

ワークフローアプリケーションに対する 計算資源割り当ての最適化

齊藤 貴文^{†1} 千葉 立寛^{†1,†2}
佐藤 仁^{†1} 松岡 聡^{†1,†3}

大規模データを扱うデータインテンシブアプリケーションは近年増加の一途を辿っており、これらの爆発的に増加するデータを分散並列環境で解析的に処理するフレームワークとしてワークフローシステムが注目されている。ワークフローは複数のタスクで構成されており、それらのタスク間での依存関係をユーザが定義して処理が並列計算環境上で実行されるが、ワークフローそれ自体は目的に応じて構成されるため、実行したいワークフローがその実行環境でどのように配置すれば最適に実行可能かを知ることは難しい。また、クラウドのように動的に計算資源をスケールアウト/イン可能なシステムにおいては、ワークフローの特徴や処理に応じてリソースを増減させることでユーザ・システム側の双方において最適な実行が可能となる。未知の環境に対してもワークフローの性能を利用する環境それぞれで最適化させるために、汎用的に適応可能な性能モデルの構築が必要不可欠である。本稿では、ワークフローアプリケーションの性能モデルを構築し、TSUBAME2.0 と Amazon EC2/S3 において天文データ解析を行う Montage に対してモデルを適応させた。小規模なワークフローの実行結果をもとに計算性能や I/O 性能を見積もり、より大規模なワークフローでの性能の予測を行った。その結果、今回構築したモデル化では表現しきれないパラメータの影響、特に中間ファイルがメモリキャッシュにのる場合やメタデータファイルへのアクセスコストにより、正確には性能を見積もることは出来なかった。

Optimization of Resource Allocation for Data-intensive Workflow Applications

TAKAFUMI SAITO,^{†1} TATSUHIRO CHIBA,^{†1,†2}
HITOSHI SATO^{†1} and SATOSHI MATSUOKA^{†1,†3}

Recently data-intensive scientific applications are becoming more popular in astronomy and high-energy physics. Scientific workflow is expected to handle these huge data efficiently. Workflow is constructed by multiple tasks and

programmers describe the dependencies between each tasks as DAG, then the workflow management systems schedule the tasks by following the DAG. Since the structure of DAG varies whenever data set changes, it is difficult to know how many resources we should prepare for the workflow. If we can estimate the performance of workflow in target execution environment, we are able to not only add but also reduce computing resources responded to the feature of workflow, so we need a performance model in order to predict workflow execution time. In this paper we propose performance model for data-intensive workflow, and we estimate the performance by adopting the execution log from test workflow to our performance model. We apply it to Montage workflow in two computing environment; TSUBAME2.0 and Amazon EC2/S3. As a result we discuss the validity of our performance model and also discuss the existing problem.

1. はじめに

高エネルギー物理学やゲノム、天文学などにおける実験や観測によって生成される大量のデータに対して解析的処理や情報抽出を行って活用していく e-サイエンスは“第4のパラダイム”¹⁾として近年注目されている。例えば、CERN の LHC で行われている ATLAS プロジェクト²⁾では、年間数ペタバイト級のオーダーでデータが生成されている。このような大規模データに対していかに効率的かつ並列にデータ処理を行うかが今後ますます重要となってきた。

現在、大規模データに対する並列分散処理を多数のノードにスケールアウトさせていくことで、より高速にデータ処理を行うことが可能なプログラミングモデルとして MapReduce³⁾ や Hadoop が利用されているが、より複雑かつ自由に処理タスクを連結させてデータ処理を次々に適用し結果を導くワークフローシステムに注目が集まっている。ワークフローでは、既存のプログラムを組み合わせているため、単一ノードからスーパーコンピュータにまで簡単に実行環境を移行できるほどプラットフォームの移植性が高く、ワークフローシステムによっては MPI などを用いた場合に必要な煩雑なデータ転送や計算の同期などがユーザから隠蔽され、容易に分散並列処理が可能となる。

