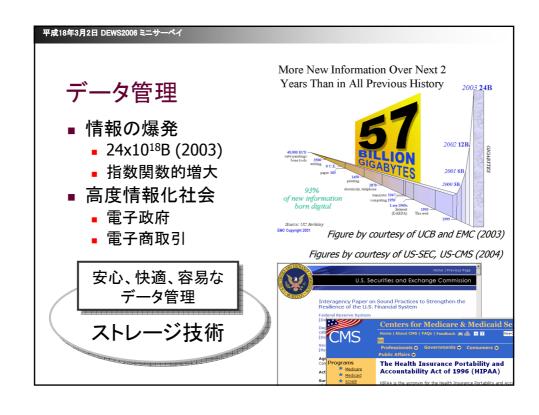
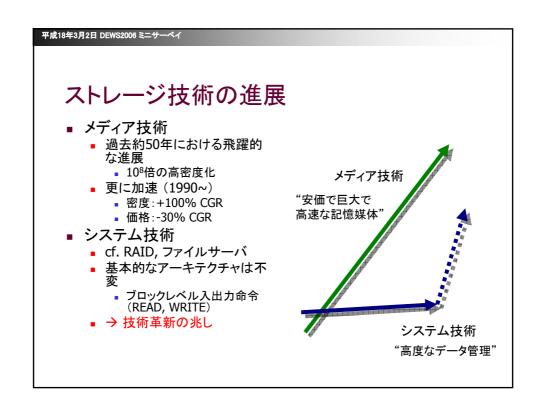
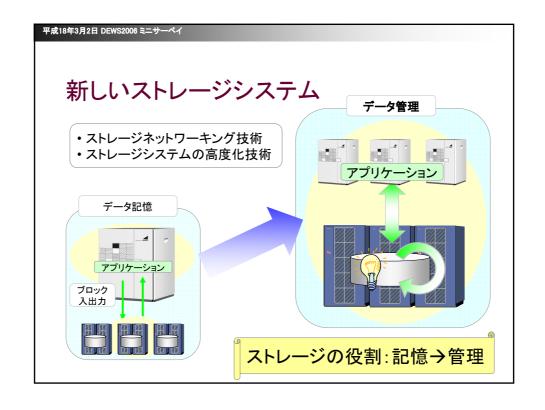
DEWS2006 ミニサーベイ/平成18年3月2日

ストレージシステムの高度化技術に関する研究動向

東京大学 生産技術研究所 産学官連携研究員 合田和生 kgoda@tkl.iis.u-tokyo.ac.jp





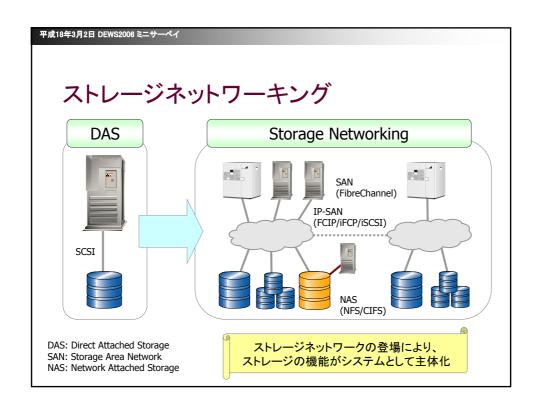


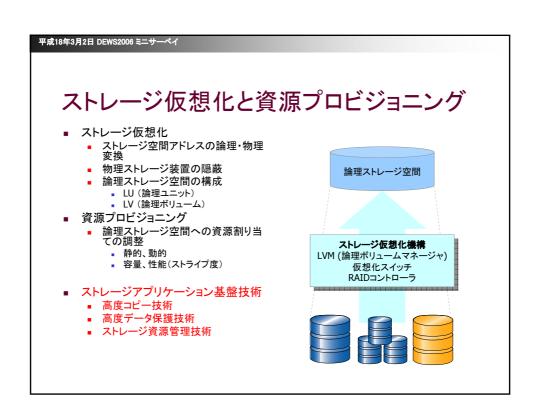
アジェンダ

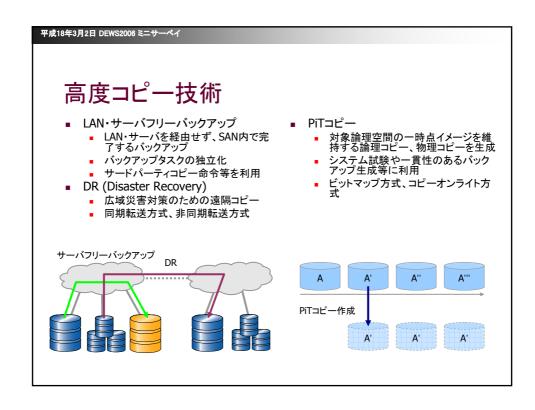
- ストレージネットワーキング
- 最近のストレージシステム高度化技術
- ストレージフュージョン技術

