

HPCI 共用ストレージ機器更新のためのデータ移行

金山 秀智* 原田 浩* 中 誠一郎† 建部 修見‡ 曾田 哲之§ 近藤 晃*

芝野 千尋* 湯川 隆広¶

hidetomo.kaneyama@riken.jp

概要

2017 年 03 月から 2018 年 04 月にかけて HPCI 共用ストレージでは機器更新、機器更新に伴うデータ移行を実施した。移行対象のデータは 2017 年 12 月時点で 9.332PB、8,436,619 ファイルであり、このデータを東京大学更新機器と R-CCS 更新機器に二重化をしながら移行した。データ移行は Gfarm の自動複製機能を用いて HPCI 共用ストレージのサービスを継続したまま実施した。東京大学と R-CCS の旧機器から各機関の更新機器へのデータ移行時の平均移行速度は 4.00GB/sec であった。データ移行中のデータ消失は発生せず、無事 2018 年 03 月 04 日にデータ移行が完了し、HPCI 共用ストレージは 2018 年 04 月から更新機器での二重化運用を開始している。

1 はじめに

HPCI 共用ストレージ^{1) 2)} (以下、HPCI ストレージ) は、革新的ハイパフォーマンスコンピューティングインフラ³⁾ (以下、HPCI) におけるデータ共有基盤として整備された大規模分散ストレージシステムである。HPCI ストレージの老朽化および増大する資源要求に対応するため、我々は 2017 年から 2018 年にかけて機器更新⁴⁾ を実施した。HPCI ストレージには、これまでの膨大な研究成果が蓄積されており、多数の HPCI 課題に対して継続してサービスを提供する必要がある。そこで我々は、HPCI ストレージのファイルシステムである Gfarm⁵⁾ の自動複製機能を活用して、サービスを継続しながら、更新機器へのデータ移行とデータ二重化を実施した。本稿では、更新機器の紹介とデータ移行方法、移行中のデータ保護方法、データ移行速度に関して論ずる。

2 データ移行方法

2.1 データ移行の要求要件

データ移行を移行を実施するにあたって、要求要件は以下の 5 つであることがわかった。

- 2018 年 03 月までにデータ移行を完了させる
- データ移行をサービス無停止で実施する
- データ移行前後でデータ破損・消失が発生していないことを確認する
- データ移行中にデータ破損・消失が発生しても復旧可能とする
- 東京大学と R-CCS にデータを二重化する

要件 (a) は、HPCI ストレージは 2012 年 11 月にサービスを開始し機材の運用年数が 5 年を超過するため、機材更新期間を 2017 年 04 月から 2018 年 03 月と定めたためである。

2017 年 04 月時点で HPCI ストレージには 11.680PB、126,004,228 ファイルのデータが保存され、68 の HPCI 課題が利用している。HPCI ストレージのサービス停止は HPCI ストレージを利用している HPCI 課題の研究活動を阻害するため、要求 (b) の通りデータ移行はサービス無停止で実施

* 理化学研究所

† 東京大学

‡ 筑波大学

§ 株式会社 SRA

¶ 株式会社 創夢

する。

HPCI ストレージは、ネットワークストレージとして東京大学と R-CCS の両機関で運用している。データ転送中のハードウェア障害やネットワーク障害により、転送データの破損が生じる可能性がある。データの破損・消失については、HPCI 課題の研究活動を阻害すると共に、HPCI ストレージの信頼を失う。要求 (c),(d) の通り、データ移行中のデータ破損・消失を検知するとともに、万が一データの破損・消失が発生しても復旧するための方法が必要である。

2018 年 04 月よりデータ二重化運用⁶⁾を開始するため、要求 (e) の通りデータ移行と併せてデータの二重化作業を行う。データ二重化運用は、HPCI ストレージの全データを東京大学と R-CCS に二重化し冗長構成をもたせる運用である。旧機器で運用していた 2017 年 04 月以前は、HPCI ストレージが逼迫しており、提供容量を確保するためデータ二重化運用を実施していなかった。機材更新によりストレージ機器の総容量が東京大学 10PB,R-CCS 10PB の計 20PB から、東京大学 42PB,R-CCS 42PB の計 84PB に増量したため、全データの二重化を実施することができた。二重化を行うことで東京大学または R-CCS のどちらか片方の機関だけで継続運用が可能になり、メンテナンスや計画停電によるサービス停止期間の削減が見込まれる。HPCI ストレージの東京大学拠点は千葉県柏市、R-CCS 拠点は兵庫県神戸市であり、HPCI ストレージ機器は地理的な離れた 2 拠点間に設置されている。このため一方の拠点で災害等影響から機器全停止やネットワーク障害が発生しても、もう一方の拠点で HPCI ストレージサービスを継続することが可能である。