^{†1} 東京工業大学

Tokyo Institute of Technology

^{†2} 日本学術振興会特別研究員 DC

Research Fellow of the Japan Society for the Promotion of Science

^{†3} 国立情報学研究所

National Institute of Informatics

様々な環境でワークフローが利用可能になるにつれて、ユーザが実行したい環境にそのワークフローをどのように配置したらよいかはワークフローの形状や実行環境それぞれの性能に強く影響されるため、どの程度のリソースを用意して実行すればワークフローの実行が最適化されるかを知ることは難しい。また、個々のワークフローを構成するタスクも大量のデータを読み書きするタスクなのか、または、長時間の計算が必要であるなど、特徴も異なっているため、それに応じて最適な資源を割り当てることは限りあるリソースを使用する上で重要な課題の一つである。ワークフローは非循環グラフ (DAG) で表現可能であり、ワークフローアプリケーションの性能を予測することが可能となれば、リソースをスケールアウトするだけでなく、必要のないフェーズではリソースを減らすことで、クラウドのようなシステムにおいて、時間的さらには金銭的成本をも削減させることが可能となる。このような性能最適化を行うためには、種々のワークフローに適用可能な汎用的なワークフロー性能モデルが必要不可欠である。

本稿では、ワークフローアプリケーションをタスクフェーズごとの特徴に対して性能モデルを構築し、小規模なワークフローをそれぞれの環境でテストして実行ログを取得し、そのデータをもとにタスクフェーズごとの計算性能や I/O 性能を見積もり、より大規模データを扱うワークフローでの性能を予測する手法を提案した。天文データ解析ワークフローの Montage⁴⁾ を対象として TSUBAME2.0 および Amazon EC2/S3 上での性能予測を行い、実性能との比較を行った。その結果、ある程度の性能予測は可能であったが、モデル化されていないパラメータの影響、特に、メモリキャッシュに中間ファイルが存在する場合やメタデータのアクセスコストが高いことが、性能に大きく影響を与えていることが分かった。

2. 関連研究

グリッドのような並列分散環境におけるワークフロー実行の最適化を行う研究はこれまでも数多くあり、最近ではクラウド上でワークフロー実行最適化に関する研究もいくつか行われている。ワークフローでは、タスク間でのデータはファイルを介して共有されるため、その環境ごとに用意されている共有ファイルシステムを用いることが多いが、データインテンシブワークフローの場合、ファイル I/O に要する時間の割合が多くなるため、Hadoop^{*1}のように分散ファイルシステムと連動してデータとの近接性を考慮した計算タスクの配置によって最適化を行うことが必要である。しかしながら、一般的に広く使われてい

*1 <http://hadoop.apache.org/>

る分散ファイルシステムやクラウドファイルシステムでのワークフロー実行した際に現実的に起こる問題点は明らかになってはいない。

論文⁵⁾では、Amazon EC2 上で様々な種類のワークフローを主にデータ共有するファイルシステム (NFS, GlusterFS^{*2}, PVFS^{*3} and S3) を変化させて性能を観測し、クラウドでのワークフロー実行における問題点を指摘している。単純な性能比較では GlusterFS がこれらのファイルシステムの中では高速であったが、細かい粒度で大量のファイルを作成するワークフローにおいては、クラウドの I/O 性能の低さによりファイル書き込みの時間の増加が顕著に現れることが報告されている。

このような分散ファイルシステム上でのワークフローの I/O 性能を解析するための研究もいくつか存在する。⁶⁾⁷⁾ 論文⁶⁾では、5つのデータインテンシブアプリケーションの Intrigger 上でのファイル I/O に関するログを取得し^{*4}、解析している。メタデータへのアクセスの最適化や大量に生成される中間ファイルをファイルシステムを介さずにタスク間で共有することの必要性について言及されている。

