

Gfarm Symposium 2020
2020年10月9日 @秋葉原

Gfarmファイルシステムの 概要と最新機能

建部修見
筑波大学

Gfarmファイルシステム

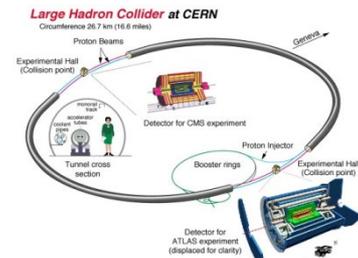


- オープンソース広域分散ファイルシステム
 - <http://oss-tsukuba.org/software/gfarm/>
- サポート
 - NPO法人つくばOSS技術支援センター(日本他)
 - Libre Solutions Pty Ltd(オーストラリア)
- 特徴
 - 性能・容量がスケールアウト
 - データアクセス局所性、ファイル複製
 - 無停止で拡張、更新可能
 - 単一障害点なし
 - 複製数維持機能、ホットスタンバイMDSサーバ
 - データ完全性を保証しサイレントデータ損傷も対応可

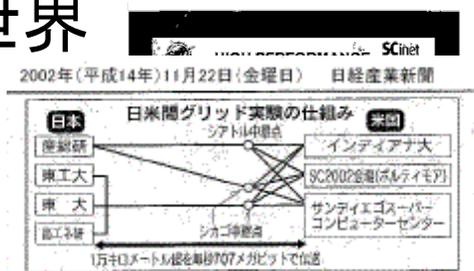
ossTsukuba
oss-tsukuba.org



Congratulation 20th anniversary!



- 2000年6月 電総研と高エネ研のGRIDミーティング
- LHC実験の要求に応えるため日本独自研究開発開始
- 2002年11月 日米間ファイル複製実験で世界最高の741Mbpsを達成!
- 2003年11月



2002年(平成14年)11月22日(金曜日) 日経産業新聞

複数結ぶ「グリッドコンピューティング」

日米間通信で最速記録

毎秒707メガバイト達成 産総研など

産業技術総合研究所(産総研)の国際学会「SC2002」で発表された日米間の通信実験で、毎秒707メガバイトという最速記録を達成した。この実験は、産総研とインディアナ大学が共同で行った。実験は、東京工業大学が中心となり、産総研、東工大、東大、高エネ研が参加した。実験は、サンディエゴスーパーコンピュータセンターと産総研のネットワークを介して行われた。実験の結果、毎秒707メガバイトという最速記録を達成した。この記録は、産総研とインディアナ大学の共同研究によるものである。実験は、産総研とインディアナ大学の共同研究によるものである。実験は、産総研とインディアナ大学の共同研究によるものである。

日刊工業新聞 2002年(平成14年)11月22日(金曜日)

産総研 日米間、毎秒707メガバイト達成

グリッドコンピューティングシステム

データ処理、パソコン利用

産業技術総合研究所は、21日、パソコンクラスターを利用した日米間を結ぶグリッドコンピューティングシステムで、これまでの記録を大幅に上回る約707メガバイト(707Mbps)の高速度で、超大規模データを処理する「SC2002」に成功したと発表された。

目標の800Mbpsには達しなかったが、1万キロ離れた日米間で、西米の7拠点に分散しているサーバと接続するコンピュータに関する研究は、国際会議「SC2002」で発表された。産総研が開発したグリッドコンピューティングシステム「データファーム」と機材を持ち込んできた。産総研の(高エネ研)千加速器研究機、東京工業大学など日米6機関が参加し、米サンディエゴスーパーコンピュータセンターなどと接続した。

産総研は、この実験で、グリッドコンピューティングの技術を実用化し、企業や官民機関、生命科学研究などの分野で、一公開し、民間での普及を図る。

産経新聞 2002年(平成14年)12月1日(日曜日)

産総研 超大規模データ転送に成功

グリッドコンピューティング 世界初の実証

独立行政法人「産業技術総合研究所」(産総研)が茨城県つくば市(つくば市)など、日米7カ所を拠点として日米十カ所のパソコンを国際的に結び、通常は大型コンピュータで行う「ハブ(CDR-ROOM)約1万台(相当)もの超大規模データ転送に成功した。

実験では、米サンディエゴスーパーコンピュータセンターや東京工業大学など日米の拠点に、それぞれ数万台のパソコンを結ぶ「グリッドコンピューティング」を実証した。実験は、産総研とインディアナ大学の共同研究によるものである。実験は、産総研とインディアナ大学の共同研究によるものである。

その結果、CDR-ROOM一枚を0.1秒で読み込める約707メガバイトの高速度処理も実現。産総研では「さらに大規模・高速化も可能。今後は欧州へも拡大して実証実験を進めていく」としている。

(伊藤壽一郎)

2.0.0公開!
m 2.7.17公開!!!

