

ユーザ透過な耐故障性を実現する MPI へ向けて

高宮 安仁[†] 松岡 聡^{†,††}

コモディティクラスタリングシステムにおける、ノード数規模の拡大、計算実行時間およびメモリ空間の急激なスケールアップに伴い、アプリケーションおよびシステムの障害発生の可能性への対処が急務となっている。しかし、クラスタ等の並列計算分野では、これまでこうした耐故障性についてのソフトウェア開発が重視されておらず、十分ではなかった。また、信頼性、ユーザ透過性、実行時オーバーヘッドの兼ね合いをユーザが指定することのできる、柔軟な耐故障性機構が求められているが、従来のクラスタ向け耐故障性システムでは、単一のポリシー/機構専用のものがほとんどであった。加えて、実アプリケーションを用いた場合のオーバーヘッドも明らかではなかった。本稿では、耐故障性機構をもつ MPI である、Parakeet システムを提案する。Parakeet システムを用いることによって、ユーザは性能を損ねることなく、容易に耐故障性、リカバリのポリシー/機構を指定できる。本稿では予備段階として、ユーザレベルチェックポイント、プロセスマイグレーション、Coordinated Checkpointing を MPICH 上にユーザ透過に実装した。予備的な評価の結果、Parakeet システムは移植性を保ちつつ効率的であり、本研究の将来的な目標であるプラグアンドブレイクラスタリングの基礎技術として有用であることがわかった。

Towards MPI with user-transparent fault tolerance

YASUHIITO TAKAMIYA[†] and SATOSHI MATSUOKA^{†,††}

Rapid increase in the number of nodes as well as the massive scale of computing in terms of both time and memory space for commodity clustering is mandating the handling the potential failure of applications and system as the norm, while inherent fault tolerance and recovery have not been integral part of software tools being developed for parallel computing on such clusters. Moreover, flexible fault tolerance mechanisms in which the user could manage the balance of reliability vs. transparency vs. execution overhead would be vital, but most previous work on cluster fault tolerance have made available only a single policy and/or mechanism, and moreover, their overhead have not been exactly measured for practical applications. Instead, we propose a new fault tolerant MPI system called Parakeet which allows various fault tolerance and recovery mechanism could be easily specified by the user, while retaining the efficiency. As a preliminary basis, we have implemented a user-level, coordinated checkpointing and migration protocol on top of MPICH in a user-transparent fashion. By specifying new protocols based on the underlying Parakeet mechanism, one could achieve Plug-and-Play management of large-scale clusters. Preliminary benchmarks show that Parakeet is portable and efficient, and could well serve as a basis for Plug-and-Play clustering.

1. はじめに

数百～数千ノード超の実用的大規模 PC クラスタリングシステムの構築において、克服しなければならない問題として、システムの信頼性と容易な利用方法の確立が挙げられる。とくに、1) 期待されるクラスタリングシステム性能のアプリケーションレベルにおける保証、2) 耐故障性等のシステムの信頼性、3) クラスタリングシステムの運用性、アプリケーションへの適応性、等が挙げられる。

一連の問題のうち、1) に関しては、研究レベルでは解決されたと言ってよい。AM¹³⁾、Fast-Messages¹⁴⁾ など、ユーザーレベル通信ライブラリの研究によって、

ほぼ理想値に近い実効値を得ることができているからである。しかし、2) 3) については、十分な成果が得られていないのが現状である。

汎用 PC を多数結合するクラスタリングシステムでは、システム全体としての信頼性は単体の PC に比べ指数的に低下する。並列処理の代表的なアプリケーションである大規模な科学技術計算では、実行時間が数日から数月に及ぶことも珍しくなく、不意のシステムダウンは頻繁に起こる。こうした障害からの復帰の技術として、MPI をはじめとする、並列タスクのための低オーバーヘッドなローバルバック/リカバリプロトコルおよびプロセスマイグレーションなどの耐故障性機構が今後不可欠である。

本研究では、MPI との互換性を保ちつつ、MPI 上でユーザ透過なチェックポイントング、マイグレーションを行うことのできる Parakeet システムについて、

[†] 東京工業大学 Tokyo Institute of Technology

^{††} 科学技術振興事業団 JST

P4 ライブラリへの改造を行う方式、および Rocks⁵⁾ とチェックポイントを用いた方式について、予備的な実装を行った。評価では、現在実装が完了している、チェックポイントと Rocks の統合による Parakeet ライブラリの基本性能について、interpositioning されたシステムコールのオーバーヘッド、及び Application Test として NAS Parallel Benchmarks 実行時のオーバーヘッドについて調査した。分散チェックポイントティング、マイグレーションについては未完成のため、評価を行っていない。