また、グリッド上でのワークフローの性能予測に関する研究もいくつかなされている。論文⁸⁾では、グリッド上で過去に実行されたワークフローの実行ログをデータベース化し、現在実行しようとしているワークフローとの類似性を様々な属性 (問題サイズ、ワークフローの形状、etc) を数値化して得られる距離関数によって定義し、それを用いることで実行時間予測を行う手法を提案している。論文⁹⁾では、グリッド上でのワークフローの実行時に起こりうるオーバーヘッドに着目し、解析的なモデルによりオーバーヘッドを見積もりワークフロー実行の最適化を行う手法を提案している。

3. 提案手法

3.1 仮定

まず始めに本稿で扱うデータインテンシブワークフローについての定義を行う。例として Montage ワークフロー (図 1) を用いて説明する。今、ワークフロー W は、 x 個のタスク $Task_i (0 \leq i \leq x-1)$ で構成されている。各タスクには入力ファイル F_{input} が与えられ、結果としてファイル F_{output} が出力される。図 1 では、9 種類のタスクを経て最終的な天文画像ファイルを作成している。全てのタスクに入力ファイルと出力ファイルが存在

*2 <http://www.gluster.org/>

*3 <http://www.pvfs.org/>

*4 http://www.intriggers.jp/wiki/index.php/Profiler_for_GXP_make

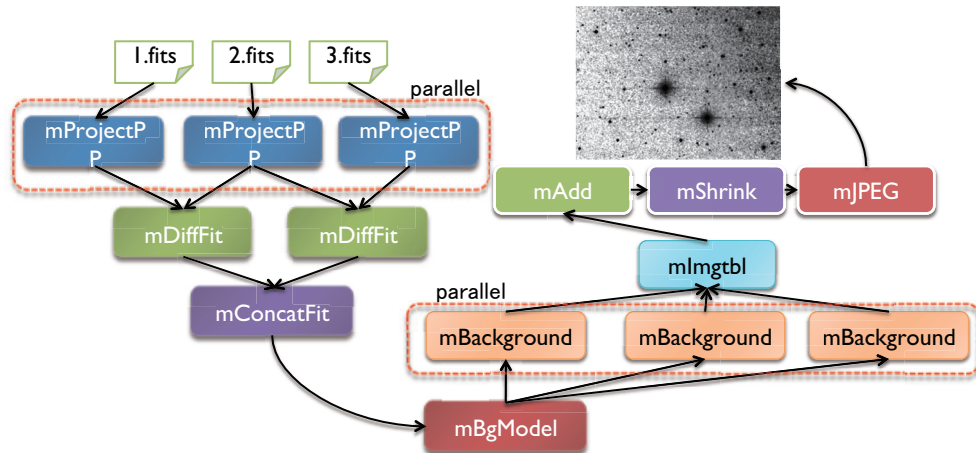


図 1 Montage ワークフローでのタスクの流れ
Fig.1 abstract of Montage workflow

し, $Task_i$ からの出力ファイルは中間ファイルとして共有ストレージ上に保存され, 次の $Task_{i+1}$ はそれを入力ファイルとして読み込み処理を行う.

タスクには, 並列に処理可能な部分 (e.g. mProjectPP) と逐次でしか処理できない部分 (e.g. mJPEG) が存在している. $Task_i$ の実行が完了しない限り $Task_{i+1}$ の実行は開始されないとする. 実行するワークフローは図 1 のようにノードをタスク, エッジをデータの流れたとした DAG で表現することが出来るが, このワークフローの形状は実行前に決定され, 処理の途中でその形状が動的に変化したりしないものとする. また, 初期データ F_{init} に応じてワークフローの形状は一意に定まるものとする.

3.2 ワークフローのモデル化

個々のタスク処理の実行時間を t_i とするとき, ワークフロー全体の実行時間を T_{total} は, 下式 (1) のように表すことが出来る.

$$T_{total} = \sum_{i=1}^n t_i \quad (1)$$

また, 個々のタスクフェーズは, 以下の式のように入出力にかかる I/O 時間とタスク処理での計算時間の和として表現することが出来る.

