

# リアルストレージワークロード特徴抽出のためのデータ 収集蓄積技術のご紹介

大江 和一\*

## 1 概要

大規模データの長期収集、蓄積、分析の1事例として、我々が実施したストレージワークロード収集蓄積技術を紹介する。当社ではストレージシステムの研究開発を進めており、その中でも我々はリアルストレージワークロードにフィットするキャッシュ・階層ストレージ制御方式の研究を進めている。この研究に用いる目的で、3,000 user 前後のアクセスが常時発生する samba ワークロード（最大 4.4TB）のデータ収集を約 1.5 年行った。

本稿では、データ収集を行ったストレージシステムとワークロードの概要、データ収集システムの概要、運用実績、収集したワークロードログより分かること、1PB キャプチャを行う場合の試算、ワークロード収集システムの現場での運用、に関して説明する。

実際に運用を行っているストレージシステムを対象に長期間ワークロード収集した事例は余りなく、我々が遭遇した課題とその解決方法が今後似たようなデータ収集に取り組む皆様の少しでも参考になればと考えている。

## 2 発表内容

発表内容は 3 章の後のスライド資料に示す。

## 3 質疑応答

### 3.1 当日の質疑応答

twitter(#spro2013)に書き込まれた質疑応答内容と我々の記録に残っている質疑応答内容を紹介する。なお、質問・回答 5,6 は我々の記録にのみ残っている質疑応答内容である。

質問 1 ワークロードの偏りは RAID の組み方やファイルシステムで変わるので？

回答 1 ファイルシステムはそうかも知れない。RAID は分からない。

質問 2 NAS 環境だとチャンクの使い回し云々があって consistent な計測はブロックレベルでは難しそうだが何とかなるのか？

---

\*株式会社富士通研究所 ICT システム研究所

- 回答 2 今回の計測では全部 Windows だがそこそこと一致していた
- 質問 3 今回の計測データを製品の改良に活かすことを考えるとどうなつか
- 回答 3 まだ取らなければいけない情報があるがなかなか難しい
- 質問 4 この話を進めていくと上で何が起きているか分からぬが DB 側で何とか最適化にしてあげる方向にいくか
- 回答 4 まだ製品に入ってないが、自分の興味としてはそういうことをしたい
- 回答 4 週ごとの差異なども見えるのでキャパシティプランニングなどに使えばよいと思っている
- 質問 5 1 サンプルの分析で大丈夫でしょうか (samba のバグなどで偏りが出ている、とか)
- 回答 5 MSR Cambridge を分析しても同様な結果が出ており、今回紹介したようなワーカーロードが少なからず存在する、と判断している。
- 質問 6 ファイルシステムレベルやもっと上位のログは必要ない？
- 回答 6 上位アプリ、ファイルシステムレベル、ロックレベルで串刺しで取るのが理想ではある。実験室レベルでは出来るが運用環境ではなかなか難しいので、今回の収集でもロックレベルのみです。逆にロックレベルの分析でうまく一般化出来れば何処へも持っていくことができます。

### 3.2 twitterへの書き込みを踏まえた補足説明

講演中に twitter に書き込んで頂いた内容を踏まえた補足説明を行う。

#### 3.2.1 LBA レベルの分析で何が分かるのか

アプリケーションごとに大まかな特徴を明確にすることは出来、特徴に基づいたストレージ制御方法を検討することが出来ると考えている。今回紹介したファイルサーバを例に取ると、負荷の偏りが発生する LBA の幅と全 IO に占める割合の傾向が把握できれば、その負荷を対象にする階層制御方式等を提案可能である。

もちろん、同一アプリケーションにおいても OS やファイルシステムが異なると LBA レベルでの偏りなども異なってしまう可能性があり、正確なアプリケーションの分析と言えないが、ストレージ制御の観点からは十分有効な分析であると考えている。