2.2 Gfarm の自動複製機能

2.1 項で求められる要件を満たすためには、Gfarm の自動複製機能⁷⁾を用いてデータ移行を行うことが最適であると判断した。

Gfarm の自動複製機能は以下の特徴を持つ。

- サービス中に自動複製が行われる

- データの複製先をホストまたはホストグループ*¹ 単位で指定可能
- データの複製数をホストまたはホストグループ単位で指定可能
- 不要な複製データの自動削除が可能、また自動削除機能を有効/無効に変更可能
- 自動複製は並列実行される、また並列数を調整が可能
- 自動複製先のホストは Gfarm により効率良く選択される

自動複製はサービス中に実行されるため、要求 (b) を満たす。自動複製機能はデータの複製先・複製数の指定が可能であり、要求 (e) も実現できる。複製データの自動削除機能を無効にすることで、旧機器に保存されている移行元データは削除されない。データ移行中にファイル・破損消失が発生しても、旧機器から復旧が可能であり、要求要件 (d) を満たす。自動複製機能では並列数が調整可能であり、複製先として空き容量が多く・負荷が低いファイルサーバが自動選択されるため、ネットワーク帯域を埋めて効率よく複製を作成することができる。このため、要求 (a) 通り 2018 年 03 月までにデータ移行が完了できると判断した。

2.3 データ保護

2.1 項の要求 (d) 通り、データ破損・消失を検知する必要がある。Gfarm では一貫性チェック機能を持ち、データ破損・消失を検知することができる。Gfarm の一貫性チェックは以下の通り行われる。

- ファイルサーバへのデータ書き込み時の一貫性チェック
 - (1) ファイルサーバにデータが書き込まれる
 - (2) 書き込まれたデータのチェックサムを計算
 - (3) メタデータのチェックサムと (2) で計算したチェックサムを比較
 - (4) チェックサムが異なる場合は、ログにアラートを出力

*¹ <http://oss-tsukuba.org/gfarm/share/doc/gfarm/html/ja/ref/man0/gfhostgroup.1.html>

- Write Verify 機能 *² による一貫性チェック
 - (1) ファイルサーバへのデータが書き込まれる
 - (2) (1) の書き込み後一定時間経過後 *³ に書き込まれたデータのチェックサムを再計算する
 - (3) メタデータのチェックサムと (2) で計算したチェックサムを比較
 - (4) チェックサムが異なる場合は、ログにアラートを出力
- レプリカチェック *⁴ による一貫性チェック
 - (1) レプリカチェックが実行される
 - (2) ファイルがファイルサーバに存在することを確認
 - (3) メタデータのファイルサイズとファイルサーバ上のファイルサイズの一致を確認
 - (4) (2) でファイルが存在しない場合、または (3) でファイルサイズが異なる場合はログにアラートを出力

2.4 データ移行方法

データ移行は図1の通り、移行データの複製先と複製数を東京大学とR-CCSの更新機材にそれぞれ1つ指定することで実施する。自動複製機能により旧機器に保存されている移行元データが、東京大学とR-CCSの更新機器へ1つずつ転送される。移行後は東京大学とR-CCS更新機材でデータが二重化される。削除機能を無効に設定することで、データ移行後も旧機器から移行元データは削除されない。データの破損チェックのため、平日に必ず1回自動複製を中断し、全データに対してレプリカチェックによるデータ一貫性チェックを実施した。