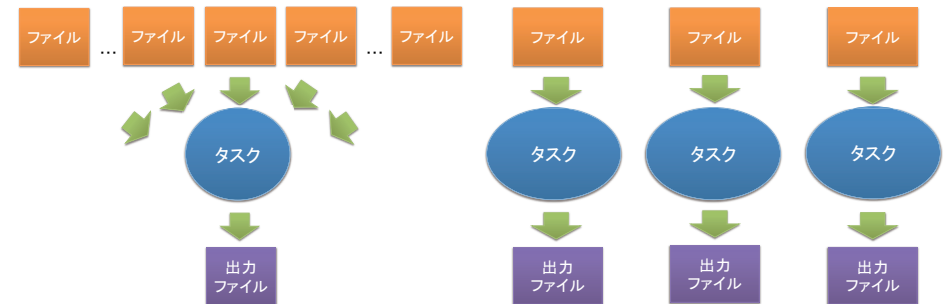


図 2 N 入力 1 出力タイプのタスク
Fig.2 task model for N-input and 1-output

図 3 N 入力 N 出力タイプのタスク
Fig.3 task model for N-input and N-output

$$t_i = T_{read}(N, P, S, J) + T_{calc}(N, P, S, J) + T_{write}(N, P, S, J) \quad (2)$$

N は利用するノード数, P はノードあたりのプロセッサ数, S はそのタスクへの入力ファイルのサイズ, J はタスクに存在するジョブ数を表す. 例えば, 図 1 の mProjectPP では 4 つのジョブが存在することになる.

3.3 タスクのモデル化

次に, 個々のタスクをワークフローの形状や特徴で以下の 4 つのモデルに分類する.

- 1-1 タイプ:** 1 個の入力ファイルに対して 1 個のファイルを出力するタスク. mJPEG などが該当する.
 - N-1 タイプ:** N 個の入力ファイルに対して 1 個の出力ファイルが伴うタスク. (図 2) mConcatFit などに相当する.
 - 1-N タイプ:** 1 個の入力ファイルに対して N 個のファイルを出力するタスク. Montage には存在しない.
 - N-N タイプ:** N 個の入力ファイルに対して N 個の出力ファイルが伴うタスク. タスク内の各ジョブは独立かつ並列に実行される. (図 3) mProjectPP などが該当する.
- 関連研究でも述べられている通り, データインテンシブワークフローでは, ファイルの入出力に要する時間が相対的に大きくなる. ファイルがそのタスクでどのように読み書きされるかに注目し, 共有ファイルシステムに対する影響を考慮するため, このような分類を行なっている. 次に, $Task_i$ を構成する $Read_i, Calc_i, Write_i$ についてのモデル化を行う.

Read/Write 時間

ファイルの Read/Write のスループットを B_{read}, B_{write} とし, 入力出力ファイルのサイ

ズを S_{in}, S_{out} としたとき, $Task_i$ での I/O にかかる時間 $T_{read|write}$ はそれぞれ, 以下のよう
にモデル化することが出来る.

$$T_{Read|Write}(Node, Proc, S_{in|out}, Jobs) = \frac{S_{in|out}}{B_{read|write}} \times \frac{Jobs}{Node \times Proc} \quad (3)$$

$B_{read|write}$ は, 実行するシステムの並列ファイルシステム性能や, 実際に読み込まれるファイルサイズによって異なってくる.

Calc 時間

$Task_i$ でのタスク処理時間を以下のようにモデル化する.

$$T_{calc}(Node, Proc, S_{in}, Jobs) = \nu_i \times \frac{Jobs}{Node \times Proc} \quad (4)$$

ν_i は 1 つのジョブを実行するのに必要な実行時間を表す係数であり, CPU の演算性能とタスクの内容によって決定される.