**FUJITSU**  
shaping tomorrow with you

## リアルストレージワークロード特徴 抽出のためのデータ収集蓄積技術のご紹介

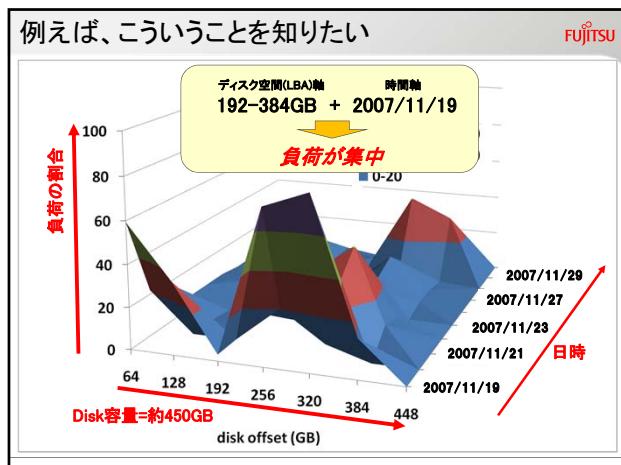
2013年8月25日  
(株)富士通研究所 大江 和一

2013年夏のプログラミング・シンポジウム  
ビューティフルデータ 講演資料

■:講演後に書き込んだ注釈

### Agenda

1. ワークロード収集の概要
2. 収集したワークロードログより分かること
3. 1PBキャプチャを行う場合の試算
4. ワークロード収集システムの現場での運用について
5. まとめ



**リアルストレージワークロード収集の意義**

目的: 時間的・空間的局所性を分析⇒一般化

ストレージ制御に関する新たなシステム提案が出来る、かも

- ・従来のstatic tieringでは性能向上しないワークロードへの方式提案
- ・新たなcacheアルゴリズムの提案
- ・省電力の観点からのSSD/HDD tiering方式提案、など

■ストレージへのデータ量: 爆発的に増大  
⇒ストレージアクセス履歴量も膨大に…

■どうやって解決するのかが共通の課題

- ストレージアクセス履歴保存に必要なディスク容量
- ストレージアクセス履歴の増加に対応可能なシステム性能  
⇒ 我々の経験が少しでも参考になれば…

**関連研究**

■ 主なストレージワークロード分析論文

- 仮想環境でのストレージワークロード分析
  - Storage Workload Characterization and Consolidation in Virtualized Environments; Ajay Gulati, Chethan Kumar, Ifran Ahmad (VPACT 2009)
  - An Analysis of Disk Performance in VMware ESX Server Virtual Machines; Ifran Ahmad, Jennifer M. Anderson, Anne M. Holler, Rajit Kambo, Vikram Makhija (6th Annual Workshop on Workload Characterization, 2003)
- Windows Serverのワークロード分析
  - Characterization of Storage Workload Traces from Production Windows Servers; Swaroop Kavalanekar, Bruce Worthington, Qi Zhang, Vishal Sharda (ISWC2008)
- UNIXディスクアクセスパターン分析
  - UNIX disk access patterns; Chris Ruemmler and John Wilkers (USENIX Winter 1993 Technical Conference)

LBAごとの幅まで読み込んだ長期間のワークロード分析論文を見つけられなかった

■ 一般公開されているストレージワークロード

- MSR Cambridgeトレースデータ(<http://iotta.snia.org/traces/388>)
  - Write off-loading: Practical Power Management for Enterprise Storage (USENIX FAST'08)で紹介  
⇒ プロジェクト開始時点では未公開

自分たちで必要なデータを採取することに

**ワークロード収集の概要**

■ 収集システム開発 & ワークロード収集を行った期間
 

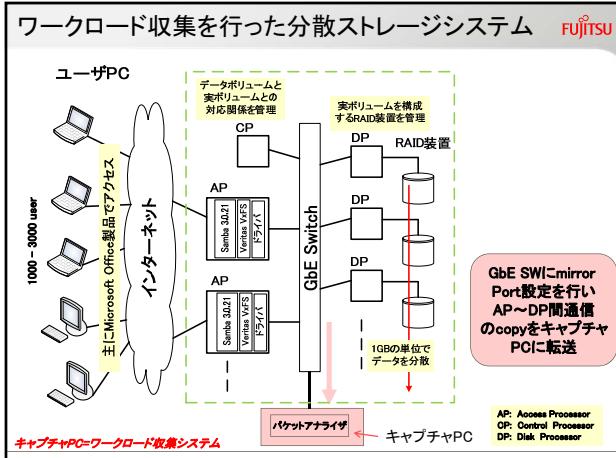
- 2007.10～2010.12  
・継続してログ収集が出来たのは、最後の1.5年間