評価の結果、Parakeet ライブラリによるオーバーヘッドは通信を行う場合でも数 % 以下であり、十分許容できるものであることがわかった。

2. 背景

MPI 上で耐故障性を実現する手法として、主に以下の 2 種類の手法が提案されている。

チェックポイントングを用いた手法^{1),12),15)}としては、既存の MPI を拡張し、並列タスク全体の無矛盾なスナップショットを保存する機構を提供するもの、クラスタ構成の変化イベントに対するリスナをユーザが定義し、イベントに応じてユーザが個々のプロセスのチェックポイントタイミングを明示的に指定する独自の MPI 等がある。これらのシステムでは、チェックポイントング、もしくはマイグレーションのみを実現している場合がほとんどである。チェックポイントングプロトコルとしては Coordinated Checkpointing が主流であり、その他のプロトコルを実装した例は無い。

一方、レプリケーションを用いた手法^{6),9)}では、コミュニケータの不整合の検出、回復のための API を提供するものや、MPI_Send() 等の引数として、レプリカプロセスを指定し、障害時には自動的にレプリカプロセスが実行を続行するシステムがある。これらの手法では、レプリケーション対象、障害からの復帰方法等を明示的にユーザが指定する必要があるため、既存のプログラムを大幅に変更する必要がある。

さらに、両者共通の問題として、Independent Checkpointing 時に起こりうる Dommino Effect の可能性の回避、ロールバックのためのメッセージ間の依存情報追跡、Log-Based Rollback Recovery のロールバック処理などをユーザが明示的に指定するのはコストが高く、各種プロトコルの深い知識が必要とされることがある。加えて、アプリケーションに対してどのプロトコルが有効であるかが一般に不明であるため、実際に複数のプロトコルの実装、比較が必要である。また、すべてのプロトコルを表現可能な汎用チェックポイントング/ロールバック API を定めるのは困難である。

3. 分散チェックポイントング

この章では、メッセージパッシングプログラミングモデルにおける、様々なロールバック/リカバリプロトコルを紹介し、比較する。この分野については、すでに理論面から数多くの研究がなされている¹⁰⁾。

並列タスクのような、通信を行う複数プロセスか

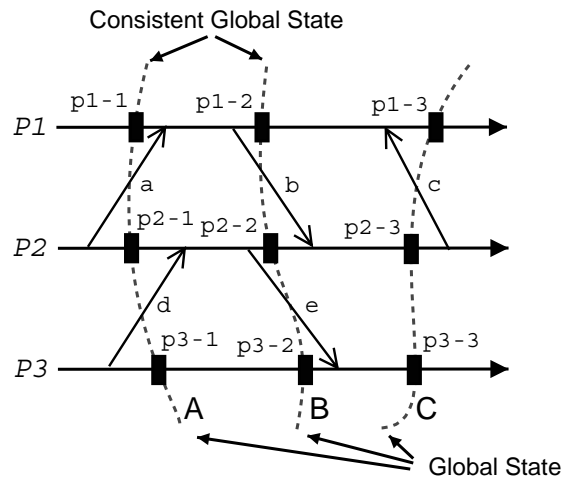


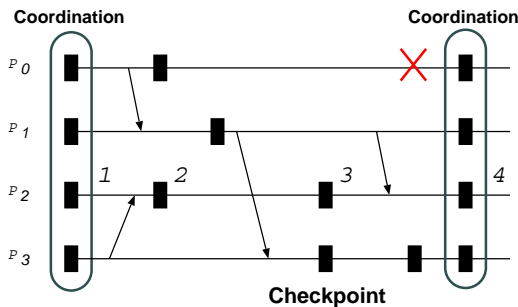
図 1 Consistent Global State

らなる系の状態は、各プロセス自体の状態 (volatile checkpoint) と、通信チャンネルの状態から成る。分散チェックポイントングとは、各プロセスの状態と、通信チャンネルの状態を Stable Storage へ保存し、障害時にリスタートする技術である。Chandy と Lamport によって、無矛盾な ロールバック/リカバリを満たすための条件が示されている³⁾。図 1 において、プロセス P1 の状態 P1-1, P2 の状態 P2-1, P3 の状態 P3-1, 加えて通信チャンネル a, d の状態を保存した Global State A は無矛盾にリスタートすることができる (Consistent Global State)。しかし、Global State C は無矛盾にリスタートすることができない。これは、Global State C からのリスタート時に、通信チャンネル C 上のメッセージが複製されてしまうからである。