3.4 性能予測および最適化手法

最後にこのモデル化されたワークフローから性能予測を行う手法について述べる. 今, 実行したい巨大なデータセットをもったワークフロー W_{large} に対して, 同様の処理を行うワークフローであるが小さいかつ異なるデータセットを利用する $W_{small1}, W_{small2}, W_{small3}$ を用意する. W_{small1} は入力データサイズが小さいため, 非常に短い時間で結果を導き出すことが出来るが, このときワークフローの実行と同時にファイルの入出力やタスク実行時間などのログを取得しておく. W_{small1} 実行時に得られたログからその環境上でのファイル I/O のスループットやタスクの計算時間を導き出し, 入力サイズの比に対する実行時間の増加率など環境依存の係数を決定する. それらの値をもとに上記のモデル式 (1)-(4) を用いて性能予測を行う.

4. 評価

本稿ではデータインテンシブワークフローとして Montage⁴⁾ を評価に用いて実験を行った. Montage は天文画像解析アプリケーションで, 2MASS⁵⁾ や SDSS⁶⁾ プロジェクトなどで撮影された大量の全天の画像データからユーザが指定した範囲に対する天体画像をデータベースから取得して生成する. 2MASS プロジェクトなどでは約 24 テラバイトもの画像

表 1 Montage に入力するデータセットの特徴
Table 1 Features of ata set of Montage workflow

データセット	入力ファイルサイズ	中間ファイルサイズ	総ジョブ数
data0	12.3MB	76.5MB	19
data1	93MB	574MB	125
data2	1.2GB	7.58GB	1542
data3	13.6GB	94.2GB	17585

データがあり, 1 枚の天体モンタージュ画像を作成するのに, 多くの画像データや中間ファイルを作成するため, データインテンシブなワークフローアプリケーションの 1 つとしてよく知られている. 表 1 は Montage で用いたデータセットの特徴を示している.

また, 本研究ではワークフローの実行ミドルウェアとして GXP Make⁷⁾ を, さらに実行したワークフローの I/O や 計算時間などのログデータを GXP make 上に実装された logfuse⁴⁾ を用いて取得した.

4.1 TSUBAME2.0 上での予測モデルの評価

まず始めに TSUBAME2.0 でにてノード数やデータサイズを変化させて実験を行った. 使用したノードは Thin ノード (Intel Xeon 2.93 GHz 6 cores * 2, メモリサイズ 96GB) であり, 占有型のノードである. 使用される初期データや実行途中の中間ファイルは全て TSUBAME2.0 上に構築されている並列ファイルシステムの Lustre を使用した.

図 4 は, データセット data2 を用いて実行した際に得られた実行時間と予測時間を比較したグラフである. ノードあたり 12 コアを使っているのので, 8 ノード使用時には合計で 96 コア使用している. 図からも分かる通り, 実測値と予測値に大きな開きが生じてしまい 2 倍程度予測時間が外れてしまう結果となった.

4.2 Amazon EC2/S3 上での予測モデルの評価

次に, Amazon EC2/S3 上でも Montage の実行をノード数を変化させて行った. Amazon EC2/S3 は, Singapore リージョンにある m1.large インスタンス 8 台用意した. m1.large は, 1 ノードあたり 2 コアの CPU, 7.5GB のメモリを有するインスタンスである. 中間ファイルは S3 に対して保存し共有されるが, S3 は POSIX で読み書き出来ず, REST API 経由でしかデータにアクセスすることが出来ない. そのため, 本稿では FUSE ベースで S3 へ

*5 <http://www.ipac.caltech.edu/2mass/>

*6 <http://www.sdss.org/>

*7 <http://www.logos.ic.i.u-tokyo.ac.jp/gxp/>

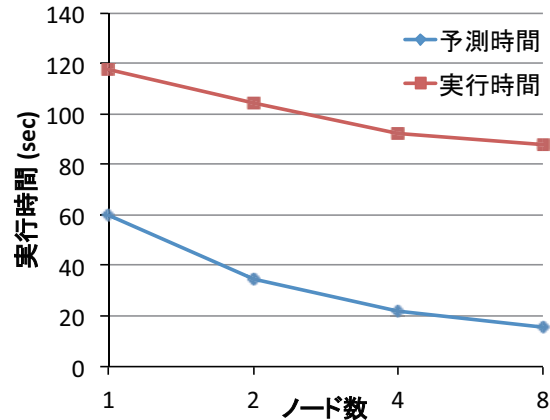


図4 Montage の TSUBAME2.0 での性能予測値と実測値 (data2)