■ 収集を行ったストレージシステム
 

- スケールアウト型分散ストレージシステム
  - ・研究所開発システムがベース ⇒ 仕様入手 & 検証が容易
  - ・実運用システム ⇒ まとまったユーザ数のワークロードが収集可能
  - ワークロード収集を行った主なストレージシステム
    - ・Samba+backup(35TB): 数ヶ月単位のワークロードを収集
    - ・Samba(4.4TB): 連続して約1.5年分のワークロードを収集

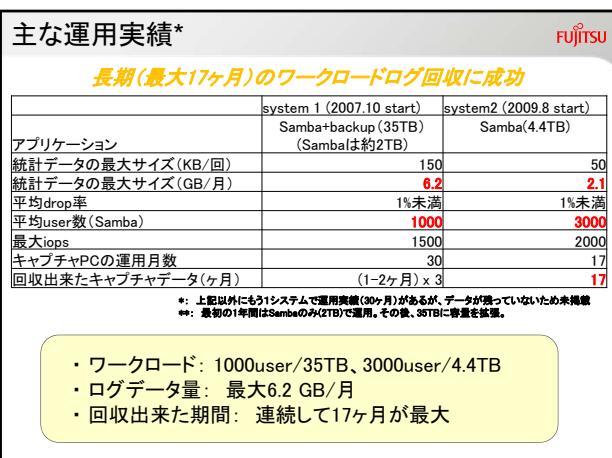
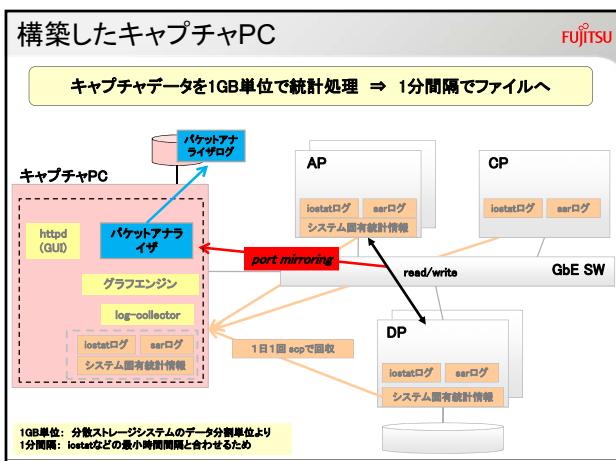
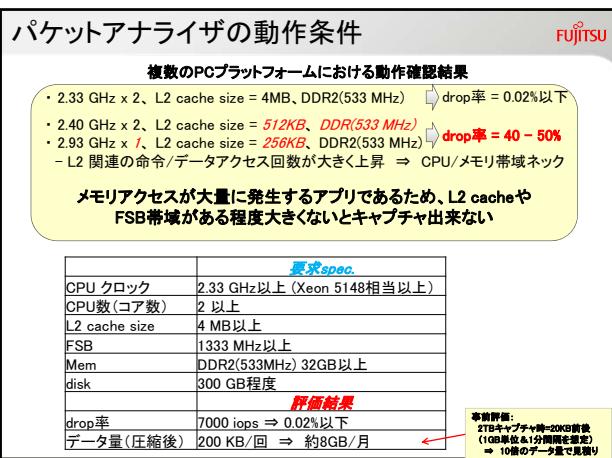
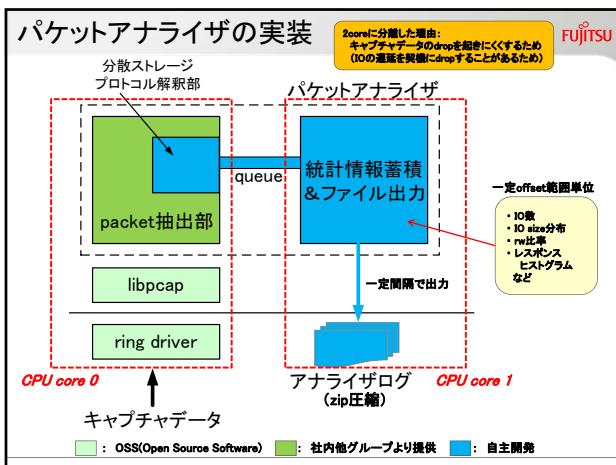
■ 収集システムの課題と解決方法
 