3 システム構成概要

3.1 東京大学 システム構成概要

東京大学のシステム構成概要は図2、更新機器の構成は表3.1の通りである。東京大学の

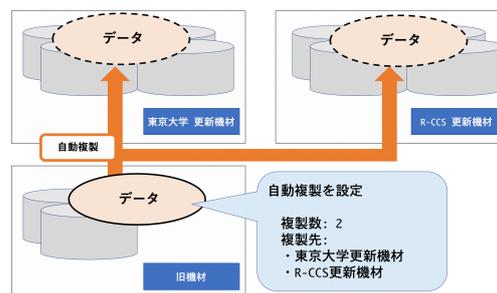


図1 自動複製によるデータ移行

ストレージ装置 DDN SFA14KXE の総容量は42PBである。ファイルサーバはストレージ装置 DDN SFA14KXE 1セットにつき4台で、DDN SFA14KXE 上の仮想マシンとして動作している。ファイルサーバはそれぞれが上位スイッチと Infiniband 4xEDR(100Gbps) で接続している。

東京大学ではデータ移行用の一時データ保存用機器を用意している。データ移行の前段階として、旧機材から一時データ保存用機器にデータ移行を実施した。一時データ保存用機器の構成は3.1の通りある。一次データ保存用ファイルサーバは2台で、それぞれが上位スイッチと Infiniband 2xFDR(56Gbps) で接続している。東京大学一時データ保存用機器と東京大学更新機器間のデータ転送は112Gbpsに律速される。

表1 東京大学の更新機器情報

機器	台数	機種
ストレージ装置	7セット	DDN SFA14KXE
ファイルサーバ	28台	DDN SFA14KXE(仮想マシン)

表2 東京大学の一時データ保存用機器情報

ストレージ装置	2セット	DDN SFA 7700X
ファイルサーバ	2台	HP Proliant DL380 Gen9

3.2 R-CCS のシステム構成概要

R-CCS のシステム構成概要は図3、更新機器情報は表3.2の通りである。R-CCSの更新ファイル

*² <https://sourceforge.net/p/gfarm/mailman/message/34771054/>

*³ HPCI ストレージは6時間後に再計算する

*⁴ <https://www.soum.co.jp/misc/gfarm/3/>

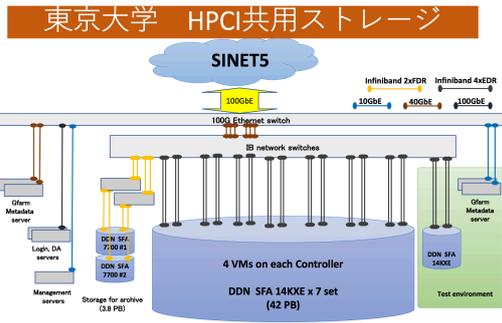


図2 東京大学 システム構成概要

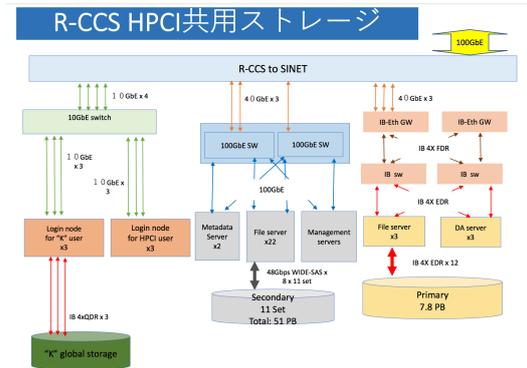


図3 R-CCS 更新機器

サーバは18台で、100GbEスイッチと各100GbE x2(bonding)で接続されている。ストレージ装置の総容量は東京大学と同様に42PBである。各ストレージ装置とファイルサーバとは192Gbpsで接続されている。ファイルサーバ1台からストレージセットへ58並列でファイル転送を行った際の、読み出し転送速度は18GB/secである。同様に書き込み転送速度は22GB/secである。100GbEスイッチから上位スイッチへは40GbE x 3のLAGで接続されており、上位スイッチとSINETへは100Gbpsで接続されている。

旧機器のファイルサーバは15台で、各10GbEでファイルサーバ用スイッチと接続されている。旧ファイルサーバ用スイッチと100GbEスイッチ間は10GbE x 4のLAGで接続されている箇所があり、R-CCS内で旧機器と更新機器間のデータ転送帯域は40Gbpsに律速される。R-CCS旧ファイルサーバ用スイッチと上位スイッチとは20Gbpsで接続されているため、R-CCS旧機器と東京大学機器間の転送帯域は20Gbpsに律速される。