Fig.4 predict and real performance of Montage workflow on TSUBAME2.0(data2)

のアクセスを可能にする s3fs^{*8}を用いて実現した。

図5は、データセット data2 を用いて実行した際の予測時間と実測時間を比較したものである。ノードあたり2コアを使用しているので、8ノード使用時には合計で16コア計算に使われている。ノード数の増加に応じて同時に実行出来るジョブの並列度が向上しているため、性能は向上しているが、TSUBAME2.0での結果と同様、予測値と実測値の差が大きく開く結果となった。

4.3 考察

性能モデルと実測値の差が生じる原因を調べるために、TSUBAME2.0およびAmazon EC2/S3での実行ログの詳細な解析を行った。その結果、いくつかモデル化されていないパラメータの影響を非常に強く受けているため性能のずれが広がっていることが確認された。

まず、図6は、EC2で実行時の個々のタスクフェーズごとでの実行時間がどのようになっているかを示したものである。並列度が向上するにつれて、並列化可能なタスクフェーズについては性能が向上している。さらに、このデータに対してどのタスクフェーズでどの程度の時間がI/Oや計算に費やされているかを示したのが表2である。これから分かる通り、meta-dataへのアクセス時間がかなり要していることが分かる。これは、TSUBAME2.0お

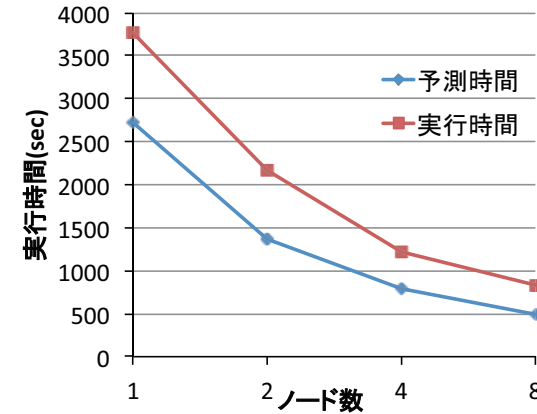


図5 Montage の EC2 での性能予測値と実測値 (data2)

Fig.5 predict and real performance of Montage workflow on Amazon EC2 (data1)

よび EC2/S3 どちらにおいても発生する。TSUBAME2.0において細かいファイルサイズのデータに多数アクセスする際、Lustreのメタデータサーバへのアクセスにおいて性能低下が発生してmeta-dataの時間が増加するものと考えられる。S3においては、FUSE経由でかつREST APIを呼び出してS3にアクセスしているため、このメタデータへのアクセスコストは非常に高くなる傾向にある。この結果から、メタデータへのアクセスに対する性能低下についてより詳しく見積もる必要があることが分かった。

また、図7は、TSUBAME2.0上で実行したmBackgroundタスクにおけるノード数を変化させたときのreadスループットがばらつきを示したものである。mBackgroundタスクフェーズで用いられる中間ファイルは、その前のフェーズで生成された中間ファイルであり、同一のノードで実行される場合、そのファイル自体がノードのメモリキャッシュにのっている可能性が高い。実際に、ログを解析してmBackgroundにおけるreadスループットを検証した結果、node1の場合、確かにほぼすべてのジョブにおいてメモリ上に存在していることを確認した。これは、TSUBAME2.0で用いたノードが96GBのメモリを有していることも影響している。さらに利用するノード数を増加させた場合、キャッシュにのらずにLustreから直接ファイルを読み込むケースも散見されるようになる。これは、ワークフローの実行ミドルウェアが中間ファイルに対するタスクとデータの近接性を考慮していないため、あるノードで実行されたmBackgroundジョブに必要な中間ファイルが他のノードの

