

グリッドコンピューティングの動向

佐藤三久

筑波大学

HPCS Lab.

並列分散システム 2005 1 High Performance Computing System Lab., Univ. of Tsukuba

グリッドコンピューティングとは

- グリッド技術とは広域の高速ネットワーク上において、**「安全に」**大量のデータ、計算資源、貴重な装置等を共有し、協調作業、資源の有効活用するネットワーク基盤技術(ソフトウェア、ネットワーク、ハードウェア)と、これを活用する応用技術

従来のインターネット技術
E-mail
www,
データの交換

→

グリッド
計算資源の
シームレスな共有
(CPU, Storage, ...)
計算資源の仮想化

Grid:
電力網(Power Grid)
のように計算資源
を使いたい

HPCS Lab.

並列分散システム 2005 2 High Performance Computing System Lab., Univ. of Tsukuba

グリッドコンピューティングとは


- なにが変わったのか？
 - 以前から、分散コンピューティングの研究はあった。
 - インターネットの急激な進歩 (bandwidth), 普及 (number of site)
 - 全世界で共通な基盤を作ろうとしている！ (Grid Forum, ..., ApGrid)
- 何につかえるのか？
 - 計算資源 (スーパーコンピュータ) を共有し、大規模計算を行う (meta-computing, MPI-G, ..., ITBL?)
 - 遠隔の計算資源と手元のPC/WSをシームレスに結合 (computing portal, GridRPC)
 - 大量のデータの処理 (data intensive computing, gfarm for ATLAS)
 - 高価な装置の遠隔共有 (電子顕微鏡, 衛星データ, 加速器?)
 - 共同作業のサポート (電子会議システム, AccessGrid)
 - 遊休の計算機の利用, SETI@home, P2P, ...

HPCS Lab.


並列分散システム 2005 3 High Performance Computing System Lab., Univ. of Tsukuba

グリッドの例


- To provide infrastructure and facilities needed for next major stages of collaborative research in:
 - genomics and bioscience
 - particle physics
 - astronomy
 - climatology
 - engineering design
 - social sciences
 - Medical Engineering




Bio Grid



VLBI, Kashima 34m telescope



JVO



Particle Physics
Large Hadron Collider at CERN
Detector for LHCb experiment
Detector for ALICE experiment

NEES Grid

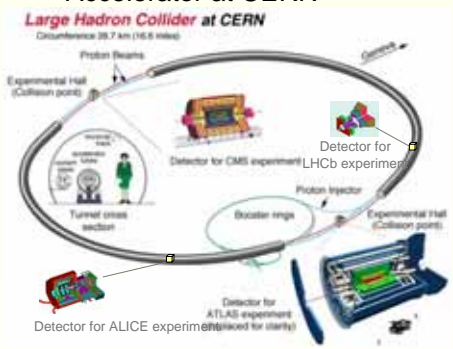
gfarmプロジェクト(産総研)

ペタスケールデータコンピューティング - Petascale Data Intensive Computing

- 大規模データ計算科学、データマイニング
 - 高エネルギー物理学、粒子物理学
 - 天文台、地球惑星
 - 生命情報工学...
- 大規模ビジネスデータベース
 - e-Japan、電子政府、電子商取引
 - データウェアハウス
 - 検索エンジン

HPCS Lab.

Example: Large Hadron Collider (LHC) Accelerator at CERN



HPCS Lab.

High Energy Physics Data Analysis

RAW ~1PB/year (1MB/event 30MB/sec)
 calorimeter-1 digits, calorimeter-2 digits, tracker-1 digits, tracker-2 digits, magnet-1 digits
 calorimeter reconstruction algorithm, track reconstruction algorithm, track reconstruction algorithm, magnetic field reconstruction algorithm

REC ~1PB/year
 calorimeter-1 energy, calorimeter-2 energy, tracker-1 position info, tracker-2 position info, magnet-1 field
 cluster reconstruction algorithm, track reconstruction algorithm

ESD ~300TB/year 100KB/event
 cluster-1 cluster-2 cluster-3, track-1 track-2 track-3 track-4 track-5
 electron identification algorithm, jet identification algorithm, Et miss identification algorithm