表3 R-CCSの更新機器

機器	台数	機種
ストレージ装置	9セット	CMS HyperSTOR Flex
ファイルサーバ	18台	DELL R730

3.3 東京大学とR-CCS間のネットワーク構成

東京大学とR-CCS間のデータ転送はSINET5⁸⁾を利用する。

東京大学では2018年02月22日までは、SINET5と東京大学機器の経路上に10GbE x 4のLAGが存在しており、R-CCS更新機器と東京大学機器間のデータ転送帯域は40Gbpsに律速されていた。2018年02月23日に東京大学機器からSINET5への接続は100Gbpsに更新された。

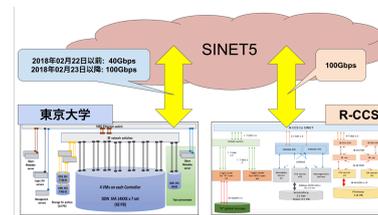


図4 東京大学とR-CCS間のネットワーク構成

3.4 各機器のデータ移行の転送速度

2018年02月22日までの各機関間の最大転送帯域は表4の通り1である。2018年02月23日以降の最大転送帯域は表5の通りであり、東京大学とR-CCS更新機器間のデータ転送帯域が広がった。

東京大学またはR-CCSとSINETの通信には、HPCIストレージへのユーザアクセスや、各機関の他のサービスやスパコンの利用があるため、実際に東京大学とR-CCS間でデータ移行を行う際に利用できる転送帯域は、最大転送帯域以下となる。

表 4 最大データ転送帯域 (~2018/02/22)

	東京大学 更新機器	東京大学 一時保存用機器	R-CCS 更新機器	R-CCS 旧機器
東京大学 更新機器	-	112 Gbps	40 Gbps	20 Gbps
東京大学 一時保存用機器	112 Gbps	-	40 Gbps	20 Gbps
R-CCS 更新機器	40 Gbps	40 Gbps	-	40 Gbps
R-CCS 旧機器	20 Gbps	20 Gbps	40 Gbps	-

表 5 最大データ転送帯域 (2018/02/23~)

	東京大学 更新機器	東京大学 一時保存用機器	R-CCS 更新機器	R-CCS 旧機器
東京大学 更新機器	-	112 Gbps	120 Gbps	20 Gbps
東京大学 一時保存用機器	112 Gbps	-	120 Gbps	20 Gbps
R-CCS 更新機器	120 Gbps	120 Gbps	-	40 Gbps
R-CCS 旧機器	20 Gbps	20 Gbps	40 Gbps	-

4 移行量

HPCI ストレージには 2017 年 12 月時点で 9.332PB、84,336,619 ファイルのデータが保存されていた。データ移行では東京大学と R-CCS のデータ二重化を併せて実施する。旧機材に保存されているデータを新機材に二重化する必要があるため、転送量は 2 倍となり、18.664TB、168,673,238 ファイルとなる。

2017 年 04 月時点は、ユーザが複製先・複製数を任意に設定できた。二重化されていないデータや、二重化されていても東京大学または R-CCS のどちらか一方に保存されているデータが存在した。表 6 の通り、「東京大学のみ」と「R-CCS のみ」に保存されているデータは合計 4.873PB であり、東京大学と R-CCS 間での転送が必要である。東京大学と R-CCS 間は 4 章の通り、東京大学内または R-CCS 内に比べてネットワーク帯域が狭く、R-CCS 旧機器から東京大学更新機材へのデータ転送は 20Gbps に律速される。

表 6 保存先別移行量

移行元データの保存先	容量 (PB)	ファイル数
東京大学のみ	2.029	29,458,968
R-CCS のみ	2.844	38,602,962
東京大学と R-CCS	4.459	16,274,689
合計	9.332	84,336,619

5 データ移行の段階的实施

データ移行は以下の理由から表 7 の通り 4 段階に分けて実施した。東京大学と R-CCS で更新機器の納入時期や運用環境への導入時期が異なったため、段階を分けてデータ移行を実施した。Gfarm は運用中のファイルサーバの追加・削除が可能であるため、更新機器は導入可能になったものから運用環境に追加しデータ移行対象に加えた。旧機器はデータ移行が完了したものから取り外した。