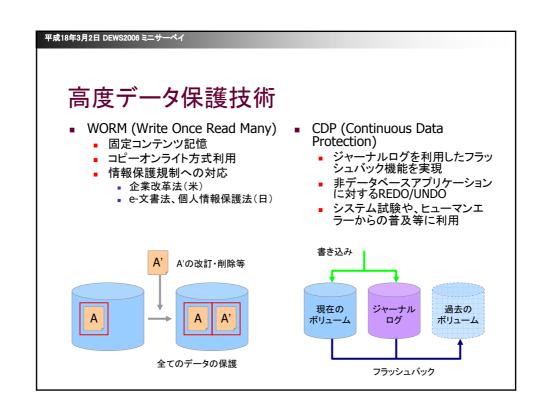
平成18年3月2日 DEWS2006 ミニサーベイ

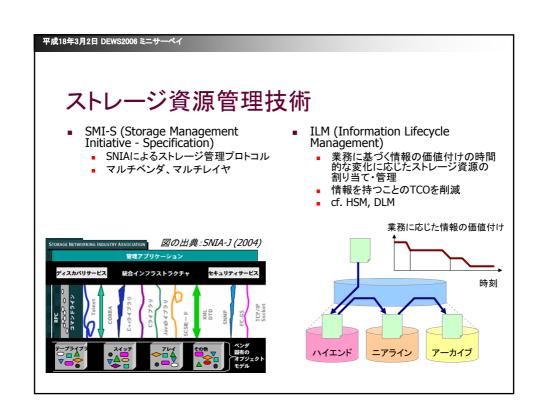
ストレージネットワーキング











ストレージネットワーキング

- SAN、NAS等のストレージネットワーク
- ストレージ仮想化とプロビジョニング
- 高度なストレージアプリケーション
- → ストレージがシステムとして機能

最近のストレージシステム高度化技術

平成18年3月2日 DEWS2006 ミニサーベイ

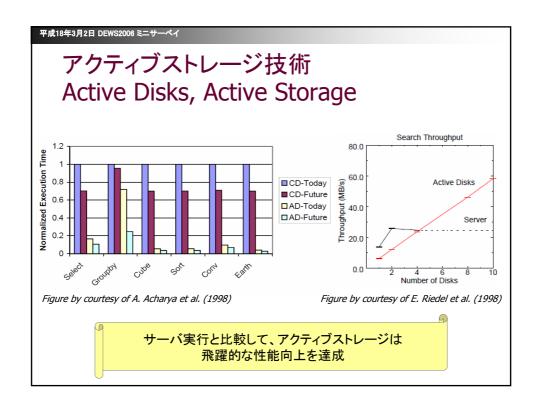
ストレージシステムの高度化技術

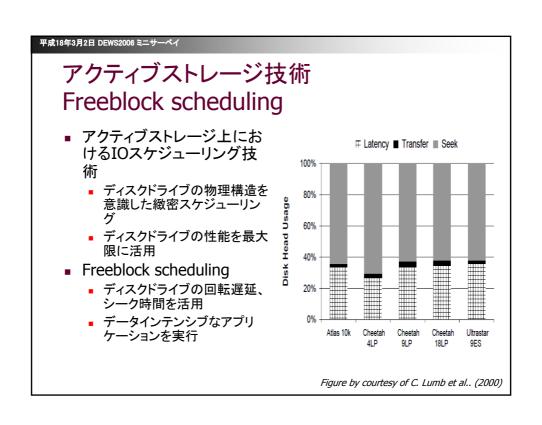
- データベースマシン(1970~80年代)
 - "Logic per track" (D. L. Slotnick, 1980)
 - 専用ハードウェアにより高性能処理を目指す
 - プロセッサ性能が不十分
 - 性能利得によるコスト正当化が不十分
- 近年
 - 情報爆発時代における巨大データ管理の社会的要請
 - ストレージ装置の計算機資源の向上
 - ストレージネットワーキングによる環境整備
- → 再び、ストレージシステムの高度化の流れ

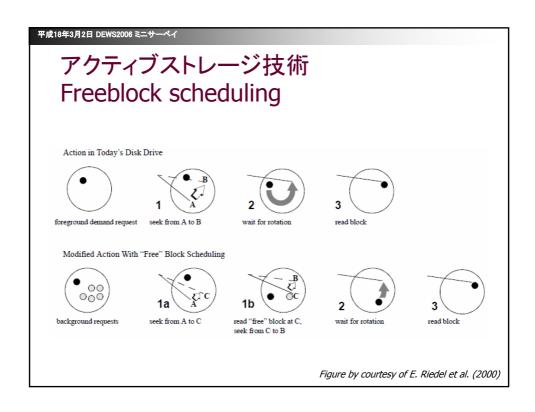
最近のストレージシステム高度化技術

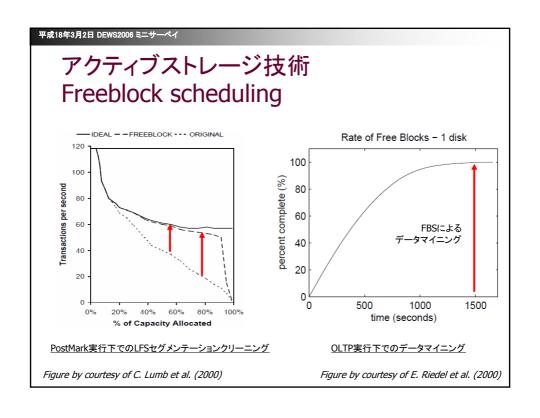
- アクティブストレージ技術
- コラボラティブストレージ技術
- オートノミックストレージ技術

平成18年3月2日 DEWS2006 ミニサーベイ アクティブストレージ技術 ストレージ上のプロセッサを用いてデータインテンシブアプリケーションを実行 データになるべく近いところでコードを実行 高い内部帯域を有効活用 高いディスクアクセス並列度 →高いIO帯域 →高いアプリケーション処理スループット Active Disks (UCSB, 1998) Active Storage (CMU, 1998) iDisks (UCB, 1998) 基本的なアイデアはデータベースマシンに類似 ブロセッサ技術の向上によって支持