*8 <http://code.google.com/p/s3fs/>

表 2 Montage の各タスクフェーズにおける実行時間の比較 (data2)
Table 2 comparison of spent time for each task phase (data2)

	read_time	write_time	write_time	meta data
mProjectPP	35.206	13.126	13.126	5.859795821
mDiffFit	6.574	5.209	5.209	24.0715874
mConcatFit	0.049	0.011	0.011	0.193256122
mBgMondel	0.012	0.00044	0.00044	0.015275258
mBackground	1.195	15.1307	15.1307	5.965266763
mImgtbl	0.0337	0.00053	0.00053	0.008552445
mAdd	1.64	7.5548	7.5548	0.3729
mShrink	0.38	0.0407	0.0407	0.041357337
mJPEG	0.0029	0.00207	0.00207	0.01630968

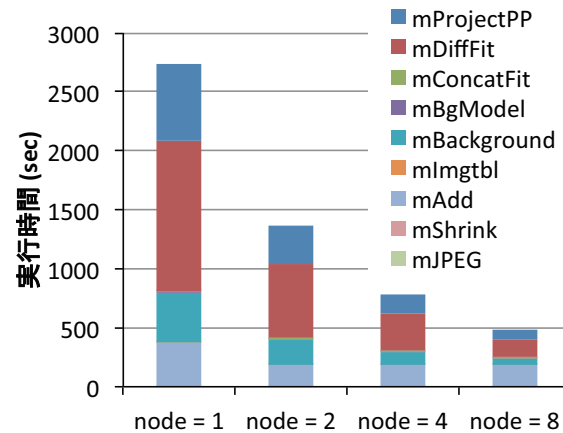


図 6 montage のタスクフェーズごとの実行時間の内訳

Fig.6 predict and real performance of Montage workflow on Amazon EC2 (data2)

キャッシュにのっておりキャッシュを利用できないことに起因する。モデル化に用いたテスト実行は、1 ノードで実行しており、このメモリキャッシュの影響を強くうけてしまっている。そのため、正確な Read のスループットを見積もることが出来ていないことが分かった。EC2/S3 においては、逆に、メモリの搭載量が TSUBAME2.0 に比べて少ないため、メモリキャッシュが性能モデルで利用するパラメータに及ぼす影響が少なかったため、実測値と予測値の差が TSUBAME2.0 での性能ほど出なかったものと考えられる。

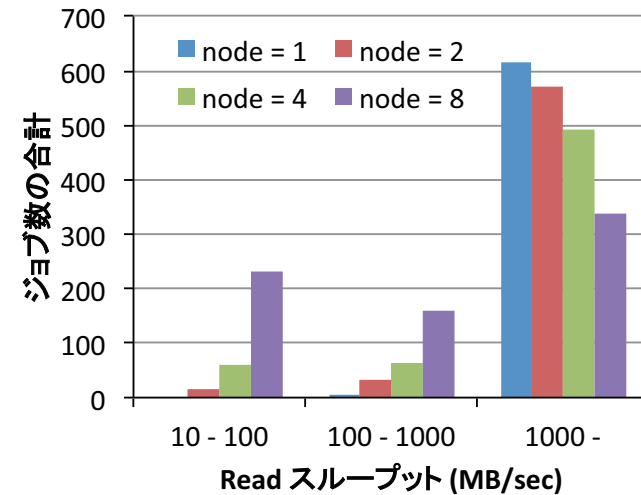


図 7 Montage の Read スループットのばらつき

Fig.7 predict and real performance of Montage workflow on Amazon EC2 (data1)