移行段階 1 は、東京大学旧機器の撤去のため、東京大学旧機器にのみ保存されているデータを退避させるため実施したデータ移行である。更新機器の導入前は、東京大学一時データ保存用機器へデータを移行し、東京大学更新機器の導入後は更新機器と一時データ保存機器へデータを移行した。

移行段階 2 は R-CCS 更新機器へのデータ移行である。

移行段階 3 は R-CCS 更新機器に加えて、東京大学更新機器へのデータ移行を実施した。移行段階 3 では東京大学と R-CCS 間のデータ二重化を優先したため、東京大学一時データ保存用機器に保存しているデータについては、東京大学更新機材へデータ移行は行わなかった。

移行段階 4 では東京大学一時データ保存用機器に保存されているデータを東京大学更新機器へ移行した。

6 移行結果

データ移行は 2018 年 03 月 04 日に完了した。データ移行中のデータ消失は発生せず、全移行データを無事更新機材へ移行することができた。各移行

表7 データ移行段階

段階	移行元機器	移行先機器
1	東京大学 旧機器	東京大学 更新機器 一時データ保存用機器
2	東京大学機器 R-CCS 旧機器	R-CCS 更新機器
3	東京大学機器 R-CCS 機器	東京大学 更新機器 R-CCS 更新機器
4	東京大学 一時データ保存用機器 R-CCS 機器	東京大学 更新機器

段階の移行期間・転送量・平均転送速度は表8の通りである。図5は移行段階2,3の移行進捗である。

表8 各移行段階の移行結果

段階	移行期間	転送量 (PB)	平均転送速度 (GB/sec)
1	2016/12/05 - 2017/12/18	4.742	-
2	2017/12/21 - 2018/01/17	5.377	2.56
3	2018/01/18 - 2018/02/17	9.422	4.00
4	2018/02/27 - 2018/03/04	1.959	7.47

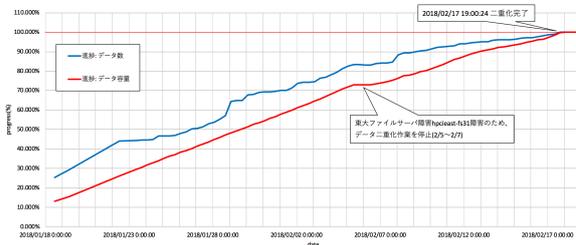


図5 移行段階2,3の移行進捗

6.1 移行段階1のデータ移行結果

移行段階1は、旧機材の撤去と更新機材の導入スケジュールに合わせて、2016年12月05日に開始し2017年12月18日に完了した。移行段階1では移行段階2~4の準備として、データ移行方法の確認や並列数の調整やデータ移行中の負荷の確認を実施し、機器の撤去スケジュールに合わせて移行元の東京大学の旧機器を限定しながら実行した。表8に以降段階1の平均転送速度を記載していないのは、移行対象とした移行元機器の台数や並列数によってバラツキが大きかったためである。

6.2 移行段階2のデータ移行結果

移行段階2のデータ移行は、平均転送速度2.56GB/secで2017年12月21日~2018年01月17日まで実施した。2018年01月18日に、東京大学の更新機器へのデータ移行の準備ができたので、移行段階3に進んだ。

図6は移行段階2のR-CCS更新機器への1日毎の移行容量の累計である。2018年01月03日から05日までの移行容量の増加がないのは、データ保護のために準備していたレプリカチェック自動停止スクリプトが実行されてレプリカチェックが中断していたためである。図7は、Zabbixで監視したR-CCSの各更新ファイルサーバへのネットワーク転送帯域の積み上げグラフである。ファイルサーバへの転送帯域であるため、データ移行の移行転送帯域に加えて、ユーザの書き込み等の帯域も加わっている。