HPCI共用ストレージ

- 大学情報基盤センターをはじめ全国からマウント可能な共有ファイルシステム(～100PB)
- スパコン間のデータ共有、共有データ格納



西拠点 (R-CCS)



東拠点 (東京大)

最新機能・状況紹介

主なリリース

日付	version	新機能、更新機能
2020/9/17	2.7.17	• MTセーフ、ROFS機能、FO強化
2019/11/30	2.7.16	• Githubへの移行！
2019/10/24	2.7.15	• Gfarmbb status, IB GRH対応
2019/9/10	2.7.14	• Gfarm/BBバーストバッファ
2016/12/8	2.7.0	• InfiniBand RDMAサポート • ディレクトリクォータ
2016/1/16	2.6.8	• 書込後ベリファイ

暗号化ファイルシステム

- <https://github.com/oss-tsukuba/gfarm/blob/master/doc/encfs.ja.md>

Gfarm暗号化ファイルシステム

EncFS(*)を用いることにより、Gfarmファイルシステムに暗号化されたデータを格納することができます。

(*) <https://github.com/vgough/encfs/blob/master/encfs/encfs.pod>

インストール

EncFSをインストールします

```
# yum install encfs
```

使い方

1. Gfarmファイルシステムをマウントする

```
$ gfarm2fs /tmp/gfarm
```

この例では、/tmp/gfarmにGfarmファイルシステムをマウントしています。

2. 暗号化ファイルシステムを作成しマウントする

ROFS – Gfarm Read only FS

- ファイルシステムを完全にread onlyに
% `gfstatus -Mm 'read_only enable'`
- Zabbixフェイルオーバースクリプトでは、split brainの可能性が残る場合にROでフェイルオーバー
 - 確認が取れ次第read onlyを解除

https://github.com/oss-tsukuba/gfarm

The screenshot shows the GitHub repository page for `oss-tsukuba/gfarm`. The browser address bar displays `github.com/oss-tsukuba/gfarm`. The repository name `oss-tsukuba / gfarm` is shown at the top, along with interaction buttons: `Unwatch` (4), `Star` (4), and `Fork`. Below this is a navigation bar with tabs for `Code`, `Issues` (138), `Pull requests` (0), `Actions`, `Projects` (0), `Wiki`, `Security`, `Insights`, and `Settings`.

The repository description reads: "distributed file system for large-scale cluster computing and wide-area data sharing. provides fine-grained replica location control." Below the description is a `Manage topics` link.

A statistics bar shows: `5,259` commits, `30` branches, `0` packages, `191` releases, and `10` contributors. A `View license` link is also present.

At the bottom of the statistics bar, there are buttons for `Branch: master`, `New pull request`, `Create new file`, `Upload files`, `Find file`, and a prominent green `Clone or download` button.

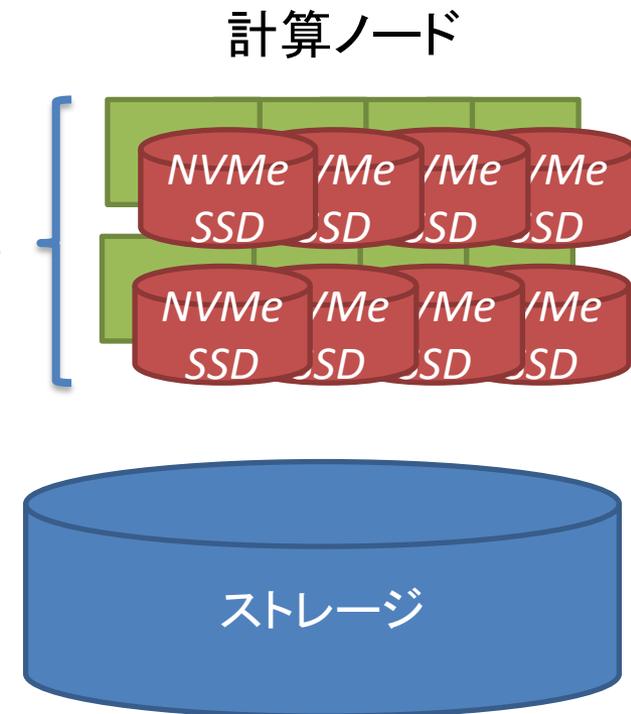
The commit history table is as follows:

