

# The Parallel Cluster File System



## BeeGFS

BeeGFS は小規模の計算環境から大規模の計算環境までハードウェアに依存しない高性能パラレルファイルシステムです。必要なレベルまで容量とパフォーマンスを簡単に向上させることができるよう設計されているため、性能とパフォーマンスを求められる環境には最適です。

ワークロードの種類によっては粒度の細かい大量のファイル I/O を効率的に処理する必要があります。いわゆるスパコンのように高価な Burst Buffer を導入すれば解決できますがコストは大幅に増大します。より廉価な解決策としてノード毎にスクラッチを持たせるのも良い方法ですが、ノードをまたいだジョブの場合には不便です。そこで各ノードに SSD を搭載し、合わせて BeeOND を使うことでこの問題の解決します。

Luster よりもコストが掛からずインストールが楽で、Windows クライアントのサポートもあります。HPC テックではお客様の要望に沿う BeeGFS サーバシステムをご提案させていただきますのでお気軽にご相談ください。

## BeeGFS の主な特徴

### ▶ Distributed File Contents and Metadata

- BeeGFS の基本コンセプトは、アーキテクチャのボトルネックを回避することです。
- 複数のストレージサーバ間でファイルを分散させ、メタデータを複数のメタデータサーバに格納します。
- BeeGFS のこのコンセプトは大規模システムやメタデータ集約型アプリケーションで大きな恩恵を得ることができます。

### ▶ High-Performance Computing

- ファイルシステムノードに HPC では一般的な Infiniband を使用する事でより高速なアクセスを実現しています。
- Infiniband と Ethernet 接続を同時に利用していて、どちらかに障害が発生した場合に自動的に冗長接続パスに切り替わります。

### ▶ Easy To Use

- BeeGFS のサーバに使用するサーバコンポーネントはユーザスペースデーモンです。
- BeeGFS のクライアントはパッチレスカーネルモジュールの追加で、必要に応じて上記サーバと共にいつでも簡単に追加できます。

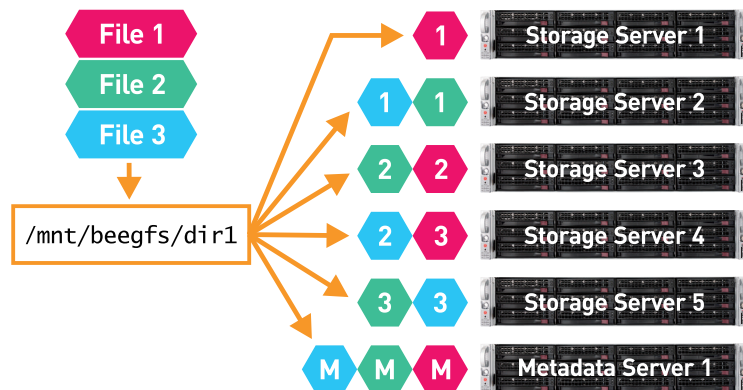
### ▶ Optimized For Highly Concurrent Access

- 一般的な共有ファイルシステムである NFS と違い、I/O 負荷が高い状況でも堅牢性とパフォーマンスを実現するよう設計されています。

### ▶ Clients and Servers On Any Machine

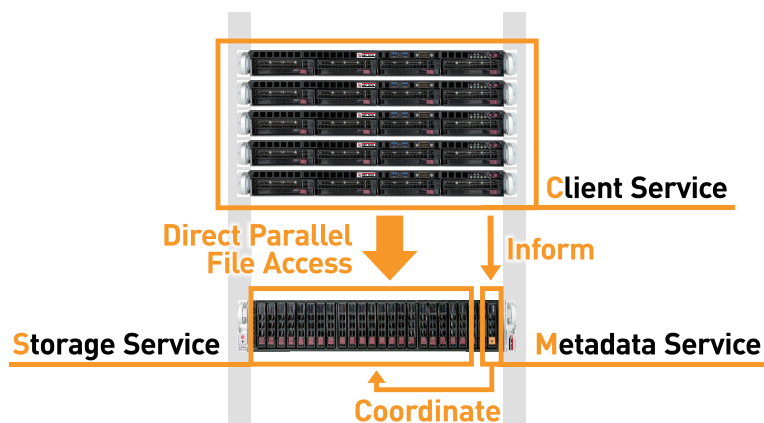
- BeeGFS は特定の OS や専用のファイルシステムパーティションを必要としません。xfs、ext4、もしくは zfs でフォーマットされた既存のパーティションを使用します。
- 大規模システムの場合、構成の異なる複数の BeeGFS ファイルシステムパーティションを作成できます。