メッセージパッシングプログラミングモデルにおけるロールバック/リカバリプロトコルは、以下の 2 種類に分類できる。

Checkpoint-based Rollback-Recovery はプロセスの復帰の際に、チェックポイントを用いて復帰するプロトコルである。上記 Consistent Global State を達成するための、Coordinated Checkpointing や、各プロセスが独立にチェックポイントングを行い、障害時には Dependency Tracking を行いロールバックする Independent Checkpointing、通信パタンによってチェックポイントングを決める Communication-Induced Checkpointing がある。

Log-based Rollback-Recovery は復帰の際に、チェックポイントに加え、非決定性イベントログを用いて復帰するプロトコルである。主なプロトコルとして、ロールバックの局所化を狙った Pessimistic Logging、ノンブロックロギングを行う Optimistic Logging、また、両者の利点を兼ね備えるが、ロールバック等がより複雑な Causal Logging 等がある。



```
mpirun -np 4 -ckpt_proto lazy interval 60
-laziness=4 -with-nonblockng ./foo
```

図 2 mpirun の引数による rollback-recovery protocol 指定. プロトコルとして Lazy Coordination, Coordination の間隔として laziness=4, Coordination の最適化として, non-blockng Coordination を指定している.

4. Parakeet システム

4.1 Parakeet システム概要

本稿で提案する Parakeet システムは, ユーザ透過なチェックポイントング, マイグレーションを用いて, MPI-1.2 プログラムへ耐故障性を提供するシステムである.

Parakeet システムは, ロールバック/リカバリ・チェックポイントング・マイグレーションを提供する Parakeet ライブラリと, 実行時のプロセス管理・モニタリング・チェックポイントングなどのイベント送信を行う Parakeet デーモンから構成される.

ロールバック/リカバリ機構は, 自動的なチェックポイントング・ロールバック・障害時のマイグレーションを提供する. 利用するプロトコルはアプリケーションや環境に依存するため, 複数のプロトコル (Checkpoint-Based や Log-Based) をサポートし, ユーザの指定により切り替えることが可能となっている. チェックポイントング/マイグレーションは, システムコールを interpositioning することで実現されている. この実装はユーザレベルで行われている.

実行時のチェックポイントング/ロールバックの起動は, ユーザによって指定されたプロトコルにより決定され, Parakeet デーモンによって各プロセスに起動シグナルが送られる.

ユーザによるプログラム中への明示的なプロトコルの挿入は不要であり, 既存の MPI プログラムを Parakeet ライブラリとリンクするだけで, Parakeet システムの提供するフォールトトレランスを利用できる.

ロールバック/リカバリプロトコルは, rc ファイルや mpirun の引数として指定できる (図 2). 指定できるパラメータには, プロトコルの種類, Coordination 間隔, チェックポイントング頻度などがある. このため, 1) プログラムのコーディング, 2) 複数プロトコルによるテスト実行, 3) 実際に用いるプロトコルの選択, 4) 実際の実行, というサイクルを迅速に行うことができる.

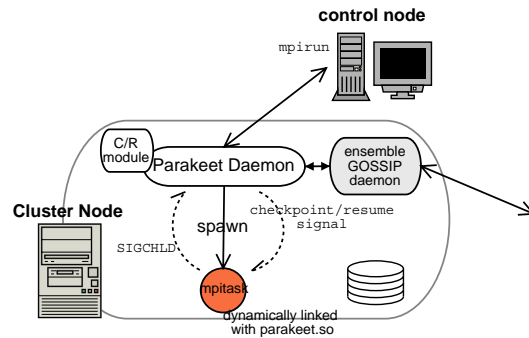


図 3 Parakeet システムアーキテクチャ

4.2 実装

Parakeet システムのアーキテクチャを図 3 に示す. 各ノード上では Parakeet デーモンが動作しており, MPI プロセスの生成, モニタリング, チェックポイントングシグナルの送信を行う. プロセスのクラッシュ, Coordination 開始, チェックポイント GC 開始といった各ノード上のイベントの通知には Ensemble GOSSIP デーモン⁴⁾ による, アトミックなグループ通信を行う.

