

# Hadoopにおけるオートパラメータチューニング手法に向けて

## Toward an Automatic Parameter Tuning Method for Hadoop

中下 和則 †      小坂 隆浩 ‡  
Kazunori Nakashita   Takahiro Koita

### 1 はじめに

近年、情報爆発時代に必要とされている新たなインフラとして、大規模分散処理に適した Hadoop[1] が注目されている。Hadoop とは、大規模なデータを効率的に分散して処理できる MapReduce 実装のプラットフォームである。Hadoop にはパフォーマンスに影響を与えるパラメータが 190 個以上存在し、これらのパラメータを適当に特定することが難しい。理由は、パラメータの数が多いため効果のあるパラメータを決めることが困難なこと、またパラメータはアプリケーションの特性によって効果が決まることが挙げられる。仮に、あるアプリケーションに対して適当なパラメータ設定を行い、短時間で、効率よく実行できたとする。そこで違う特性を持つ新たなアプリケーションに対し、同じパラメータを用いても同じような結果を得られるとは限らない。アプリケーションの特性が違うためである。

本研究では、Hadoop のパラメータ値をアプリケーションの実行前にかつ自動的に設定するオートパラメータチューニング手法の実現を目的とする。オートパラメータチューニング手法を目指すにあたって、パラメータを効果のあるパラメータ群と効果のないパラメータ群に分けることとアプリケーションの特性を調査することが重要となる。パラメータチューニングにおいて、アプリケーションの特性が予めわからないため、多数のアプリケーションを調査し、特性をある程度推測することにより、過去の適当なパラメータセットを利用することが可能である。

### 2 Hadoop のパラメータ

ある特定のアプリケーションに対し、Hadoop を適用し、オートパラメータチューニング手法を設計するための予備評価実験を行う。

本章では、有効な Hadoop のパラメータを検討するための方法について議論する。Duke 大学の Shivnath Babu ら [2] は、Hadoop の全パラメータの中の 25 個がパフォーマンスに大きな影響を与えるとする結果を発表した。ここでのパフォーマンスとは、処理に必要な時間である。これにより、全パラメータの中の 25 個を対象とする。以下に Babu らがパフォーマンスに影響を与えたとしたその内のいくつかのパラメータを示す。

- `mapred.map.tasks`
- `mapred.reduced.tasks`
- `io.sort.factor`
- `io.file.buffer.size`

- `mapred.job.shuffle.input.buffer.percent`

例えば、`mapred.map.tasks` や `mapred.reduce.tasks` について、パフォーマンスへの影響を考える。これらは `map` や `reduce` のタスクの数を指定するためのパラメータである。`map`、`reduce` のタスクの数はデフォルトで、それぞれ 2 個と 1 個である。実際にデータ処理を行うマシンの台数を増やしても、これらのパラメータによるタスクの制限が存在するため、用意したマシンの台数分の役割を果たしていない。従って、`mapred.map.tasks` や `mapred.reduced.tasks` の値はパフォーマンスに与える影響が大きい。

また、`io.sort.factor` や `io.file.buffer.size` がパフォーマンスに与える影響について述べる。MapReduce は、各 reducer への入力キーでソートされている。このソートプロセスはシャッフルと呼ばれている。これらのパラメータはこのシャッフル過程で必要になるパラメータである。`map` が出力を始める際、つまり `map` からシャッフルに移る際、その出力過程は複雑である。単純にディスクに書き込まれるわけではなく、メモリへの書き込みのバッファリングなどのために、前もって多少のソートが行われる。シャッフル過程のパラメータの値を変化させることで、パフォーマンスに大きな影響を与える。

### 3 予備実験

予備評価実験の概要について述べる。本研究では、Amazon Web Services[3] の一つである AmazonEC2 を用いた分散データ処理を対象とする。AmazonEC2 ではインスタンスの数を自由に設定することができる。使用するアプリケーションは、Hadoop に添付されているサンプルアプリケーションである。実験の手順を説明する。インスタンスの数は 1,2,4,8 とする。前述のパラメータを実際に変化させ、処理に必要な時間を調べる。実験結果から、これらのパラメータが実際に有効なパラメータであるのかを確認する。

### 4 まとめと今後の課題

本稿では、Hadoop を新しいアプリケーションに適用した際に、短時間で、効率よく実行させるため、オートパラメータチューニング手法の実現を目指すことを示した。今後は、違う特性を持つアプリケーションをいくつか用意し、本稿で得たパラメータに関する知見を元にアプリケーションの特性に基づいた指標を作成することを課題とする。

### 参考文献

- [1] Apache hadoop. <http://hadoop.apache.org/>.
- [2] Shivnath Babu: Towards Automatic Optimization of MapReduce Programs, Proc of the SoCC'10, pp.137- pp.142, 2010.
- [3] Amazon Web Services. <http://aws.amazon.com/>.

† 同志社大学大学院 工学研究科 情報工学専攻

‡ 同志社大学 理工学部 情報システムデザイン学科