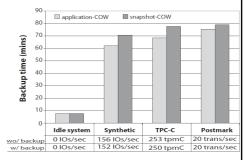






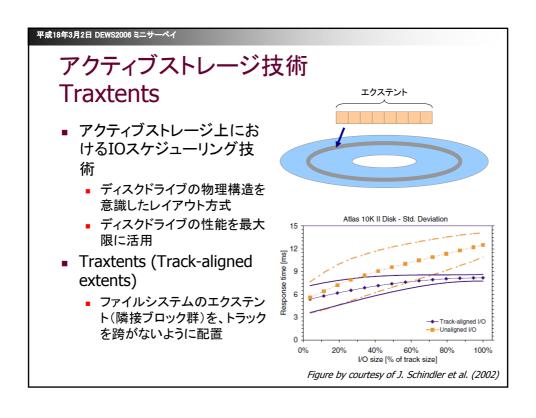
ではは年3月2日 DEWS2006 ミニサーペイ アクティブストレージ技術 Freeblock scheduling

- 一般フレームワーク化
 - フォアグランドアプリケー ション
 - → 同期IO
 - バックグラウンドアプリケーション
 - → 非同期IO (FBS対象)
- データ管理ジョブに利用
 - PiTコピー作成
 - バッファクリーニング
 - レイアウト再編成



Freeblock scheduling APIを用いた スナップショット作成

Figure by courtesy of E. Thereska et al. (2003)



平成18年3月2日 DEWS2006 ミニサーベイ アクティブストレージ技術 Adjacent blocks Adjacent blocks ■ 多次元データの格納に利用 回転遅延・シーク時間から一 ディスクドライブの物理構造を意識した多次元データフブロックの 定時間で到達可能なブロック 候補群 マップ (MultiMap) 各次元方向の走査に「程々の」 スループットを期待 余計な回転遅延を削減 4D earthquake dataset. Maxtor Atlas 10k III Naive Hilbert MultiMap z (semi-seq) 6000 y (semi-seq) 4000 x (seq) 10⁴ Figure by courtesy of S. Scholosser et al. (2003) Figure by courtesy of S. Scholosser et al. (2003)

平成18年3月2日 DEWS2006 ミニサーベイ

アクティブストレージ技術 まとめ

- ストレージ上のプロセッサを用いてデータインテンシ ブアプリケーションを実行
 - 高いディスク並列を以って高いIO内部帯域を活用
 - 物理情報を利用したIOスケジューリングとデータ配置技術
 - → 高いアプリケーション処理スループットを目指す
- 課題
 - プロセッサ性能
 - インターフェースの標準化
 - 専用ハードウェアの開発コスト

コラボラティブストレージ技術

アプリケーションはサーバ上で実行

平成18年3月2日 DEWS2006 ミニサーベイ

- ストレージがサーバ上のア プリケーションを意識して動作
 - ヒント情報 専用のインターフェースを通じ て情報交換
 - アプリケーション予測 ストレージ側でアプリケーション動作を高度に予測
- ストレージアーキテクチャの 大幅な改編を回避

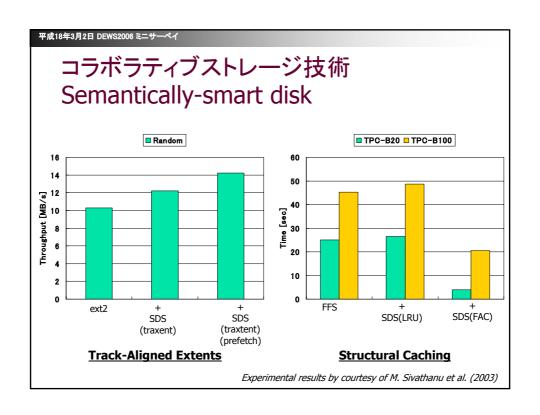


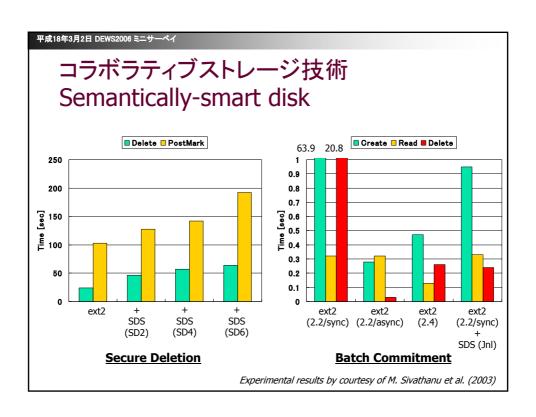
平成18年3月2日 DEWS2006 ミニサーベイ

コラボラティブストレージ技術 Semantically-smart disk

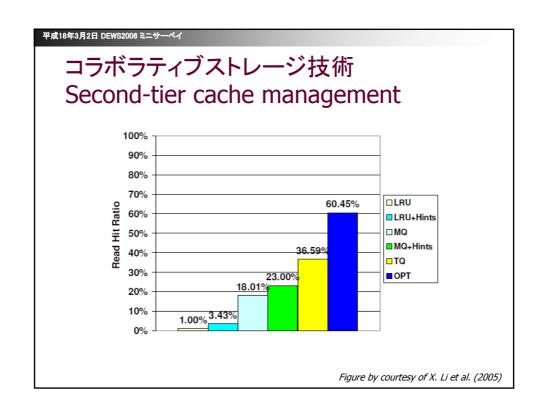
- 上位のファイルシステムを 理解して動作する「賢い」 ディスク
 - アプリケーションコード SCSIインターフェースは不変
 - ディスクのプロセッサ、メモリを 活用
 - ファイルシステムを補助
- → 高い性能、高い可用性、 高い保安性
- データベースにも拡張