今回, MPI 上へのチェックポイントング, マイグレーションの実装として, Rocks (Reliable Sockets)⁵⁾ を用いた方法と, p4 ライブラリへの改造を行う方法の 2 種類を実装した. まず, チェックポイントングのために, ソケット等の状態も復元可能なユーザレベルチェックポイントを実装した.

Rocks は, accept(), write() 等のシステムコールを interpositioning することにより, IP アドレスの変化, 物理層の切断といった障害後も, ユーザ透過に復帰できる仕組みを提供している. 実装としては, TCP/IP スタック中の Send/Receive バッファの他に, in-flight なメッセージを保存するバッファを設け, 送信側/受信側で送信/受信したパケット数をカウンティングすることによって, in-flight メッセージの損失を防いでいる.

今回の Rocks を用いた実装では, チェックポイントと Reliable Sockets を改造し, 双方が interpositioning しているシステムコールの部分を統合した. Coordinated Checkpointing 要求の際には, 各プロセスについてチェックポイントング終了したプロセスから順次実行をストップする. このことによって, チェックポイント後のプロセスから他プロセスへのメッセージ送信が起こらないため, Global State を逆向きにカットするメッセージが発生せず, Consistent Global State が実現できる. 全プロセスがストップ後, Parakeet デーモンから SIGCONT が送信され, リスタートする. チャンネル中の in-flight メッセージの再現については, Reliable Sockets のバッファ中のメッセージがチェックポイントングライブラリによってバッファごとチェックポイントされ, リスタート後に再送される.

p4 への改造を行う実装については, CoCheck⁷⁾ が行っているように, p4 チャンネルの flush 機能等

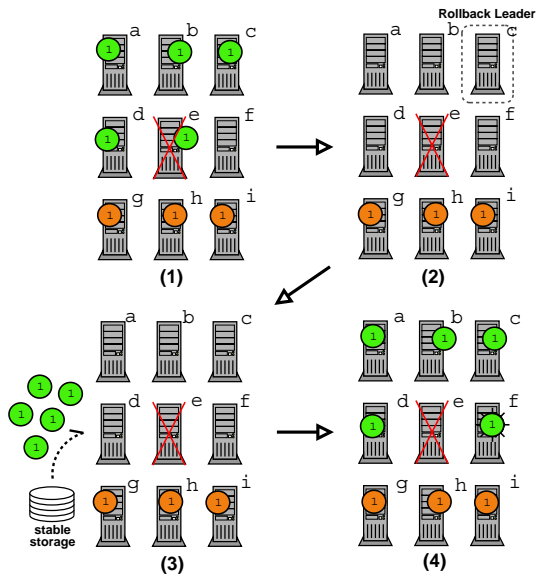


図 4 fail 時の rollback 機構

を追加した．現段階では実装が完全ではないため，評価は行っていない．

4.3 ユーザ透過なロールバック/マイグレーション
 ユーザ透過なロールバック/マイグレーションプロトコルは次のとおりである．(Coordination を行うプロトコルの場合)．

- (1) 図 4 (1) において，ノード a, b, c, d, e では並列タスク 1 を，ノード g, h, i では並列タスク 2 を実行しているとする．ここで，ノード e が突然クラッシュすると，MPI タスク全体が異常終了する．
- (2) 生き残っているノード上の Parakeet daemon は SIGCHLD を受けとり，生き残った Parakeet daemon 間で Leader Election を行う．Leader となったノード上の Parakeet daemon は，rollback メッセージを他の Parakeet daemon へブロードキャストする (図 4 (2)) ．
- (3) Stable Storage に保存されている前回 Coordinated Checkpoint 時のイメージを，対応する各ノードに戻し，クラッシュしたノードで動作していたプロセスイメージは，他の健全なノードへマイグレーションする (図 4 (3)) ．
- (4) マイグレーションしたプロセス以外のプロセスは，チェックポイントされた状態からのリカバリ時，マイグレーションを行ったプロセスとのチャンネルについて，ソケットディスクリプタを更新し，リスタートする．

Coordinated Checkpointing を行う場合，並列タスク全体がロールバックすることに注意されたい．この例

現在の実装では，Stable Storage は各ノードのローカルディスクとなっている．ノードクラッシュ等の障害からすみやかに回復するためには，Stable Storage を専用のチェックポイントサーバ上に設けるが⁷⁾，近隣ノードのディスクもしくはメインメモリ上に設ける必要がある．