5. 結論と今後の課題

本稿では、ワークフローアプリケーションをタスクフェーズごとの特徴で分類して、それぞれについてモデル化を行い、それらを組み合わせることで性能を予測する手法を提案した。小規模なデータセットに対する実行のログを取得して環境ごとでのパラメータを決定し、それらをモデル式に当てはめて予測を行っている。データインテンシブワークフローの Montage を TSUBAME2.0 および Amazon EC2/S3 上で実行し、その結果についての考察を行った。実行ログを解析した結果、メモリの容量や共有ファイルシステムの特徴によって大きく予測結果とがずれることがわかった。

計算時間に関しては実測と予測の差は小さいが、I/O 時間に関しては TSUBAME2.0, Amazon EC2/S3 ともに実測と予測の差が大きくなってしまっている。TSUBAME2.0 においては、ノードが備えるメモリが 96GB と非常に大きく、ワークフロー実行途中に生成される中間ファイルがメモリキャッシュに常にのっている状態となり、ノード数が少ない場合は高速にワークフローが実行できる一方、ノード数が多くなるとタスクが分散するためノードのメモリキャッシュにのらないため Lustre ファイルシステムの I/O 性能の影響によって

大きく実行時間が変わることがわかった。一方、Amazon EC2/S3においては、REST 経由で共有ストレージである S3 にアクセスするため、TSUBAME2.0 に比べてメタデータへのアクセスの影響が大きくなっていることが分かった。

今後の課題としては、メモリへのキャッシュヒット率やシステムの I/O スループットなどをより詳細に考慮してモデル式に反映させていく必要がある。また、Montage 以外のアプリケーションにおいてもモデル化が適用可能かを調べる。

謝辞 本研究の一部は、文部科学省科学研究費補助金 (特別研究員奨励費 20・8911, 特定領域研究 18049028) の支援によって行われた。

参 考 文 献

- 1) Gray, J.: *The Fourth Paradigm: Data-Intensive Scientific Discovery*, Microsoft, chapter Jim Gray on eScience: A Transformed Scientific Method, pp.xix–xxxiii (2009).
- 2) A Toroidal LHC ApparatuS Project (ATLAS): <http://atlas.web.cern.ch/>.
- 3) Dean, J. and Ghemawat, S.: MapReduce: Simplified Data Processing on Large Clusters, *Operating Systems Design and Implementation (OSDI '04)*, pp.137–150 (2004).
- 4) Berriman, G.B., Laity, A.C., Good, J.C., Jacob, J.C., Katz, D.S., Deelman, E., Singh, G., Su, M.H. and Prince, T.A.: Montage: The Architecture and Scientific Applications of a National Virtual Observatory Service for Computing Astronomical Image Mosaics, *Proceedings of Earth Sciences Technology Conference* (2006).
- 5) Juve, G., Deelman, E., Vahi, L., Metha, G., Berriman, B., P.Berman, B. and Meachiling, P.: Data Sharing Options for Scientific Workflows on Amazon EC2, *In Proceedings of the 2010 ACM/IEEE International Conference for High Performance Computing (SC '10)*, pp.1–9 (2010).
- 6) Shibata, T., Choi, S. and Taura, K.: File-Access Patterns of Data-Intensive Workflow Applications and their Implications to Distributed Filesystems, *The Third International Workshop on Data Intensive Distributed Computing (DIDC '10)* (2010).
- 7) Dun, N., Taura, K. and Yonezawa, A.: ParaTrac: A Fine-Grained Profiling for Data-Intensive Workflows, *In Proceedings of the 19th ACM International Symposium on High Performance Distributed Computing (HPDC '10)*, pp.37–48 (2010).
- 8) Nadeem, F. and Fahringer, T.: Predicting the Execution Time of Grid Workflow Application through Local Learning, *In Proceedings of the Conference on High Performance Computing Networking, Storage and Analysis (SC '09)*, pp.33:1–33:12 (2009).
- 9) Prodan, R. and Fahringer, T.: Overhead Analysis of Scientific Workflows in Grid

Environments, *IEEE Transactions on Parallel and Distributed Systems*, vol. 19, no. 3, pp., Vol.19, No.3, pp.378–393 (2008).