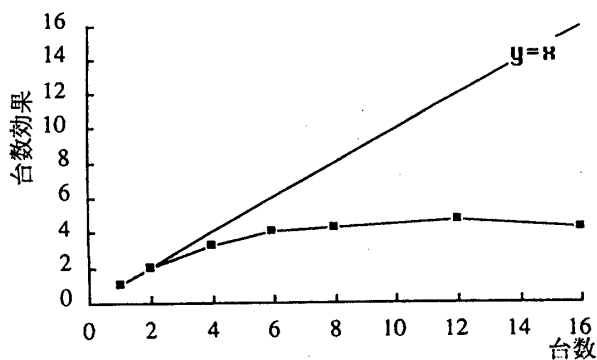
各論理モジュールは、システムコールと、システムコール間のユーザー時間を記録したデータを読み込み、それに従って動作を行なうことにより、アプリケーションの動作をシミュレートする。このデータは、プローブを埋め込んだOS上で実際にアプリケーションを動作させて得たものを使用している。

コード量はシミュレータの基本部分が1500行、ファイルシステムのシミュレーションを行なうためのコードが3000行程度である。

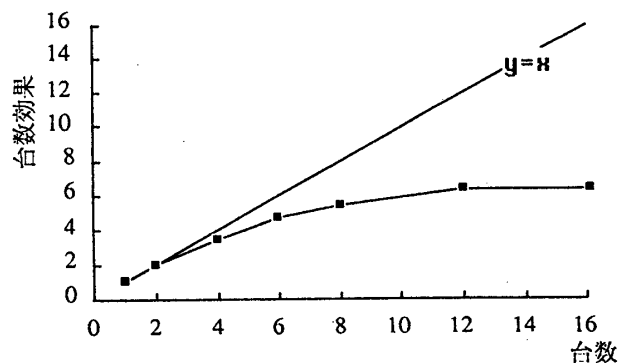
5. シミュレーション結果

Xのライブラリをネットワーク上で分散makeをシミュレーションした結果を示す。

システムの構成として、25 MIPS程度の計算機がイーサネットに接続されている場合を想定している。ファイルのライト共有がないので、一貫性プロ



グラフ1 分散makeのシミュレーション (/tmp共有)



グラフ2 分散makeのシミュレーション (/tmpローカル)

トコルは使用していない。シミュレーションに使用したキャッシュのヒット・ミス時のCPU時間・実時間などのパラメータは、以前測定したベンチマーク [1] での実測値を使用している。

グラフ1がテンポラリファイルのためのディレクトリ(/tmp)をサーバー上に置いた場合、グラフ2は/tmpを各クライアントがローカルに持っている場合のシミュレーション結果である。

/tmpを共有した場合は6台、共有しなかった場合は10台程度から台数効果が出なくなることが示されている。

これは我々が得ているSUNOSでの実際の測定値より多少悪い。このシミュレーションではディスク内のブロック配置の最適化を行っていないために、実際のシステムに比べてディスクアクセスに時間がかかっていることが原因だろう。

6. おわりに

分散ファイルシステムのシミュレーションに適したシミュレータを作成した。論理モジュールと物理モジュールを区別したことにより、ネットワーク形状やキャッシングアルゴリズムを変更してシミュレーションを行なう場合でも、柔軟に対応することが可能である。

今回報告した分散makeのシミュレーションだけではなく、ファイル内容のキャッシングのシミュレーションも現在進行中である。

これからの課題としては、大規模な分散システムもシミュレートできるような、並列実行を行うことや、現在の版では行っていない誤差の解析も行なうシミュレータを開発していくことが挙げられるだろう。

参考文献

[1] 金沢, 岸本, 湯原: "UNIXファイルシステムの性能測定", 第43回情報処理学会全国大会.