Commit	Message	Time
<code>otatebe</code>	<code>libpq-devel for RHEL8</code>	Latest commit <code>cf6140d</code> 6 days ago
<code>bench</code>	merge <code>r10991</code> from the <code>replica_check_2</code> branch (via <code>r10993</code> from the 2.7...	12 months ago
<code>debian</code>	remove <code>gftool</code>	4 years ago
<code>doc</code>	regen for #1081 - <code>gfmv(1)</code> between different directory quota settings	16 days ago
<code>gftool</code>	fix data race	6 days ago
<code>include/gfarm</code>	assign log message numbers	2 months ago
<code>lib</code>	fix out-of-bound access	6 days ago
<code>linux</code>	assign log message numbers in linux	3 years ago

Gfarm/BBノバーストバッファ [建部 2020]

- ノードローカルNVMe SSD等高速ストレージによる一時的な分散ファイルシステム
- アクセス性能の向上
 - ファイルディスクリプタパッシングによるgfsdを経由しない直接アクセス
 - RDMAアクセス
- メタデータ性能の向上
 - メタデータの永続性、冗長性なし
 - ジャーナル書込み、バックエンドDB、スレーブgfsdのオーバヘッドの削減
- 冗長性オーバヘッド削減
 - ファイル複製によるデータの冗長性なし
- ファイルシステム構築、撤去の高速化