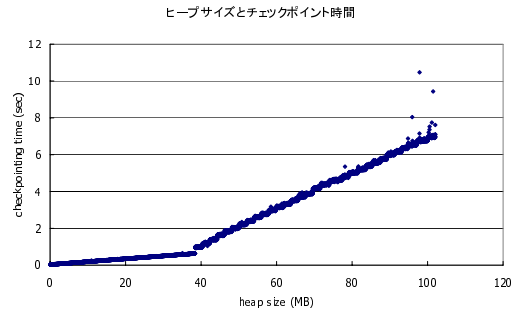


図 5 ヒープサイズとチェックポイント時間

の場合では，Coordination 間隔が 1 時間であった場合，ロールバック時には 1 時間前の状態からリスタートすることになる．

チェックポイントライブラリはチェックポイント時にソケットをクローズ，リスタート時にコネクションを張り直す機構になっている，(4) でのソケットディスクリプタの更新は，この時に Parakeet デモンへマイグレーションプロセスのマイグレーション先を問い合わせ，書き換えを行っている．Reliable Sockets による実装ではコネクションの uniformity を提供しているのので，この操作は行われていないが，Rocks の制限により，同時にマイグレーションできるプロセスの数は 1 つという制限がある．

5. 実験

実験に用いる環境として，東京工業大学松岡研究室の ion クラスタ²⁾ (CPU: Mobile Celeron 600MHz, Memory: 128MB, HDD: 20GB, OS: Linux 2.2.17 with PAPI 1.1.5 \times 16 nodes, Network: 100base-T switching hub) を用いた．計時関数として，PAPI¹⁸⁾ の PAPI_get_real_usec(3) を用いた．実験では，Parakeet ライブラリの基本性能と，実アプリケーションを用いた性能の評価を行った．なお，マイグレーション，Coordinated Checkpointing については，実装が完了していないため，評価を行っていない．

5.1 Kernel Tests

Parakeet ライブラリの基本性能評価として，

- sbrk() を連続して行うプログラムについて，ヒープサイズとチェックポイント時間 (図 5)
- ヒープサイズとチェックポイントファイルサイズ (図 6)
- チェックポイントファイルサイズとリスタートに要する時間 (図 7)
- netpipe-2.4 を用いた，Parakeet ライブラリリンク時のスループット (図 8)
- 逐次版 NAS Parallel benchmarks のチェックポイントライブラリのリンクによるオーバーヘッド (図 9)