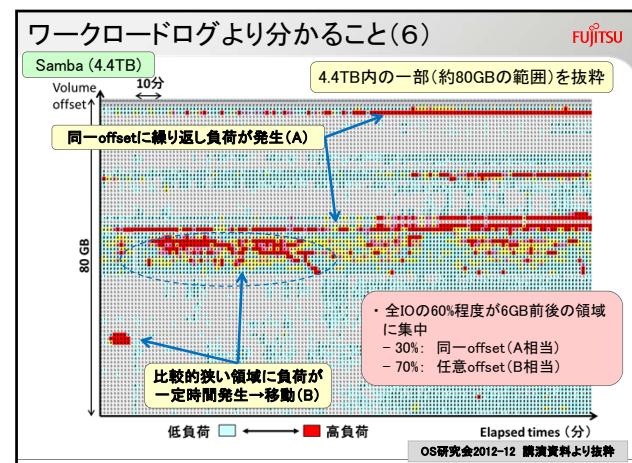
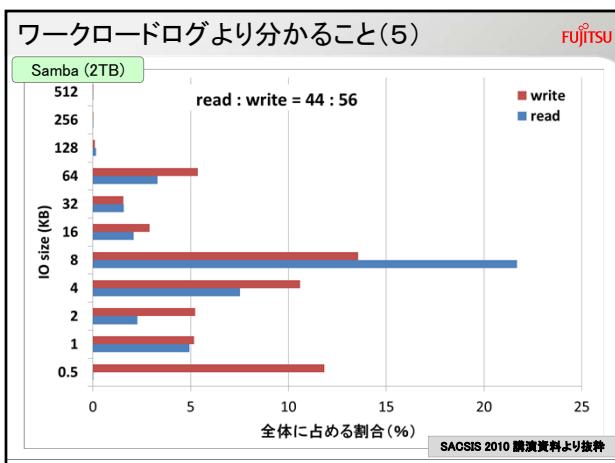
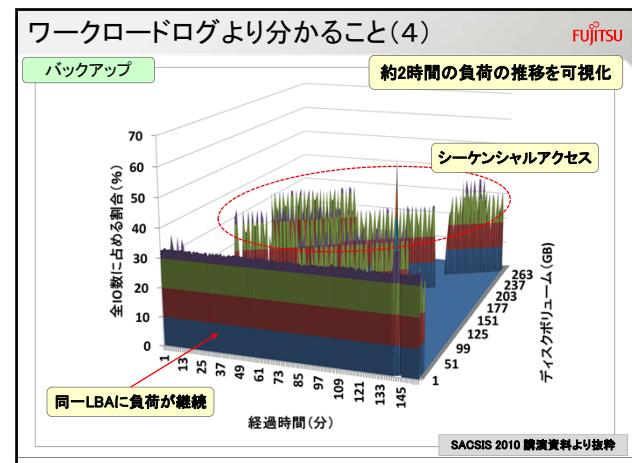
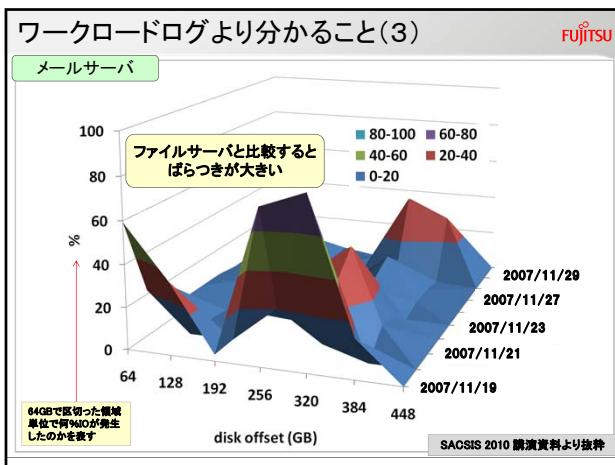
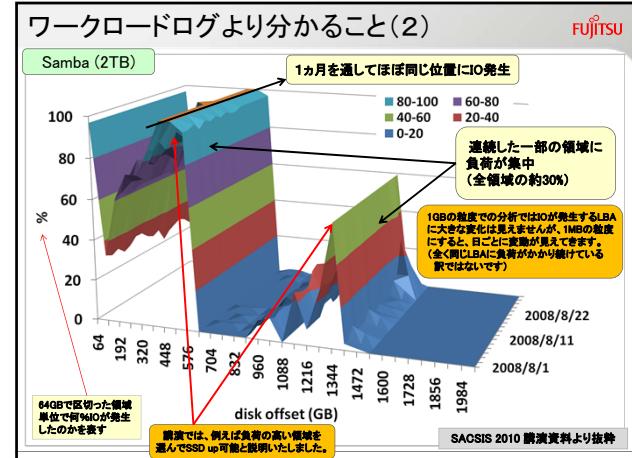
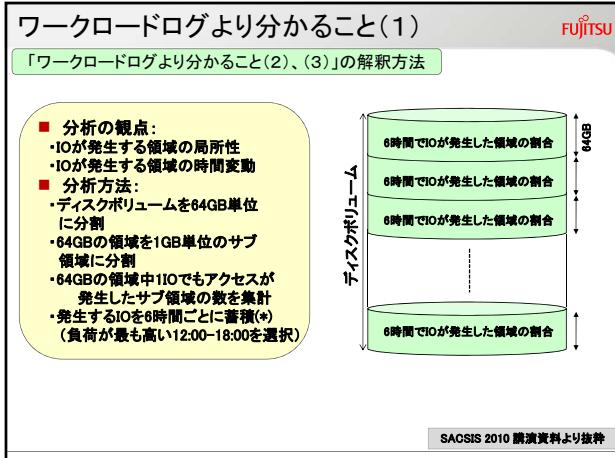
- 障害原因とならないこと
  - ⇒ GbE SW mirroring機能を用いた主系とは別PC上のパケットキャプチャ
- ごく普通の1U IA serverで運用 & 頻繁にデータ回収出来ない
  - ⇒ 生トレースの保存は諦め、1GB・1分単位の統計情報を圧縮・保存

