

DMTCP によるノード構成の動的変更に対応した並列分散処理環境の検討

杉山 裕紀[†] 澤田 祐樹^{††} 大津 金光^{††} 横田 隆史^{††} 大川 猛^{††}[†]宇都宮大学工学部情報工学科 ^{††}宇都宮大学大学院工学研究科情報システム科学専攻

1 はじめに

近年、モバイル端末の高性能化が進み、新しい並列分散アプリケーションの基盤として注目されている。我々は Android 端末を利用して、並列分散アプリケーションを実行する Android クラスタコンピュータシステムを開発している [1]。本クラスタコンピュータは持ち運びが容易なモバイル端末をノードコンピュータとして利用しているため、並列分散アプリケーション実行中に一部のノードが脱退した場合、処理を途中で中断しなければならない。このため、本システムは並列分散アプリケーション実行中にチェックポイントデータを取得し、これを用いて処理を再開できるチェックポイント/リスタート機能を実装している。また、リスタート時において、最適な負荷分散を行うために、クラスタ内でのノード構成を把握し、その構成に応じて最適な配置でアプリケーションの実行をリスタートするノード構成変更機構を備えている。しかし、現時点でノード構成の動的変更機能は、一部手動で行われており、その完全自動化に必要となる、リスタート処理の開始条件等が未決定のままである。そこで、本機能の完全自動化を実現するための検討を行う。

2 Android クラスタシステム

本クラスタシステムでは、ノード間の通信には高速かつ容易に利用可能な Wi-Fi を用いている。また、並列分散処理のフレームワークとして Open MPI を用いている。本システムは、移動体である Android 端末をノードコンピュータと用いるため、クラスタを構成する一部のノードコンピュータの脱退や、新たなノードの参入が起こりえる。MPI アプリケーションの実行開始には、並列処理に使用するノードの情報を記述したマシンファイルが必要である。MPI アプリケーションを起動するノード（ホストノード）がマシンファイルを参照することで、ホストノード以外のノード（リモートノード）にプロセスを投入する。しかし、本クラスタシステムではノードの構成が動的に変化するため、並列処理に使用可能なノードをあらかじめ把握することが難しいため、正確に把握する必要がある。

MPI アプリケーション実行中にノードが脱退した場合、脱退したノードと通信ができなくなるため実行

中のアプリケーションを継続することは難しい。そのため、本クラスタコンピュータには、MPI アプリケーションの処理を途中から再開する機能が実装されている [2]。

この機能は並列実行途中にチェックポイントングを行うことで実現されている。チェックポイントングに DMTCP [3] を用いている。DMTCP の管理下でアプリケーションを起動するために `dmtcp.launch` コマンドが用意されている。これを使ってアプリケーションを起動すると、ホストノードでは DMTCP の管理プロセスである `dmtcp_coordinator` が起動する。また、起動したアプリケーションの実行プロセス内には、アプリケーションの実行スレッドと DMTCP のチェックポイントスレッド (CT) が生成される。`dmtcp_coordinator` が CT にチェックポイントデータ (以降 `ckpt` データと呼ぶ) を取得する要求メッセージを送信することで、チェックポイントングを開始する。また、ユーザが `dmtcp_coordinator` に指示を送るために、`dmtcp_command` コマンドが用意されている。これを利用することで、チェックポイントング等を行える。取得した `ckpt` データにはプロセスを復元するために必要な情報があり、取得時に生成されたヘルパスクリプトを実行することで、`ckpt` データからプロセスを復元しアプリケーションをリスタートすることができる。

現在の Android クラスタシステムは、ノード構成の変化によるアプリケーションの中断を自動的に検知して、自動的に `ckpt` データから処理をリスタートすることができない。そこで、DMTCP を利用してノード構成の変更を検知して、自動的にアプリケーションをリスタートする機能を検討する。

3 ノード構成の動的変更への対応

MPI アプリケーションの実行開始時には、ノード構成を把握してマシンファイルを作成する必要がある。また、アプリケーションのリスタート時においても、`ckpt` データからアプリケーションをリスタートするためには、変更後のノード構成を把握する必要がある。これを実現する方法として並列実行アプリケーション全体を管理する存在である `dmtcp_coordinator` が動的なノード構成の把握と、リスタート処理の制御を行う方法が考えられる。しかし、`dmtcp_coordinator` は管理下にあるプロセスが全て終了した後に自身も終了する。そのため、`dmtcp_coordinator` が停止している間にも、クラスタ全体を制御するプロセスが存在している必要がある。そこで、`dmtcp_coordinator` とは別にクラスタ全体を制御するためのプロセスを立ち上げ、

Consideration on Parallel and Distributed Processing Environment Allowing Dynamic Change of Node Configuration with DMTCP

[†]Hiroki Sugiyama, ^{††}Yuki Sawada, ^{††}Kanemitsu Ootsu, ^{††}Takashi Yokota and ^{††}Takeshi Ohkawa

Department of Information Science, Faculty of Engineering, Utsunomiya University ([†])

Department of Information Systems Science, Graduate School of Engineering, Utsunomiya University (^{††})

