

## Hadoop におけるジョブトラックの耐故障機能の検討

黒松 信行<sup>†</sup> 置田 真生<sup>†</sup> 萩原 兼一<sup>†</sup>

### 1. はじめに

クラスタ環境において MapReduce プログラム<sup>1)</sup> を実行するためのマスタ・ワーカ型フレームワークとして Hadoop<sup>2)</sup> が注目を集めている。Hadoop ではマスタであるジョブトラック (以降, JT) がタスクをスケジューリングし, ワーカであるタスクトラック (以降, TT) にタスクを割当てることでジョブを実行する。Hadoop はジョブの実行に数十から数千の TT を用いる。したがって, TT の耐故障機能が必要となる。Hadoop はタスクを冗長化することで TT に耐故障性を備えている。

一方, JT は耐故障性を備えておらず, 単一故障点である<sup>2)</sup>。ジョブの実行中に JT が故障すると, ユーザは JT を再起動し, ジョブを再実行する必要がある。

本研究の目的は, ジョブの実行中に JT が故障した場合, ユーザに対して透過的に JT を復元しジョブを継続することである。想定する故障の種類は, 故障停止のみとする。以降, 故障停止を単に故障と表記する。

### 2. 提案手法

#### 2.1 透過的な耐故障機能

透過的な JT の復元のためには, 故障時にジョブを継続する機構および故障検知が必要である。

##### 2.1.1 故障時にジョブを継続する機構

故障時にジョブを継続する機構にはチェックポイントイング (以降, CP) およびロールバックを用いる。JT を実行するノードと異なるノード上でバックアッププロセス (以降, BP) を起動する。JT は, タスク実行が一定数終了する毎にスナップショット (以降, SS) を作成し, BP に送信する。SS は JobTracker クラスの全フィールドを持つ。BP は JT の故障通知 (後述) を受信すると, SS を基に JT を復元 (ロールバック) しジョブを継続する。

ただし, CP を用いた復元には情報消失の可能性が

ある。最後の SS 作成から障害発生までの処理に関する情報は失われる。その結果, タスクの重複割当および永遠に終了しないタスクが発生する。

しかし, Hadoop は重複割当を許容する機能を備える。また, 長時間終了しないタスク実行を破棄し, 再度割当てて機能を備える。したがって, 情報消失が発生しても問題にならない。

##### 2.1.2 故障の自動検知

故障検知の機構には, タイムアウトを用いる。Hadoop において, TT はタスクの進捗報告としてハートビート<sup>2)</sup> (以降, HB) を定期的に JT へ送信している。そこで, TT は HB の送信に一定回数失敗すると JT が故障したと判断し, BP に故障通知を送信する。JT の復元まで待機した後, TT は HB の送信先を新しい JT に変更する。

ただし, この故障検知には誤検知の可能性がある。例えば, ネットワークの一時的な障害が発生すると TT は HB を送信できない。この場合, JT が故障していないにもかかわらず, 一部の TT は BP に故障通知を送信し, HB の送信先を新しい JT に変更する。その後直ちにネットワークが復旧すると, 誤検知を行った TT と他の TT では HB の送信先が一致しない。2 つの JT が独立にタスクを割当てている状態になるため, タスクの実行効率が低下する。

この問題を解決するため, BP は新しい JT を起動する前に元の JT に停止要求を送信する。これにより, 誤検知が発生しても数秒後に JT の数は 1 つとなる。

#### 2.2 オーバヘッドの見積

提案手法によるオーバヘッドは CP に伴うものおよび復元に伴うものの 2 つである。

##### 2.2.1 CP に伴うオーバヘッド

CP に伴うオーバヘッドの原因は, SS 送信の際 TT においてタスクの割当て待ちが発生するためである。

最悪時の CP に伴うオーバヘッド  $T_c$  を予測する。タスク実行が終了すると, TT は HB を用いて JT に終了通知を送信し, 次のタスクの割当てを要求する。この際, SS 作成の条件を満たすと, JT において SS 作成のために排他制御が必要となる。 $T_c$  が最悪となる

<sup>†</sup> 大阪大学 大学院情報科学研究科  
The Graduate School of Information Science and Technology, Osaka University

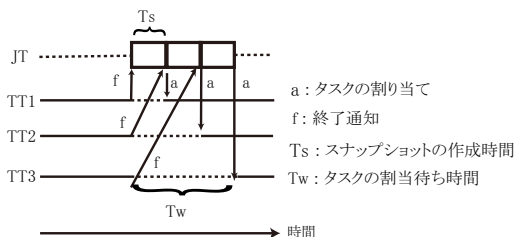


図 1 最悪のオーバーヘッドが発生する実行状況 ( $J = 1, P = 3$ )