について計測した．

図 5, 7, 6 に見られるように，ヒープサイズ/チェックポイント時間，ヒープサイズ/チェックポイントファイルサイズ，およびチェックポイントファイルサイズ/

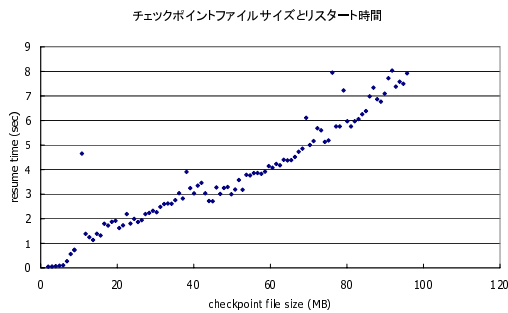


図 6 チェックポイントファイルサイズとリスタート時間

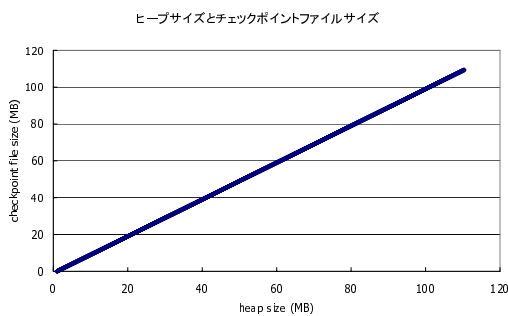


図 7 ヒープサイズとチェックポイントファイルサイズ

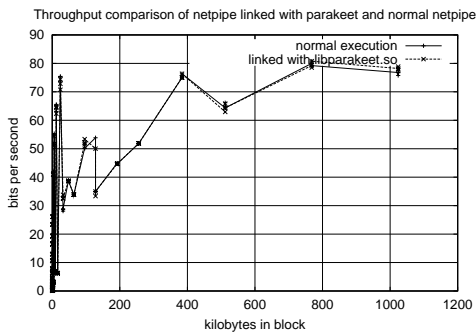


図 8 netpipe-2.4 による, Parakeet ライブラリリンク時, および通常実行時のスループット

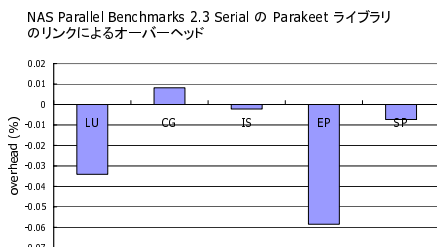


図 9 NAS Parallel Benchmarks 2.3 Serial の Parakeet ライブラリのリンクによるオーバーヘッド

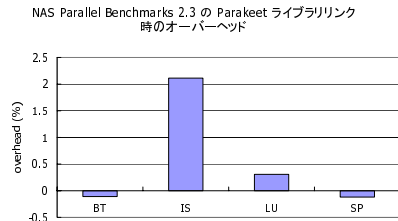


図 10 NAS Parallel Benchmarks 2.3 の Parakeet ライブラリリンク時のオーバーヘッド

リスタート時間はほぼ比例している．図 8 について，どのブロックサイズでもスループットはリンク有/無時でほとんど変わらないことがわかり，Rocks によるオーバーヘッドが低いことがわかる．図 9 では，いくつかのベンチマークについて，リンク無時の性能よりもリンク後の方が性能が良い．これは，interpositioning 対象のプログラムで実行されるシステムコールよりも，interpositioning により置き換わった後のシステムコールの方が高速であり，チェックポイントのオーバーヘッドが低いことがわかる．

5.2 Application Tests

NAS Parallel Benchmarks を用いて，Fail Free 時の Parakeet ライブラリのリンクによるオーバーヘッドを計測した (図 10)．NAS Parallel Benchmarks のような，`read()`，`write()`，を内部で頻繁に実行するプログラムでは，Rocks の影響が大きいと考えられる．しかし，計測結果より，オーバーヘッドは最大でも 2% 程度と低く，Rocks のオーバーヘッドはほとんど無いことがわかる．

6. 関連研究

2 章で述べたように，MPI にチェックポイントング，マイグレーションを実現しようとする研究は多数行われているが，未実装であったり，PC クラスタで主に用いられる Linux 等に対応したものは少ない．

CoCheck⁷⁾ は MPI に `sync-and-stop` を実装し，coordinated checkpointing およびマイグレーションを実現したものである．CoCheck は Condor 上での並列タスクチェックポイントングをサポートする目的で開発されたが，`sync-and-stop` 等のオーバーヘッドが高いため，実用には至っていない．

Starfish¹⁾ は Parakeet と同様，Ensemble グループ通信ライブラリを用いており，ノードの追加/削除等に対するイベントリスナを定義し，状況に応じたプログラミングを行える．Starfish は Ocaml のバイトコード実行であるため，そのオーバーヘッドが大きい．チェックポイントング API により，Coordinated Checkpointing 以外のプロトコルも対象としているが，API の詳細，各プロトコルの表現能力等は不明である．

SCore-D¹⁶⁾ には，時分割スケジューリングで用いられるギャングスケジューリングのために，ネットワークチャンネルの状態の退避，復元を行うネットワークプリエンブション機構がある．西岡らはこの機構を

用いて, sync-and-stop 方式に類似した方法で, 通信プロトコル透過にコンシステントチェックポイントを実現している¹⁵⁾. ただし, SCORE-D におけるネットワークコンテキストはネットワークポロジータ的な絶対位置に依存しているため, マイグレーションについては実現していない.

MPI/FT¹⁹⁾ は nMR (n-Modular Redundancy) とチェックポイントングを用いており, MPI の種類 (MPI-1.2, MPI-2), アプリケーションのクラス (SPMD, Master/Slave) 等に応じて nMR のモジュール数, rank 数等の nMR のパラメータ, およびチェックポイントングの有無を決定する. MPI/FT は fail-stop ではなく, fail-through を目的としており, そのための障害検知レイヤを提供している.

FT-MPI⁹⁾ は MPI-2 を拡張し, MPI_Send() 等の関数が fail した場合にコミュニケータの状態が invalid となった場合, コミュニケータを新たに作りなおすことにより, 復帰するための関数を定義している. Implicit Fault Tolerance⁶⁾ はプログラミングモデルを SPMD Master-Worker のみに限定し, FT-MPI で行われるレプリケーション/復帰をある程度自動的に行う.

