

並列データ解析処理基盤の I/O 性能評価に関する一考察

山田 浩之† 合田 和生† 喜連川 優†

† 東京大学生産技術研究所

1 はじめに

企業が保持する情報は大規模化しつつあり、データ解析による大規模データからの価値創出への期待は高い。企業における大規模データ解析は、大規模ストレージに接続されたハイエンドサーバ等の集約システムにより従来から行われてきた。他方、近年では多数のコモディティサーバ、OS、DBMS を一体化したデータウェアハウス・アプライアンスの登場により、並列処理による解析が行われ始めている。データウェアハウス・アプライアンスは、多数のコモディティサーバを利用した並列データベースアーキテクチャを採用し、構造化データ処理に対して優位であることを特徴としている。データ解析技術としては 1990 年代までの並列データベース技術を応用したものに過ぎないと考えられるが、データ解析システムを安価にパッケージ化することで、大規模データ解析が一般企業へと普及する大きなきっかけとなった。

並列処理による大規模データ解析が普及しているなか、並列データ解析処理基盤の Hadoop[1] により、この流れは加速しつつある。Hadoop は Google 社の MapReduce[2] のオープンソース実装であり、非構造化データ処理の柔軟な扱いを特徴としている。また、Yahoo や Facebook では数千ノードでの使用実績があるなど、スケーラビリティに対しても優位性が高い。さらに、Cloudera 等のディストリビューションにより、その導入や管理コストも大きく低下していることから、ウェブ企業でのデータ解析におけるデファクトスタンダードとなりつつある。各社データウェアハウスベンダにおいても、独自のデータ解析処理基盤と Hadoop との連携が進められていることもあり、エンタープライズ分野においてもその活用が進んでいる。

近年の Hadoop による大規模データ解析においては、I/O 帯域を確保するために、各ノードに多数のディスクドライブを搭載する流れが進んでいる [3, 4]。しかし、現状の Hadoop は必ずしも I/O 処理能力が高くない [5]。本稿では、TPC-H データセットを用いた並列解析処理基盤 Hadoop の I/O 性能評価を行うことにより、Hadoop の I/O 処理能力を検証するとともに、その実行アーキ

テクチャについて考察する。

2 実験環境

実験に用いたハードウェア構成を以下に述べる。Hadoop 用サーバとして IBM System x3850 M2 を 1 台と、データ用ストレージとして AMS500 を用いた。当該サーバは、Intel Xeon X7350 を 4 ソケット、メモリを 48GB(DDR2 ECC SDRAM RDIMM)、ローカル HDD として SAS ディスク (10Krpm) を 4 台搭載している。ローカル HDD は OS と Hadoop の一部の管理データの保存用に用いた。また、ストレージは 32 台の FC ディスク (10Krpm) を使用した。2 台の FC ディスクを LU(RAID0) とし、16LU を構築した。サーバ・ストレージ間は 4 本の 4Gbps ファイバチャネルにより接続されている。

次に、実験に用いたソフトウェア構成を以下に述べる。OS は Linux(2.6.18)、Hadoop はバージョン 0.20.2 を使用した。HDFS のブロックサイズは 128MB、I/O のバッファサイズ (io.file.buffer.size) は 32MB とし、MapReduce における Map タスク数を 8 とした。実験においてはソフトウェアのレベルで複数の LU を扱う必要があるが、今回は LVM によるストライピングを用いた。

実験データは TPC-H を用いた。スケールファクタを 100 とし、LINEITEM 表 (約 75GB) と PART 表 (約 2.3GB) を使用した。実験クエリは、以下の 2 つの問い合わせを実施した。

- Job1: 1Map ジョブ (map 処理は空)
- Job2: 1Map ジョブ (map 処理は選択処理)

Job1 は、LINEITEM 表を入力とし、map 処理は何も行わない。つまり、Hadoop のフレームワークを介して LINEITEM 表を全スキャンするのみの処理となる。Job2 は、LINEITEM 表を入力とし、map 処理は shipdate カラムに対する選択処理 (選択率 1%) を行う。選択結果の HDFS への書き込みは行わない。

3 実験と考察

実験は、1LU, 2LU, 4LU, 8LU, 16LU の 5 つのディスク構成に対して実施した。図 1 には、Job1、Job2 に対する、各ディスク構成における読み取りのディスク転送レートを示している。比較として、LVM による論理