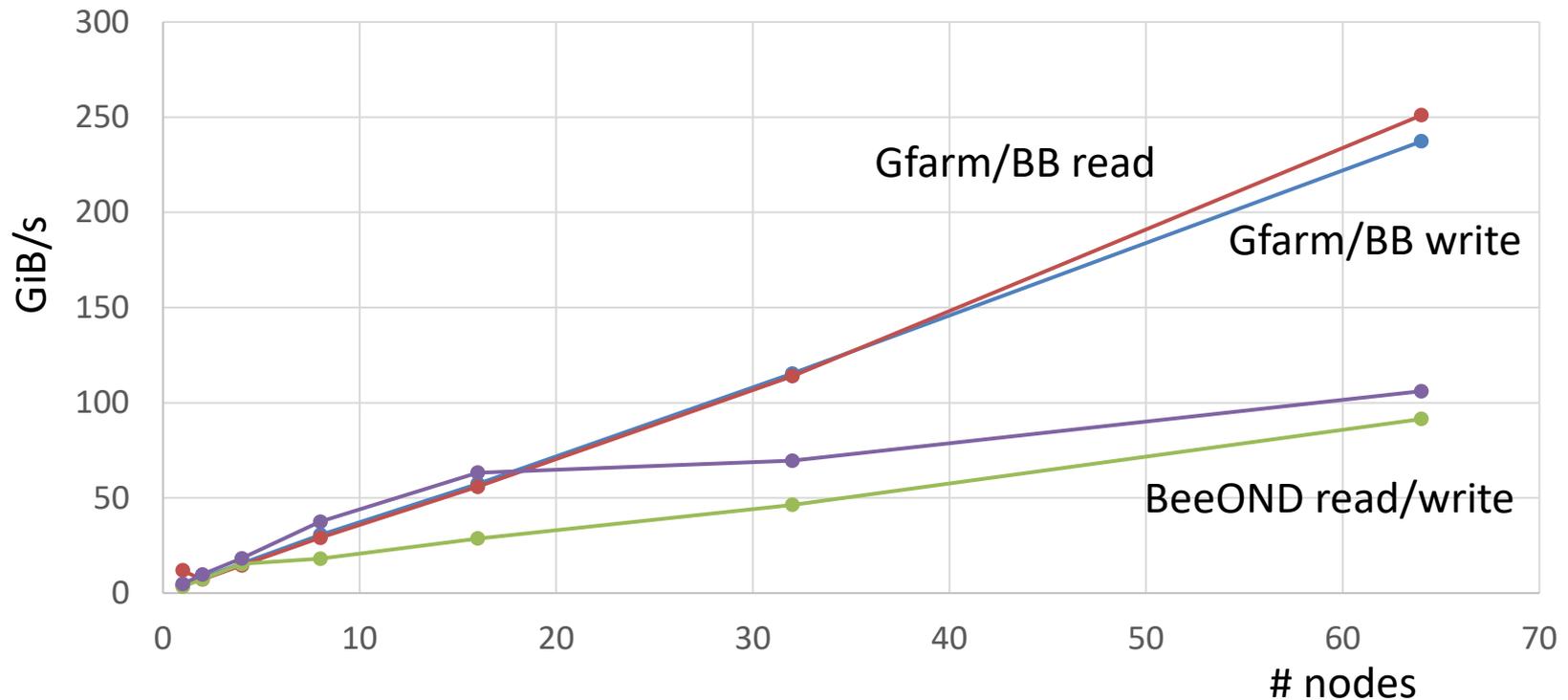
Gfarm/BB



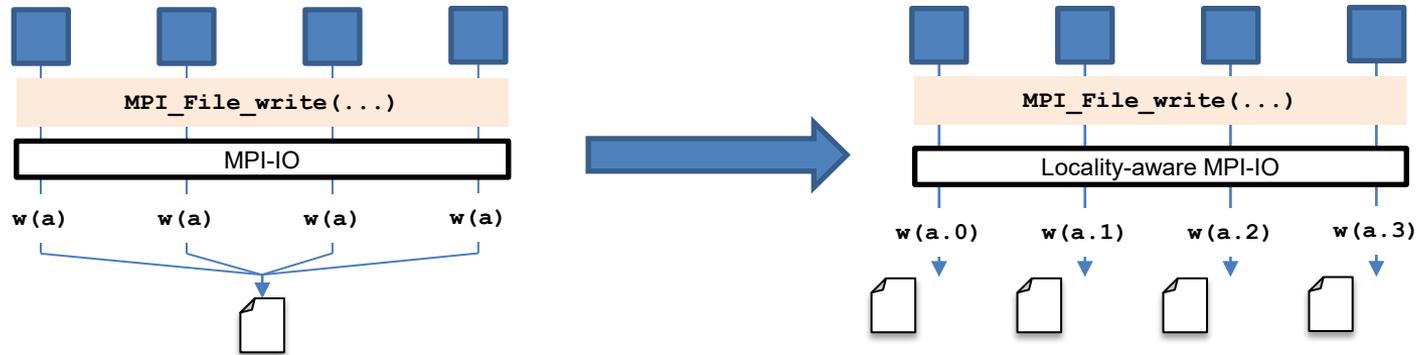
Gfarm/BBノバーストバッファ(2)

```
gfarmbb -h hostfile -m mount_point start  
...  
gfarmbb -h hostfile stop
```

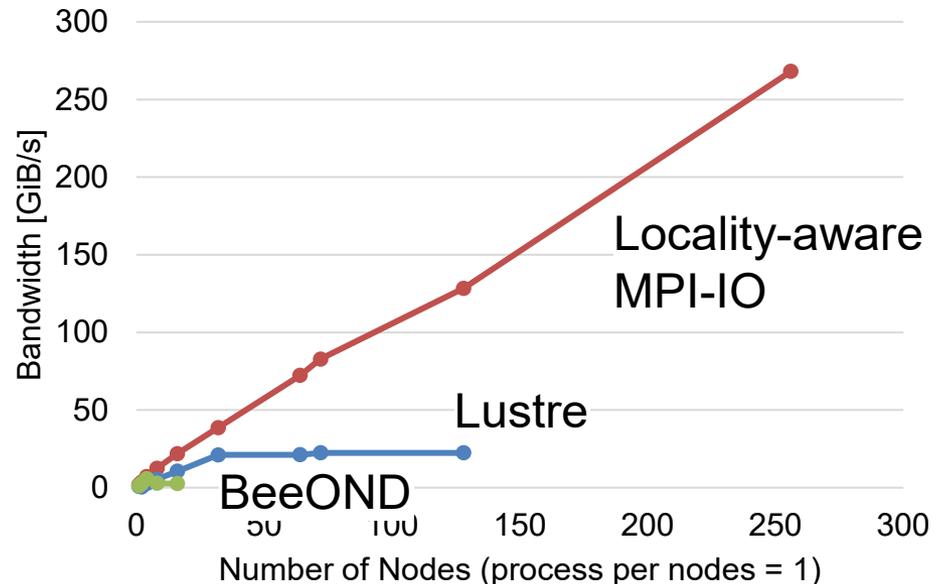
IOR – file-per-process read/write bandwidth on Cygnus
supercomputer



Locality-aware MPI-IO [杉原 2020]



- MPI-IOで分割セグメントの内部形式で格納して、N-1アクセス→N-Nアクセスに変換
- IOR N-1アクセス



まとめ

- Gfarmファイルシステム
 - NPO法人つくばOSS技術支援センターによるサポート
 - Gfarm 2.7.17を9/17にリリース
 - <https://github.com/oss-tsukuba/>
- Gfarm/BBバーストバッファ
- InfiniBand RDMA、ディレクトリクオータ機能
- データ完全性、サイレントデータ損傷対応
- HPCI共用ストレージ、JLDGなど実運用実績
- 進行中
 - IPv6対応 (Gfarm 2.8)