## A Hardware-Independent Parallel File System



## BeeGFS Architecture

- ▶ Client Service
- ▶ Storage Service
- ▶ Metadata Service



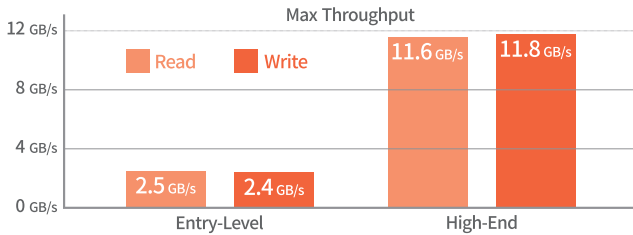
## エントリーレベルとハイエンドレベルの最大スループット比較例

### Entry-Level Building Block

- 8CPU Cores @ 2GHz
- 2x 400GB SSD for OS and BeeGFS Metadata
- 24x 4TB HDDs
- 128GB RAM

### High-End Building Block

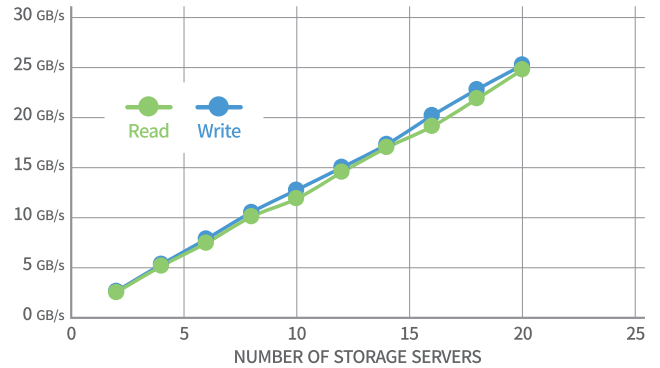
- 36CPU Cores
- 24x NVMe drives with ZFS
- 1TB RAM



## ストレージサーバの台数に応じて性能が伸びる

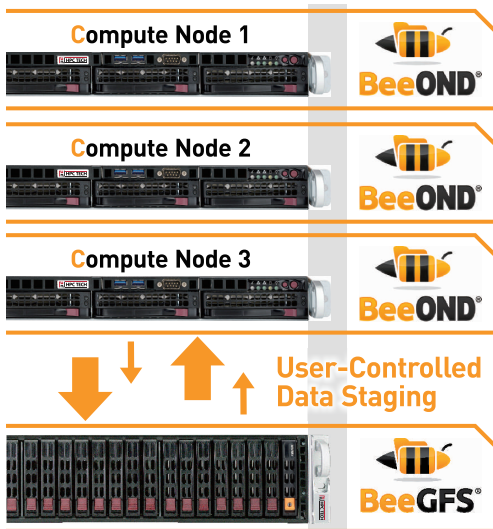
### Sequential Read / Write

up to 20 Servers, 160 Application Processes



## BeeGFS vs Lustre

BeeGFS	vs	Lustre
Support Ubuntu on Servers	OS依存	Support fewer Linux Distribution
簡単	構築	難しい
依存なし	カーネル依存	カーネルの変更が必要
1ノードからOK 小規模システム向けに計算ノードを ストレージサーバとして設定することが可能	ノード数	4ノード以上 Metadata 用サーバ 2台 + Data ストレージ 2台
Yes	Storage Pool	No
OK	Mirroring	NG



## BeeOND - BeeGFS On Demand

BeeOND はジョブに割り当てたノード上にジョブ内からのみアクセス可能な BeeGFS 領域を動的に生成します。この領域をグローバルクラッチとして扱うことで、大量のファイル I/O を効率的に処理する事を可能とします。

- ▶ 並列ファイルシステムインスタンスをその場で作成
- ▶ 簡単な起動/停止コマンド
- ▶ 用途：クラウドコンピューティング、テスト環境、クラスタ計算ノード等
- ▶ Jobスクリプトに記述しスケジュールによる作成が可能 (例：Univa Grid Engine)
- ▶ 一般的な使用例：計算ノードのローカルディスクを利用したクラッチ領域
  - ジョブ実行中の計算ノードがアクセス可能なクラッチ領域を拡大
  - グローバルストレージの利用
  - 性能の出にくい I/O パターンを高速化