- ファイルシステムの構造・操作を類推
 - FFS, ext2fsを仮定
 - ブロックの同定(データ、iノード、ビットマップ、ジャーナル、メタデータ)
 - 直接的ブロック分類: プロー ブプロセスを利用
 - 間接的ブロック分類: オンラインでのブロック監視
 - ブロック間関連の抽出
 - 高レベルファイル操作の類推
- ファイルシステムの操作を 補助





平成18年3月2日 DEWS2006 ミニサーベイ コラボラティブストレージ技術 Second-tier cache management 2つのキャッシュが競合する Database Clients ことの無駄を排除 バッファキャッシュ(サーバ) Transaction ディスクキャッシュ(ストレージ) サーバ側でWRITEコマンド にヒント付け Read/Write sync. write async. replacement write Storage Server async. recoverability write TQ (type queue) キャッシュ管理アルゴリズム ■ ストレージ側でWRITEヒントを 活用 Figure by courtesy of X. Li et al. (2005)

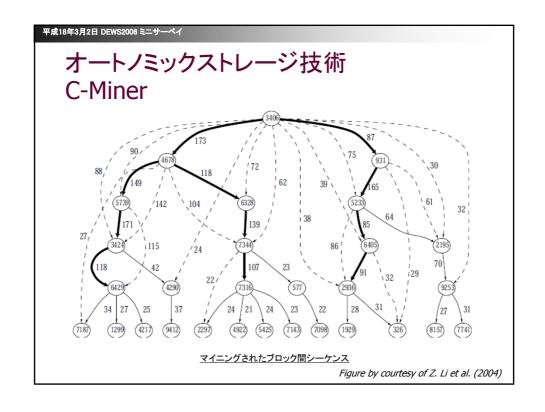


コラボラティブストレージのまとめ

- ホスト上のアプリケーションとの緩やかな連携
 - ストレージシステムが、アプリケーションを意識することにより自己調整
- ストレージアーキテクチャの大幅な改編を回避
- 課題:
 - インターフェースの共通化
 - ヒント情報 vs. アプリケーション予測



オートノミックストレージ技術 C-Miner 「ブロックレベルの典型的な IO振舞を抽出し、IO制御に 利用 IOトレースベースの学習方式 「プロック間相関抽出マイニン グ技術 "C-Miner" (frequent sequence mining) JUNグターンをキャッシュ管理に活用 CDP (Correlation-Directed Prefetching) Layout reorganization Figure by courtesy of Z. Li et al. (2004)

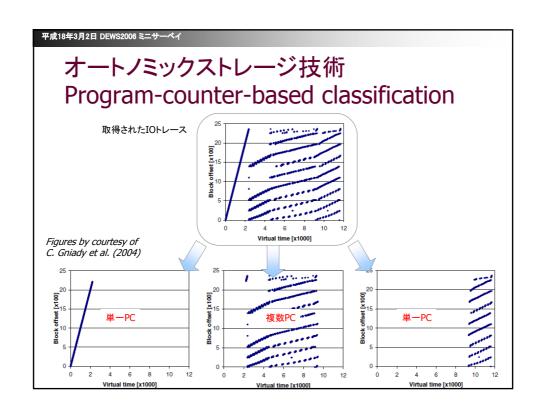


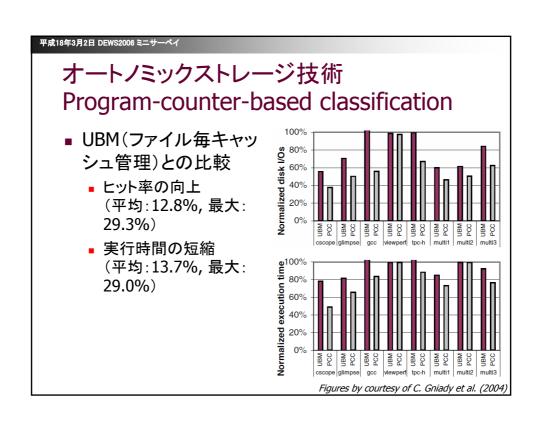
オートノミックストレージ技術 C-Miner IO応答性能の削減効果 (Cello92/96, TPC-C, OLTP) ・線形プリフェッチによる改善: 7-20% CDP+Layoutによる改善: 12-25%

平成18年3月2日 DEWS2006 ミニサーベイ

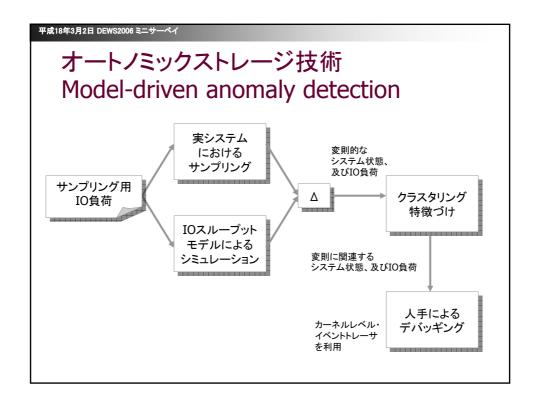
オートノミックストレージ技術 Program-counter-based classification