AOD ~10TB/year 10KB/event
 electron, photon, electron2, jet 1, jet 2, Et miss

Osaka-U Cybermedia Center (下条) GRID FS and App Development

Bioinfo App MEG Data Analysis
 DV Movie Data
 Remote Medicine
 High Data-rate Analysis
 Parallel Computing on Grid
 AIIST West Center Life Electronics Lab
 Centre for Multimedia and Network Technology NTU, Singapore
 Bioinformatics Institute Singapore

Grid-based Astronomical Analysis (NAO&国立天文台)

- Large fraction of astro-papers based on archives
- HST archive uses growing faster than archive
- Distributed database retrieval and analysis on Grid

Data Server interface
 Data Server (Mirror), Data Server (Mirror), Data Server, Data Server
 Tools Archive Server (Analysis, Mining, Visualization)
 Catalog Server - Pre-matched Catalog
 Data Analysis Servers
 Distributed Analysis GRID computing
 Common Query Language
 Upload
 USER

Already more retrieval than ingest!

graphics from US NVO project

グリッドのソフトウェア

- Globus: 最もmajorになりつつあるgridコンピューティングのツールキット(ミドルウェア)
 - ジョブ起動、セキュリティ、通信のためのライブラリ
 - globusrun: globus版のrsh
 - MPICH-G: globus版のMPICH
 - gftp: globus版のftp
- Nimrod-G(efusion): クラスタ、Grid向けのjob dispatcher. jobレベルのパラメータサーチを行うツール(GUI 付き)
- Condor-G: 遊休の計算機に対するジョブスケジューラ。プロセスの移送もサポート
- GridRPC: Grid上の遠隔手続き呼び出し
 - Ninf-G, NetSolve, OmniPRC

並列分散システム 2005 10 High Performance Computing System Lab, Univ. of Tsukuba

Globus

- Resource management ; GRAM
 - 計算資源の割り当て・ジョブ起動実行制御
- Communication Infrastructure Globus IO
 - 様々なプロトコルをサポートする通信レイヤ
- Metacomputing Directory Service MDS
 - Grid上のGlobusで利用できる計算資源の情報提供
 - LDAPを使った情報提供
- Globus Security Interface GSI
 - 認証などのセキュリティ機構
 - X509 Certificate-base の認証
 - Single-Sign-ON
- Heartbeat Monitor HBM
 - システムの状況モニタ
- Remote data access GASS (Globus Access to Secondary Storage)
 - ファイルへのリモートアクセスサービス
- executable management GEM
 - 実行ファイルの構築、転送

並列分散システム 2005 11 High Performance Computing System Lab, Univ. of Tsukuba

グリッドと並列アプリケーション

- 典型的なアプリケーション
 - パラメータ検索: 同じ計算を膨大な計算資源で実行
 - master-slave型の並列プログラム
- 典型的なグリッド計算資源
 - 複数のクラスタが利用可能
 - 計算資源の状況が動的に変化

Our target

PC cluster, PC cluster, PC cluster

並列分散システム 2005 12 High Performance Computing System Lab, Univ. of Tsukuba

特定領域研究 A05-01: 計算物理学分野のGridアプリケーションと並列プログラミングシステムの研究

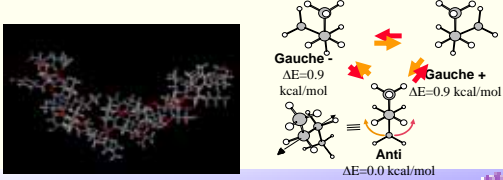
- グリッド並列プログラミングミドルウェアOmniRPCの開発
 - グリッドでの Master-workers 型の並列プログラムをサポート
 - Gridのend-pointの計算資源として、クラスターを対象とする
 - プライベートアドレスのクラスターについてもサポート
 - Firewall内のクラスターにもアクセス可能
 - 既存のプログラムに対し、簡単な並列プログラミング環境を提供
 - 並列プログラミングの記述、制御はOpenMPで行うことができる
 - クラスターでは "ssh", globus環境では "GRAM", "ssh" 環境でも使用できる、
 - 各計算資源の管理ポリシーを考慮したジョブの起動をサポート
 - サイトのスケジューラ (PBS, SGE, ...)
 - 大規模なグリッド環境 (upto 1000 hosts!) のサポート
 - http://www.omni.hpcc.jp/OmniRPCで公開中
- グリッドアプリケーション
 - CONFLEX-G: 網羅的分子探索プログラムのグリッド並列化
 - グリッド上に分散している大規模なクラスター計算資源を利用して、計算が可能
 - OmniRPCで複数のクラスターで実行、有効性を検証
 - HMCS-G (Grid-enabled Heterogeneous Multi-computer System)
 - 広域ネットワーク環境において、貴重な計算リソースである重力専用計算機GRAPE-6を共有し、汎用計算機と融合計算を行う。
 - OmniRPCを用いて、セキュリティ、認証、プログラミングインタフェースを提供

HPCS Lab.