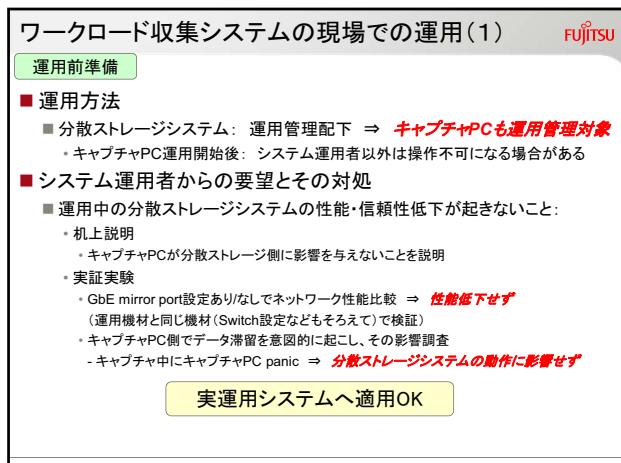
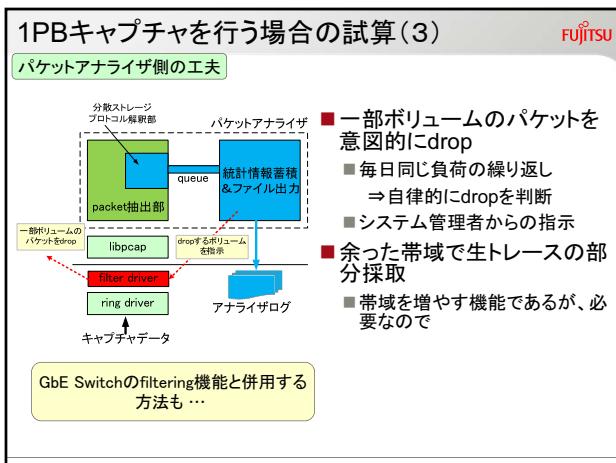
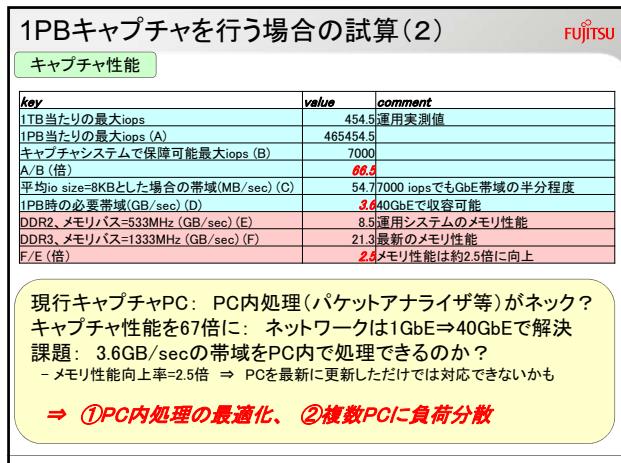
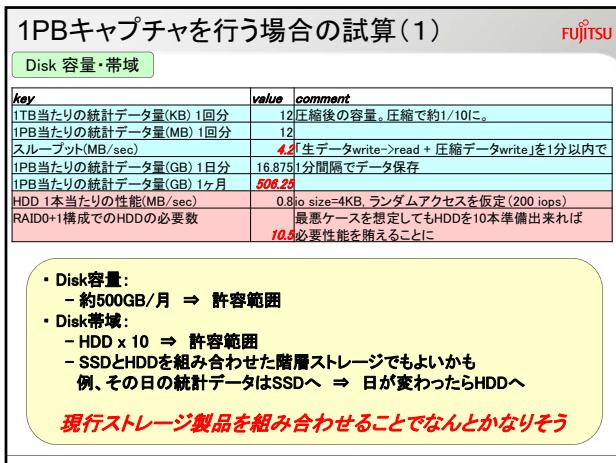
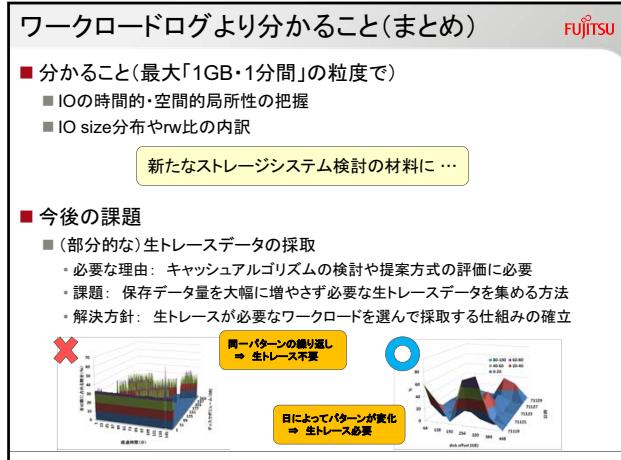
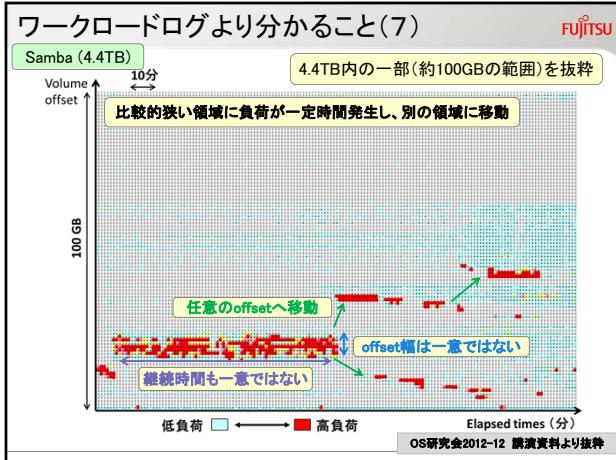