- サーバからのIO発行情報を活用
 - PC(プログラムカウンタ)を利用してIOパターンを 解析
 - PCからソフトウェアサブルーチン(IOが発行されたソフトウェア文脈)を類推
 - 頻出IOパターンを抽出
 - sequential, loop, others
- → バッファキャッシュ管理に利用





平成18年3月2日 DEWS2006 ミニサーベイ オートノミックストレージ技術 Model-driven anomaly detection ■ IO負荷とスループットモ Workload characteristics System I/O throughput デルに基づく変則検出 ■ 4層のモデル化 OSキャッシュモデル throughput' workload' 先読みモデル OS configuratio I/O prefetching model ■ スケジューリングモデル throughput" ストレージ装置モデル I/O scheduling ■ → IOシステム全体を捉 model えるデバッギング方式 throughput" workload" Storage properties Figure by courtesy of K. Shen et al. (2005)



オートノミックストレージ技術 Model-driven anomaly detection

- Linuxカーネルに実在する バグを発見
 - 先読み(1)
 - 予測スケジューリング(2)
 - エレベーションスケジューリン グ(1)
- 全工程
 - サンプル取得: 6分/サンプルル → 40時間/400サンプル
 - クラスタリング、特徴付け: <1分
 - 人手によるデバッグ: 1~2日 /バグ → ~6日/全4バグ

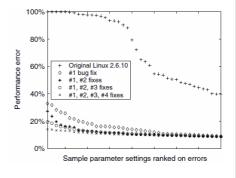


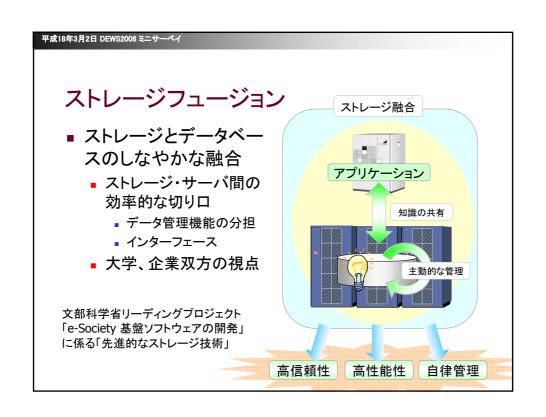
Figure by courtesy of K. Shen et al. (2005)

平成18年3月2日 DEWS2006 ミニサーベイ

オートノミックストレージ技術 まとめ

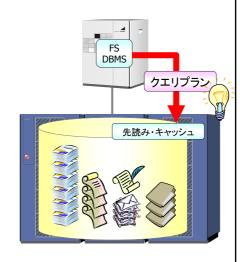
- IOトレースを広範囲に活用
- 新しい解析技術
 - データマイニング技術
 - ソフトウェアフックによる文脈予測技術
 - モデル依存変則検出技術
- ■課題
 - 実システムに対する検証

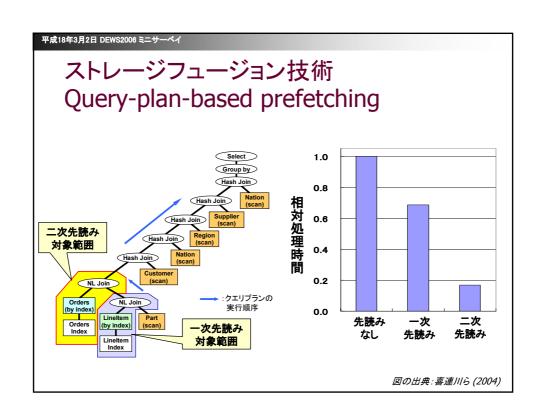
平成18年3月2日 DEWS2006 ミニサーベイ ストレージフュージョン技術

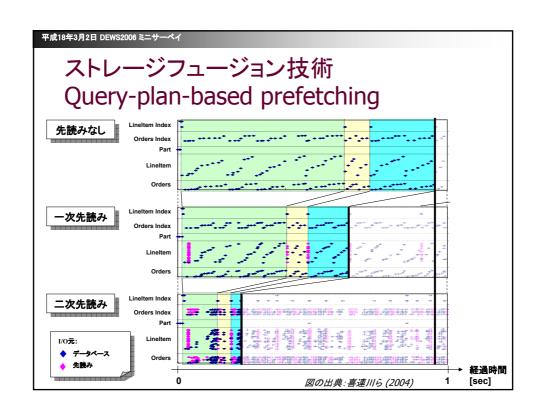


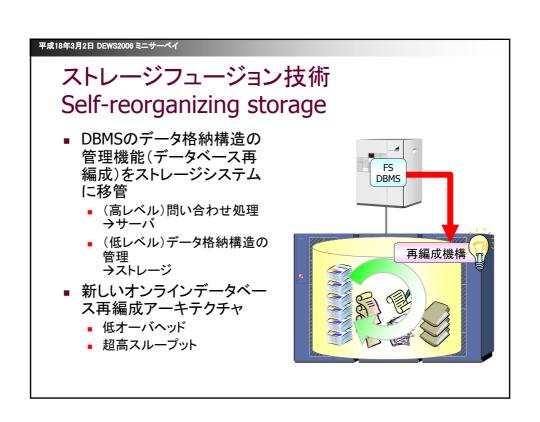
平成18年3月2日 DEWS2006 ミニサーベイ ストレージフュージョン技術 Query-plan-based prefetching

- 問い合わせ実行計画を利用 して、データベース構造に基 づいた先読みを実現
 - QPをストレージに注入し、索引構造に基づき先読み
 - ストレージ装置の有するキャッシュを活用
- cf. 線形先読み (sequential prefetch)