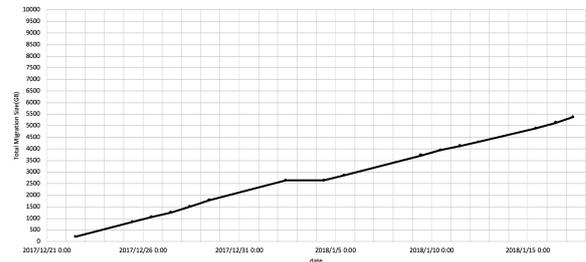


図6 移行段階2のR-CCS更新機器への移行容量の累計

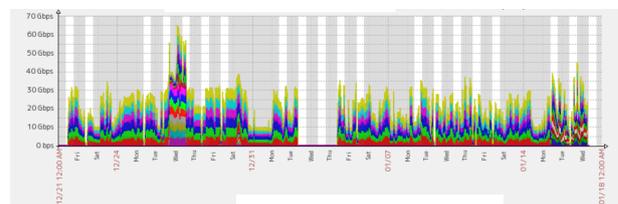


図7 移行段階2のR-CCS更新ファイルサーバへの転送帯域

6.3 移行段階3のデータ移行結果

移行段階3のデータ移行は、平均転送速度4.00GB/secで2018年01月18日~02月17日

で実施した。移行段階 3 のデータ移行完了後は、R-CCS 更新機器へ全データのデータ移行が完了した。全データが東京大学の更新機器と一時データ保存用機器に 1 つ、R-CCS の更新機器に 1 つの計 2 つ保存されたことで、データの二重化が完了した。図 8 は移行段階 3 の東京大学更新機器と R-CCS 更新機器の移行容量の累計である。図 9 は、Zabbix で監視した R-CCS の各更新ファイルサーバへのネットワーク転送帯域の積み上げグラフある。

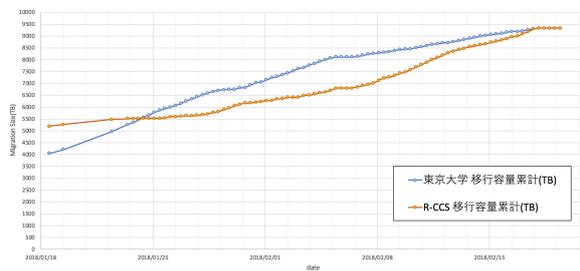


図 8 移行段階 3 の東京大学・R-CCS 更新機器への移行容量の累計

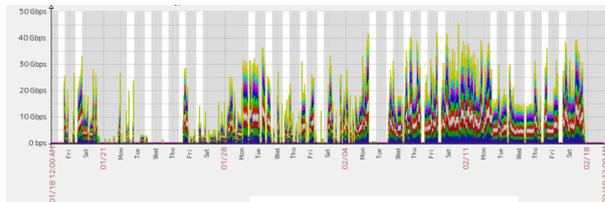


図 9 移行段階 3 の R-CCS 更新機器への転送帯域

6.4 移行段階 4 のデータ移行結果

移行段階 4 のデータ移行は、平均移行速度 7.47GB/sec で 2018 年 02 月 27 日～03 月 04 日の 5 日間で完了した。データ移行完了後は、東京大学更新機器へ全データのデータ移行が完了した。

移行段階 4 は、東京大学更新機器と SINET 間が 100Gbps で接続されたため、R-CCS 更新機器から東京大学更新機器へのデータ転送帯域が 100Gbps となった。図 10 は、Zabbix で監視した R-CCS の各更新ファイルサーバから東京大学更新機器へのネットワーク転送帯域の積み上げグラフある。データ移行中の R-CCS 更新機器から東京大学更新機器