7. まとめと今後の課題

我々は, MPI 上へのユーザ透過なチェックポイントング/マイグレーションの実現を, p4 ライブラリへの改造と Rocks を用いた方法によって試みた. 本稿では予備段階として, ユーザレベルチェックポイントと Rocks を統合し, その基本性能について, interpositioning されたシステムコールのオーバーヘッド, 及び NAS Parallel Benchmarks 実行時のオーバーヘッド等について調査した. 結果, オーバーヘッドは最大 2% と低く, 十分許容できるものであることがわかった.

今後の課題として, マイグレーション機構等の未完成部分について, 実装を完了させる必要がある. ボトルネックであるチェックポイントングの高速化のため, Checkpointing, Incremental Checkpointing, Checkpoint Compression¹¹⁾ 等の最適化を実装する必要がある. また, 残る Rollback-Recovery プロトコルについても実装を行い, 比較する必要がある.

Parakeet を用いたクラスタの動的な再構築, スケジューリングに関する試みとして, 我々はテストベッドとして ion プラグアンドプレイクラスタ²⁾ を構築している. 今後は, Dynamite PVM⁸⁾ で行われているような, マイグレーションによるロードバランシングの試み, および高速クラスタインストーラ Lucie²⁰⁾ とマイグレーションによる, クラスタ運用中のノードの動的な追加, 削除といった可用性についても調査する.

参考文献

- 1) A. Agbaria et al.: Starfish: Fault-Tolerant Dynamic MPI Programs on Clusters of Workstations., HPDC (1999)
- 2) ion Cluster Web Page, <http://matsu-www.is.titech.ac.jp/sig/cluster-team/ion/ion.html>

- 3) K.M. Chandy et al.: Distributed Snapshots: Determining Global States of Distributed Systems., TOCS Vol. 3-1,(1985)
- 4) M. Hayden: The Ensemble System, Technical Report TR98-1662, Department of Computer Science, Cornell University (1998)
- 5) Rocks: <http://www.cs.wisc.edu/%7Ezandy/rocks/>
- 6) P. Evripidou et al.: Implicit Fault Tolerance, http://www.cs.ucy.ac.cy/%7Eeskevos/html/ft_for_mpi.html.
- 7) G. Stellner: CoCheck: Checkpointing and Process Migration for MPI, Proceedings of the International Parallel Processing Symposium (1996)
- 8) Dynamite: <http://www.hoise.com/dynamite/>
- 9) G. Fagg et al.: FT-MPI: Fault Tolerant MPI, supporting dynamic applications in a dynamic world, LNCS 1908, (2000)
- 10) E. Elnozahy et al.: A survey of rollback-recovery protocols in message-passing systems, Department of Computer Science, Carnegie Mellon University, October, (1996)
- 11) J.S. Plank et al.: Libckpt: Transparent Checkpointing under Unix, Usenix Winter 1995 Technical Conference, pp. 220-232 (1995)
- 12) D. H. J. Epema et al.: A worldwide flock of Condors: load sharing among workstation clusters, Delft University of Technology, Department of Technical Mathematics and Informatics, Technical Report DUT-TWI-95-130 (1995)
- 13) A.M. Mainwaring et al.: Active Messages: Organization and Applications Programming Interface. (1995)
- 14) S. Pakin et al.: Fast Messages (FM): Efficient, Portable Communication for Workstation clusters and Massively Parallel Processors., IEEE Concurrency, Vol. 5, No2, pp.60-73 (1997)
- 15) 西岡 et al.: クラスタにおけるコンシステントチェックポイントの実現, JSPP 1999, pp. 229-236 (1999)
- 16) Atsushi Hori et al.: An Implementation of Parallel Operation System for Clustered Commodity Computers, Cluster Computing Conference (1997)
- 17) The NAS Parallel Benchmarks, <http://www.nas.nasa.gov/Software/NPB/>
- 18) PAPI: <http://icl.s.utk.edu/projects/papi/>
- 19) Batchu et al: MPI/FT: Architecture and Taxonomies for Fault-Tolerant, Message-Passing Middleware for Performance-Portable Parallel Computing, CCGRID (2001)
- 20) Lucie: <http://matsu-www.is.titech.ac.jp/%7Eetakamiya/lucie/>