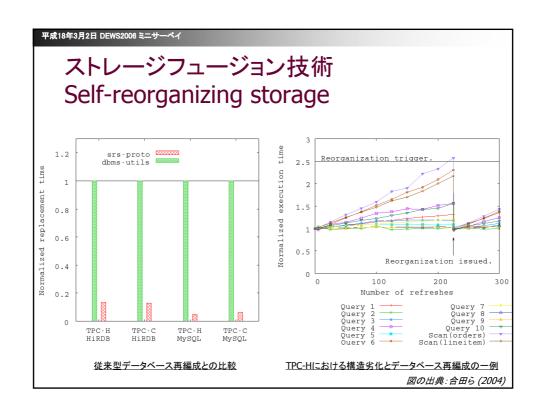


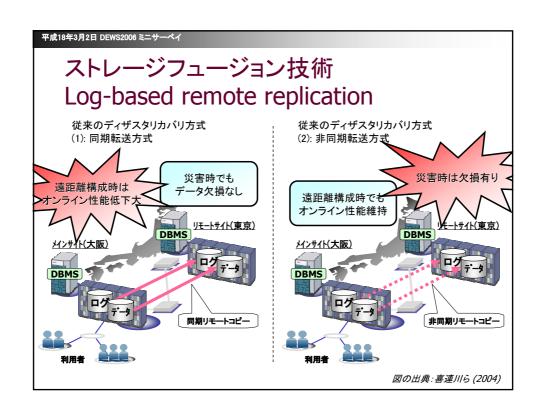


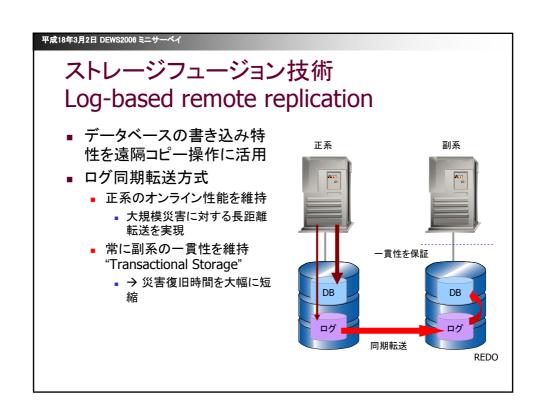


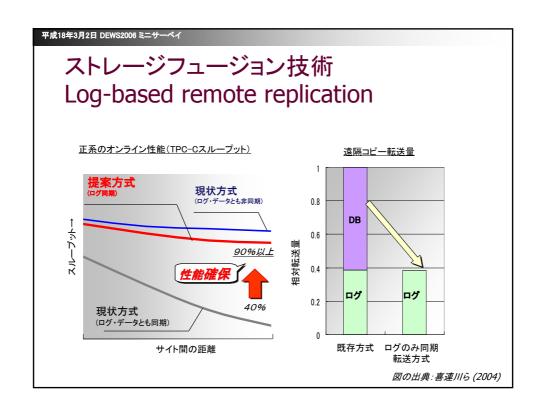


平成18年3月2日 DEWS2006 ミニサーベイ ストレージフュージョン技術 Self-reorganizing storage ■ 10TB級の巨大データベース の再編成を想定 サーバ側 サーバ側 同時実行再編成 分離再編成 ■ サーバ-ストレージ間IO帯域 がボトルネック ストレージシステムの高い 内部IO帯域を活用 ■ ストレージ側分離再編成戦略 ■ IO指向オンライン再編成処理 ■ 物理アドレスレベルIOスケ ジューリング パイプライン化データベース 再編成 高速ログ適用処理 差分適用









ストレージフュージョン技術のまとめ

- サーバ上のアプリケーションとしなやかに融合するストレージ の高度化技術
 - 大学、企業双方の視点
 - 新しいストレージ融合技術を提案、検証
- 新しいアプローチによる高性能、高可用、低管理コストの実現
 - Query plan based prefetching
 - データベース中の索引情報をストレージのキャッシュ管理に活用
 - Self-reorganizing storage
 - データベースの格納構造管理をストレージに移管
 - Log-based remote replication
 - データベースの書き込み特性を遠隔コピー操作に活用

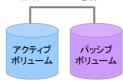
その他の話題

平成18年3月2日 DEWS2006 ミニサーベイ その他の話題 ストレージの自律分散化 OceanStore (UCB) 分散ハッシュに基づくP2Pストレー ジシステム ■ 冗長符号化、ビザンティンプロトコ 自律ディスク (東工大) ディスクプロセッサを利用した高 性能分散ファイルシステム ログによる自律的複製取得 Figure by courtesy of S. Rhea et al. (2003) ■ Fat-B木によるデータ空間配置の 自律的バランス化 Host₂ Host 1 FAB (HP Lab.) ブリック(PC)を統合した大規模 ディスクアレイ 多数決原理に基づく一貫性管理 プロトコル Figure by courtesy of H. Yokota. (1999)

その他の話題 ストレージの省電力化

- MAID (Copan Syst.)
 - アクティブボリュームとパッシブボ リュームの2階層ストレージ構成
 - 常時通電のアクティブボリュームをキャッシュとして利用
 - 低負荷時にパッシブボリュームの給 電を停止
 - 主に、アーカイブストレージとして利用
- AutoMAID (Nexsan Tech.)
 - 低頻度アクセス時に、全ディスクをス タンパイ化
 - 複数のスタンバイモードをサポートアーカイブストレージ以外に、ニアラインストレージ、オンラインストレージとしても利用可能
- Hibernator (HP Lab.)
- Dynamic RPM (Pa. State Univ.)
 - 動的に回転速度を変更可能な次世代 ディスクドライブの活用を目指す

2階層ストレージを構成



"always powered on" "sometimes powered off" キャッシュとして機能 本体ストレージとして機能

60% savings – after 5 minutes of idle time, drives are programmed to unload heads and slow down to 4000 RPM. Very fast recovery time.