への最大ネットワーク帯域は 42Gbps であった。

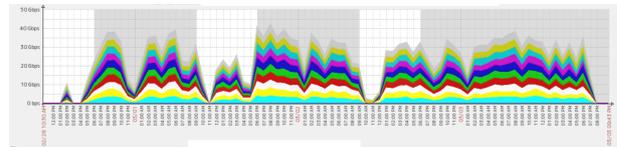


図 10 移行段階 4 の R-CCS 更新機器から東京大学更新機器への転送帯域

7 考察

移行段階 2 では、(1)R-CCS 旧機材から R-CCS 更新機器へのデータ移行と、(2) 東京大学機器から R-CCS 更新機器へのデータ移行が実行された。3.4 項の通り最大転送帯域は (1) と (2) ともに 40Gbps である。(1) と (2) の転送が同時に行われた際の転送帯域は 80Gbps なので、移行段階 2 のデータ転送速度は最大 10GB/sec であるが、実際のデータ移行での平均データ転送速度は 2.56GB/sec であった。図 7 の転送帯域を見ると、60Gbps を超過するような転送帯域でデータ移行が行われている期間と 10Gbps 程度で転送している期間があり、転送帯域は安定していなかった。60Gbps を超過するような転送帯域でデータ移行が行われている期間は、期待通り (1) と (2) の両方のデータ転送が行われたものと考えられる。10Gbps の期間は Gfarm の自動複製の特性により、転送帯域が狭まっていたものと推測している。Gfarm の自動複製は inode 番号の順に実行される。東京大学または R-CCS のどちらか一方に保存されているデータが連続した inode 番号に割り振られている場合、データ転送は東京大学または R-CCS のどちらか一方からしか行われていなかったものと推測している。小さいファイルサイズのデータ転送連続で実行されている場合も、東京大学と R-CCS 間の RTT^{*5}や Gfarm の読み込み・書き出し処理がデータ転送の時間に比べ大きく比重を持つためデータ転送速度が悪化すると考えられる。

東京大学と R-CCS 更新機器へのデータの二重化

*5 東京大学と R-CCS 間の RTT は約 10msec

が完了し、2018年04月より二重化運用を開始した。二重化運用では東京大学またはR-CCSの一方の拠点が停止してもサービスを継続することが可能である。大規模障害やメンテナンス、計画停電等で、東京大学またはR-CCSのどちらか一方へのアクセスができなくなった場合でも、片方の機関でサービスを継続することができる。『次期HPCI共用ストレージにおけるサイト間データ冗長運用によるサービス継続性向上策⁶⁾』によると東京大学とR-CCS間のデータ二重化により「サービス停止サーバ時間積の81%、サービス停止回数の87%が改善され、サービス継続が可能であると見込まれる」(原田浩,2016,p4)とされる。2018年04月以降のHPCIストレージのサービス提供時間は、2018年03月以前に比べ向上するものと見込んでいる。今後は、二重化運用によりサービス提供時間が増加することを検証する。

8 まとめ

データ移行は2018年03月04日に完了し、HPCI共用ストレージは2018年04月より更新機器での運用とデータ二重化運用を開始した。

データ移行は、Gfarmの自動複製機能を利用して実施した。移行データは9.332PB、84,336,619ファイルで、東京大学とR-CCSの更新機器にそれぞれ1つずつ二重化されるようにして転送した。

データ移行は、更新機器導入のスケジュールに合わせて4段階に分けて実施し、このうち東京大学更新機器とR-CCS更新機器へのデータ移行の平均転送速度は4.00GB/secであった。

参考文献

- 1) 實本英之, 建部修見, 佐藤仁, 石川裕: 広域分散環境を提供するHPCIシステムソフトウェア基盤の設計概要と共有ストレージ構築, 研究報告ハイパフォーマンスコンピューティング(HPC), Vol. 2011-HPC-130, No. 67, 情報処理学会, pp. 1-6 (2011)
- 2) 原田 浩, 建部 修見, 平川 学, 藤本 大 輔, 蛭原 純, 實本 英之, 宮崎 洋, 佐島 浩之, HPCI

共用ストレージの構築と運用、大学ICT推進協議会2013年次大会論文集、HPCIテクノロジーT3G-5、pp422-427、2013.

- 3) 革新的コンピュータインフラストラクチャ http://www.mext.go.jp/a_menu/kaihatu/jouhou/hpci/1307375.htm
- 4) 平成29年度機材更新のお知らせ <https://www.hpci-office.jp/info/pages/viewpage.action?pageId=39230294>
- 5) Osamu Tatebe, Kohei Hiraga, Noriyuki Soda, Gfarm Grid File System, New Generation Computing, Vol28, pp257-275, 2010. 新的コンピュータインフラストラクチャ
- 6) 原田浩, 建部修見, 埜敏博, 中誠一郎, 平川学, 金山 秀智, 近藤 晃, 次期HPCI共用ストレージにおけるサイト間データ冗長運用によるサービス継続性向上策、AXIES2016年次大会HPC、2016.
- 7) 建部修見 Gfarm ファイルシステムの概要と実装 Gfarm Workshop 2015 資料 <http://oss-tsukuba.org/wpcontent/uploads/2015/09/Gfarm-overview.pdf>
- 8) SINET5 https://www.nii.ac.jp/userdata/openforum/PDF/2015/1_setsumeikai2015_sinet5_20151023.pdf