場合は、JT に終了通知の到着が集中し、排他制御が連続する場合である。最悪時の実行状況の例を図 1 に示す。SS の作成時間を  $T_s$ 、TT の数を  $P$ 、SS の作成条件をタスク実行が  $J$  個終了する毎とすると、タスクの割当待ち  $T_w$  の最大値は  $T_w = T_s * (P/J)$  と表せる。さらに、ジョブを構成するタスク数を  $N$  とすると 1 つの TT に対するタスクの割当は約  $N/P$  回発生する。よって  $T_c = T_w * (N/P) = T_s * (N/J)$  と予測できる。

### 2.2.2 復元に伴うオーバーヘッド

復元に伴うオーバーヘッドの原因は、故障から復元までの間にタスクの割当待ちが発生するためである。この時間  $T_r$  の最大値は、故障検知に要する時間と JT の復元に要する時間の和となる。これらの時間は  $P, N, J$  によって変化するが、その変化量はジョブの実行時間と比較して十分小さい。そのため、 $T_r$  はほぼ一定とみなせる。

## 3. 評価実験

故障時に JT が透過的に復元することを確認する。さらに、提案手法に伴うオーバーヘッドを評価する。実験には Hadoop のバージョン 0.20.2 を使い、ギガビットイーサネットで接続された 16 台のクラスタ上で実行した。各ノードマシンのスペックは OS が CentOS 5.2、CPU が Intel Xeon 3.4 GHz Dual、メモリが 2 GB である。

### 3.1 透過的な復元

故障時の状況として、次の 3 つを想定する。

- (1) 復元時に情報消失が発生しない
- (2) 情報消失が発生する
- (3) 故障を誤検知する

JT の故障はジョブの実行中に JT を kill コマンドにより強制終了することで再現した。

まず (1) の場合、BP が JT を復元し、ジョブを継続できた。その後ジョブは正常に完了し、得られた結果は正しかった。次に (2) の場合、JT は未割当のタスクの報告を受けると、タスクを破棄し新たにタスク

表 1 オーバヘッド  $T_c$  の予測値と実測値の比較 ( $N = 163, P = 16$ )

$J$	6	12	18
全体の実行時間 (s)	515.7	505.3	498.6
$T_c$ の予測値 (s)	48.9 (9.4%)	24.5 (4.8%)	16.3 (3.2%)
$T_c$ の実測値 (s)	24.1 (4.6%)	13.7 (2.7%)	7.0 (1.4%)

括弧内は全体の実行時間に占める  $T_c$  の割合

を割当てた。また、10 分 (Hadoop における設定で変更可能) 経過しても終了しないタスクが存在した場合、JT はそのタスクを他の TT に割当てた。これらの結果、ジョブが完了し正しい結果を得られた。最後に (3) の場合、全ての TT が新しい JT に HB の送信先を切り替えてジョブを継続し、正常に完了できた。

### 3.2 オーバヘッドの評価

$J$  を変化させながら Hadoop におけるベンチマークの 1 つである WordCount プログラムを実行した。入力には 22 種類の単語を繰り返し並べた 10 GB のテキストファイルを用いた ( $N = 163$ )。  $T_s$  の実測値は 1.8 秒であった。また、 $T_r$  の実測値は 19 秒であった。表 1 にオーバーヘッド  $T_c$  の予測値と実測値を示す。表 1 が示すように、 $T_c$  の実測値は予測値と比較して約半分であった。

この原因は、実際には各タスクの終了通知が分散するためである。分散の増大に従って  $T_w$  は減少する。1 回のタスク割当待ちによるオーバーヘッドが減少したため、全体のオーバーヘッド  $T_c$  が減少した。

## 4. まとめ

Hadoop において単一故障点である JT に対して CP を用いて耐故障性を向上した。JT の故障検知には HB によるタイムアウトを用いた。提案手法に伴うオーバーヘッドは最大 4.6% であった。

謝辞 本研究の一部は科学研究費補助金 (基盤研究 (B)2330007) の支援を受けた。また、有益な助言を頂いた大阪大学の土屋 達弘 准教授に感謝します。

## 参考文献

- 1) Dean, J. and Ghemawat, S.: MapReduce: Simplified Data Processing on Large Clusters, *Communication of the ACM*, Vol. 51, No. 1, pp. 107–113 (2008).
- 2) White, T.: *Hadoop*, オライリー・ジャパン (2010).