80% savings – after 5 minutes of idle time, drives are programmed to spin down. Recovery time on the order of 30 seconds

Configuration examples by courtesy of Nexsan. (2005)

平成18年3月2日 DEWS2006 ミニサーベイ

まとめ

- 安心、快適、容易なデータ管理の要請
- ストレージシステムにおける技術革新
 - ストレージネットワーキング
 - 最近のストレージシステムの高度化技術

参考文献

平成18年3月2日 DEWS2006 ミニサーベイ

参考文献

(ストレージネットワーキング)

- 喜連川 優(編). ストレージネットワーキング. オーム社. 2002. 喜連川 優(編). ストレージネットワーキング技術—SNIAストレージ技術者認定プログラム準拠. オーム社. 2005.
- SNIA. A Dictionary of Storage Networking Terminology. http://www.snia.org/education/dictionary/.
- SNIA-J. SNIA用語集. http://www.snia-j.org/dictionary/.
- SNIA Technical Council. Shared Storage Model: A framework for describing storage architectures. 2003.
- Rob Peglar. Storage Virtualization I What, Why, Where and How? SNIA Tutorials. 2005.
- Rob Peglar. Storage Virtualization II Effective Use of Virtualization. SNIA Tutorials. 2005.
- Nik Simpson. ILM: Tiered Storage and the Need for Data Classification. SNIA Tutorials. 2005.
- Storage Management Initiative, SNIA. SMI-S: A Standard For Managing Storage. 2005.
- ファイバチャネル協議会(編). ファイバチャネル技術解説書. 論創社. 2001.
- JDSFファイバチャネル技術部会. ファイバチャネル技術解説書 II. 論創社. 2003.

参考文献 (データベースマシン)

- D. Slotnick. Logic per Track Devices. Advances in Computers 10, pp. 291-296. 1970.
- S. Su and G. Lipovski. CASSM: a cellular system for very large database. Proc. VLDB 1975, pp. 456-472. 1975.
- E. Ozkarahan, S. Schuster and K. Smith. RAP Associative Processor for Database Management. Proc. AFIPS 1975, pp. 379-387. 1975.
- D. DeWitt and P. Hawthorn. A Performance Evaluation of Database Machine Architectures. Proc. VLDB 1981, pp. 199-214. 1981.

平成18年3月2日 DEWS2006 ミニサーベイ

参考文献 (アクティブストレージ技術)

- E. Riedel, G. Gibson and C. Faloutsos. Active Storage for Large-Scale Data Mining and Multimedia. Proc. 24th VLDB 1998, pp. 62-73. 1998.
- A. Acharya, M. Uysal and J. Saltz. Active Disks: Programming Model, Algorithm and Evaluation. Proc. 8th ASPLOS 1998, pp. 81-91. 1998.
- K. Keeton, D. Patterson and J. Hellerstein. A Case for Intelligent Disks (IDISKs). SIGMOD Record 27(3), pp. 42-52. 1998.
- Erik Riedel, Christos Faloutsos, Gregory R. Ganger, David Nagle. Data Mining on an OLTP System (Nearly) for Free. Proc. SIGMOD Conference 2000: pp. 13-21. 2000.
- Christopher R. Lumb, Jiri Schindler, Gregory R. Ganger, David Nagle, Erik Riedel.
 Towards Higher Disk Head Utilization: Extracting "Free" Bandwidth from Busy Disk Drives. Proc. OSDI 2000, pp. 87-102. 2000.
- Christopher R. Lumb, Jiri Schindler, Gregory R. Ganger. Freeblock Scheduling Outside of Disk Firmware. Proc. FAST 2002, pp. 275-288. 2002.
- Jiri Schindler, John Linwood Griffin, Christopher R. Lumb, Gregory R. Ganger:. Track-Aligned Extents: Matching Access Patterns to Disk Drive Characteristics. Proc. FAST 2002, pp. 259-274. 2002.
- Steven W. Schlosser, Jiri Schindler, Stratos Papadomanolakis, Minglong Shao, Anastassia Ailamaki, Christos Faloutsos, Gregory R. Ganger. On Multidimensional Data and Modern Disks. Proc. FAST 2005, pp. 225-238. 2005.

参考文献 (コラボラティブストレージ技術)

- Xuhui Li, Ashraf Aboulnaga, Kenneth Salem, Aamer Sachedina, Shaobo Gao. Second-Tier Cache Management Using Write Hints. Proc. FAST, pp. 115–128. 2005
- Muthian Sivathanu, Vijayan Prabhakaran, Florentina I. Popovici, Timothy E. Denehy, Andrea C. Arpaci-Dusseau, Remzi H. Arpaci-Dusseau. Semantically-Smart Disk Systems. Proc. FAST 2003. 2003.
- Muthian Sivathanu, Lakshmi N. Bairavasundaram, Andrea C. Arpaci-Dusseau, Remzi H. Arpaci-Dusseau: Database-Aware Semantically-Smart Storage. Proc. FAST 2005, pp. 239–252. 2005.