## キャプチャPCの課題と解決方法 FUJITSU

- ストレージシステム側との関係:
    - いかなるケースでもストレージシステム側の障害原因とならないこと  
⇒ *GbE SW mirroring機能を利用することで別PCで分析*
  - キャプチャPCの要件:
    - 設置面積: 1U程度に収まること
    - マシンスペック: 比較的安価なサーバを用いてもログ収集できること
    - データ回収: 数ヶ月程度はスタンダードで動き続けること  
⇒ *生トレースの保存は除外、1GB・1分単位で統計情報を抽出し、圧縮・保存*
  - パケットアライザ\*開発方針:
    - 性能: 収集ワークロードに耐える程度の性能  
(むやみに高性能を狙わない)
    - 実装: 自主開発コードは必要最小限に&長期間安定動作  
⇒ *収集ワークロード負荷は分からないので、考えられる対策を行った上であとは実際に運用してみて対策を考えることに…*
- \*: キャプチャPC上で動かすソフトウェア。このソフトウェアが収集したワークロードの統計処理・ファイルへの統計情報出力を行う







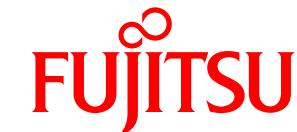
ワークロード収集システムの現場での運用(2) FUJITSU

**運用**

- インストール
  - 運用者のインストール作業に立ち会ってのインストール作業サポートも必要
    - OSのカスタマイズ、GbE Switchの設定、パケットアナライザの初期化、など
- 統計データ(アナライザログ)の回収
  - 運用者の作業に同行して回収
    - 管理用ネットワーク帯域が細い、あるいは、潤沢に使えない場合がある
    - 遠地に設置してある場合には、頻繁な回収は困難
    - 運用者側の業務スケジュールなどにも依存するので事前の計画が重要に
- 日々の運用
  - 運用手順を簡略化できないと、設定誤りのため正しく動作しないケースも
    - 例、ボリュームを追加 ⇒ 運用者がパケットアナライザの設定変更が必要に  
⇒ 出来る限り自動化しておく必要性を痛感
  - バージョンアップ
    - 運用者のバージョンアップ作業に立ち会う必要性がある場合も
    - 3年間ほぼ安定して動作させることができた

まとめ FUJITSU

- ストレージワークロード収集の概要
  - 4.4TBと35TBの分散ストレージシステムが対象
  - 最大17ヶ月間のワークロードログ収集に成功(4.4TB)
    - 「1GB・1分間」の粒度 ⇒ 最大6.2 GB/月のデータ量
- 収集したワークロードログより分かること
  - 時間的・空間的局所性の特徴抽出が可能
    - ⇒ 新たなストレージシステム検討可能
- 1PBキャプチャを行う場合の試算
  - 4.4TB分散ストレージシステムのログ量が単純にスケールする前提
  - PCのメモリ・バス帯域ネックになる可能性
    - ⇒ ①キャプチャソフトウェアの最適化、②複数PCに分散してキャプチャ
- ワークロード収集システムの現場運用
  - 本番(分散ストレージシステム)系の障害原因とならないこと
  - 「エラーケースまで含めた検証済み手順の確立 + 現場で容易な運用」が重要



shaping tomorrow with you