An Experimental Study on I/O Performance in Parallel Data Processing Platform

†Hiroyuki YAMADA †Kazuo GODA †Masaru Kitsuregawa

†Institute of Industrial Science, The University of Tokyo

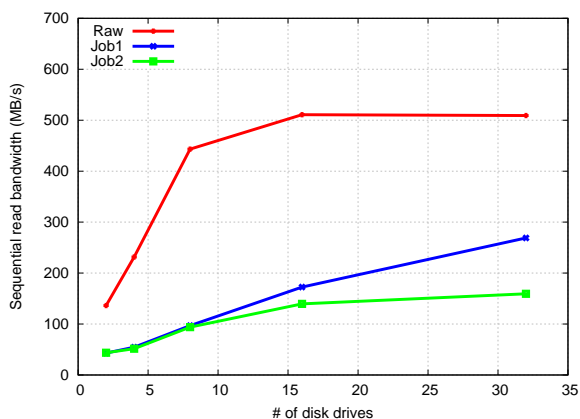


図 1: Job1, Job2 におけるディスク転送レート

ディスクドライブを Raw デバイスとして扱った場合の、単一ストリームによる読み取りディスク転送レートを示している。Raw によるディスク転送レートがディスクドライブ 16 台以上で飽和しているのは、ストレージ側のコントローラの性能によるものである。全体として、Hadoop の実行におけるディスク転送レートは、ディスクドライブ本来の性能よりも大幅に下回った結果となっていることがわかる。これは、Hadoop の I/O 処理に伴う CPU 負荷が高いため、同期的な読み取り I/O では十分な I/O を発行できていないことが原因だと考えられる。実行トレースから得られた CPU 使用率を見ると、ディスクを 2 台使用した場合で平均約 160%^{*}前後となり、I/O 処理量の増加に伴い CPU 使用率も増加し、32 台のディスク使用時に平均約 700%前後となっている。つまり、16 台以降でほぼ性能が飽和しているのは 8 つの Map タスクが各々 1 つの CPU コアを 100% 近くまで使用していることが原因であると考えられる[†]。図 2 に 32 台のディスク使用時の Job2 の CPU 利用率の内訳を示す。Hadoop の I/O 処理の CPU 使用率が高い原因は、汎用的な入力機構や Java による実装、実装の洗練度が低いことなどが考えられるが、詳細な原因の検証は別稿に譲りたい。

実験の結果、現状の Hadoop は必ずしもディスクドライブがもつ I/O 帯域を有効に活用できないという問題を確認した。また、図 1,2 からわかるように、CPU リソース及び I/O 帯域に余裕がある状況下においても性能が飽和してしまうという問題があることが明らかになった。

^{*}16CPU コアを最大限に使用する場合を 1600%とした時の割合となる

[†]Map 処理はブロックごとにプロセスの起動を行うため、各 Map タスクが使用する CPU コアは固定されない。よって、8Map タスクで全体の半分程度の CPU 使用率を使っていることはこのように解釈される。

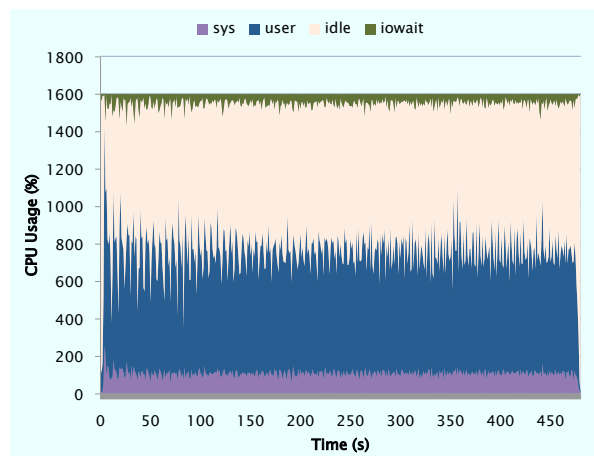


図 2: 32 台のディスクドライブ使用時の Job2 における CPU 使用率の内訳

4 まとめ

近年の並列データ解析処理基盤による大規模データ解析においては、I/O 帯域を確保するために、各ノードに多数のディスクドライブを搭載する流れが進んでいる。本稿では、そのような背景を鑑み、TPC-H データセットを用いた Hadoop の I/O 性能評価を実施した。実験の結果、現状のアーキテクチャでは CPU リソース及び I/O 帯域に空きがあるにも拘わらず性能が飽和してしまうという問題がある。多数の CPU コアを最大限に活用するようにデータ処理アーキテクチャの変更が求められ、今後、当該問題の解決に向けて取り組んでいきたい。

参考文献

- [1] Hadoop, <http://hadoop.apache.org/>.
- [2] Jeffrey Dean, Sanjay Ghemawat, “MapReduce: Simplified Data Processing on Large Clusters” Proc. of OSDI, pp. 137–150, 2004.
- [3] Alex Loddengaard, “Cloudera’s Support Team Shares Some Basic Hardware Recommendations”, <http://www.cloudera.com/blog/2010/03/clouderas-support-team-shares-some-basic-hardware-recommendations/>, 2010.
- [4] “Greenplum Data Computing Appliance”, <http://www.greenplum.com/products/greenplum-hd>.
- [5] “The MapR Distribution for Apache Hadoop”, <http://www.mapr.com/resources/maprforapachehadoopwp>