平成18年3月2日 DEWS2006 ミニサーベイ

参考文献 (オートノミックストレージ技術)

- Zhenmin Li, Zhifeng Chen, Sudarshan M. Srinivasan, Yuanyuan Zhou. C-Miner: Mining Block Correlations in Storage Systems. Proc. FAST 2004, pp. 173-186. 2004.
- Zhenmin Li, Zhifeng Chen, Yuanyuan Zhou. Mining block correlations to improve storage performance. TOS 1(2), 213-245. 2005.
- Chris Gniady, Ali Raza Butt, Y. Charlie Hu. Program-Counter-Based Pattern Classification in Buffer Caching. Proc. OSDI 2004, pp.395-408. 2004.
- Ali Raza Butt, Chris Gniady, Y. Charlie Hu. The performance impact of kernel prefetching on buffer cache replacement algorithms. Proc. SIGMETRICS 2005, pp. 157-168, 2005
- Kai Shen, Ming Zhong, Chuanpeng Li. I/O System Performance Debugging Using Model-driven Anomaly Characterization. Proc. FAST 2005, pp. 309–322. 2005.

参考文献

(ストレージフュージョン)

- 向井 景洋, 根本 利弘, 喜連川 優. 高機能ディスクにおけるアクセスプランを用いたプリフェッチ機構に関する評価. 電子情報通信学会第11回データ工学ワークショップ (DEWS2000), 3B-2 2000
- 出射英臣, 茂木和彦, 西川記史, 大枝高. クエリプランを利用した先読み技術の開発と初期評価. 電子情報通信学会第16回データエ学ワークショップ/第3回DBSJ年次大会(DEWS2005), 5B-o1, 2005.
- 河島徹,河村信男,山口浩太,藤原真二,鈴木芳生,大枝高.ストレージのリモートコピー機能を利用したDBのディザスタリカバリ方式,第66回情報処理学会全国大会,2004.
- 合田和生, 喜連川優. データベース再編成機構を有するストレージシステム. 情報処理学会 論文誌データベース, 46(SIG 8,TOD 26), pp. 130-147. 2006.
- 喜連川優. 先進的なストレージ技術およびWeb解析技術. e-Societyシンポジウム2003講演 資料集. 2003.
- 喜連川優. 先進的なストレージ技術およびWeb解析技術. e-Societyシンポジウム2004講演 資料集. 2004.

平成18年3月2日 DEWS2006 ミニサーベイ

参考文献

(自律分散技術)

- Ben Y. Zhao, John Kubiatowicz, Anthony D. Joseph. Tapestry: a fault-tolerant wide-area application infrastructure. Computer Communication Review 32(1), p. 81. 2002.
- Sean C. Rhea, Patrick R. Eaton, Dennis Geels, Hakim Weatherspoon, Ben Y. Zhao, John Kubiatowicz. Pond: The OceanStore Prototype. Proc. FAST 2003. 2003.
- Haruo Yokota. Autonomous Disks for Advanced Database Applications. Proc. DANTE 1999, pp. 435-442. 1999.
- Daisuke Ito, Haruo Yokota: Automatic Reconfiguration of an Autonomous Disk Cluster. Proc. PRDC 2001, pp. 169-172. 2001.
- Svend Frølund, Arif Merchant, Yasushi Saito, Susan Spence, Alistair C. Veitch. FAB: Enterprise Storage Systems on a Shoestring. Proc. HotOS 2003, pp. 169-174. 2003.
- Yasushi Saito, Svend Frølund, Alistair C. Veitch, Arif Merchant, Susan Spence. FAB: building distributed enterprise disk arrays from commodity components. Proc. ASPLOS 2004, pp. 48-58. 2004.

参考文献 (省電力化技術)

- Dennis Colarelli, Dirk Grunwald, Michael Neufeld. The Case for Massive Arrays of Idle Disks. Proc. FAST 2002, WiP report. 2002.
- Copan Systems. http://www.copansys.com/.
- Fred Moore and Aloke Guha. Introducing COPAN Systems MAID architecture (Massive Arrays of Idle Disks). White Paper, Copan Systems and Horison Information Strategies. 2004.
- Nexan Technologies. http://www.nexsan.com/.
- Qingbo Zhu, Zhifeng Chen, Lin Tan, Yuanyuan Zhou, Kimberly Keeton, John Wilkes:. Hibernator: helping disk arrays sleep through the winter. Proc. SOSP 2005, pp. 177-190.
- E. Pinheiro, and R. Bianchini. Energy conservation techniques for disk array-based servers. Proc. ISC 2004, pp. 68-78. 2004.
 Sudhanva Gurumurthi, Anand Sivasubramaniam, Mahmut Kandemir, and Hubertus Franke. Reducing Disk Power Consumption in Servers with DRPM. Computer 36(12),
- H. Yada, H. Ishioka, T. Yamakoshi, Y. Onuki, Y. Shimano, M. Uchida, H. Kanno, and N. Hayashi. Head positioning servo and data channel for HDDs with multiple spindle speeds. IEEE Transactions on Magnetics, 36(5), pp. 2213–2215. 2000.

平成18年3月2日 DEWS2006 ミニサーベイ

謝辞

- 本講演資料を作成するにあたり、東京大学生産技術 研究所の喜連川優教授にご指導をいただきました。
- 本講演資料において紹介する「ストレージフュージョ ン技術」は、文部科学省リーディングプロジェクト「e-Society 基盤ソフトウェアの開発」に係る「先進的なス トレージ技術」よるものです。当該プロジェクトの推進 にあたっては、協力企業である株式会社日立製作所 より多くの有益なコメントを頂戴しました。