並列分散システム 2005 13 High Performance Computing System Lab., Univ. of Tsukuba

CONFLEX-G: Grid-enabled Molecular Conformational Search Program

- JST-ACTによるプロジェクト
- CONFLEX
 - 分子の配座(3次元構造)を網羅的に計算するプログラム
 - 豊橋技術科学大学 後藤らが開発
 - 同じ分子でもエネルギー状態が異なる3次元構造をたくさん持つ
 - 大量の構造を高速に配座探索する必要性

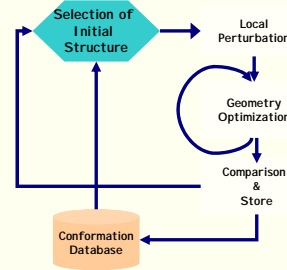


HPCS Lab.

並列分散システム 2005 14 High Performance Computing System Lab., Univ. of Tsukuba

CONFLEXの配座探索アルゴリズム

- 分子の配座空間を探索するプログラム
 - 分子が取り得るすべての立体配座を自動的に発生させ、化学的に重要な配座異性体の最適化構造を探索
 - 2D drug database 3D structure database
- 分子力学 (Molecular Mechanics) を使い構造最適化を行う
 - 全原子に対応した計算が可能
- 配座発生と構造最適化の機能を合わせ持つ
 - 配座発生時に使用するアルゴリズムにより優れたパフォーマンス
 - Reservoir-Filling (貯水池注水) アルゴリズム



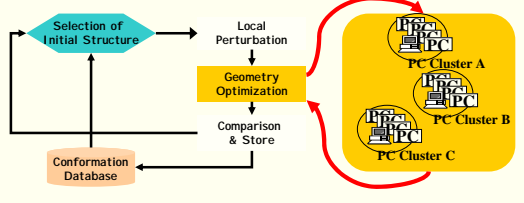
すでに保存されている構図の中から初期構造の選択
 試行構造の生成
 それぞれの構造を最適化
 すでに得られている構図と比較し、新しい配座を保存

HPCS Lab.

並列分散システム 2005 15 High Performance Computing System Lab., Univ. of Tsukuba

CONFLEX-G: Grid enabled CONFLEX

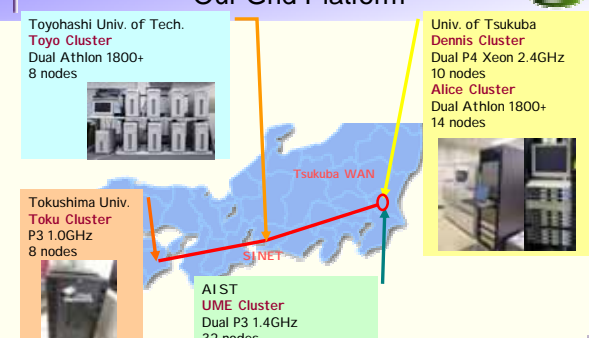
- 全体の処理の90%以上を占める構造最適化の処理を Master/Worker パラダイムで並列化
- 巨大な計算機資源を利用することができる
- OmniRPCのモジュール再初期化機能を使用
 - RPC呼び出しごとに初期化が不要



HPCS Lab.

並列分散システム 2005 16 High Performance Computing System Lab., Univ. of Tsukuba

Our Grid Platform



- Toyohashi Univ. of Tech. **Toyo Cluster**
Dual Athlon 1800+
8 nodes
- Tokushima Univ. **Toku Cluster**
P3 1.0GHz
8 nodes
- AIST **UME Cluster**
Dual P3 1.4GHz
32 nodes
- Univ. of Tsukuba **Dennis Cluster**
Dual P4 Xeon 2.4GHz
10 nodes
- Univ. of Tsukuba **Alice Cluster**
Dual Athlon 1800+
14 nodes

HPCS Lab.

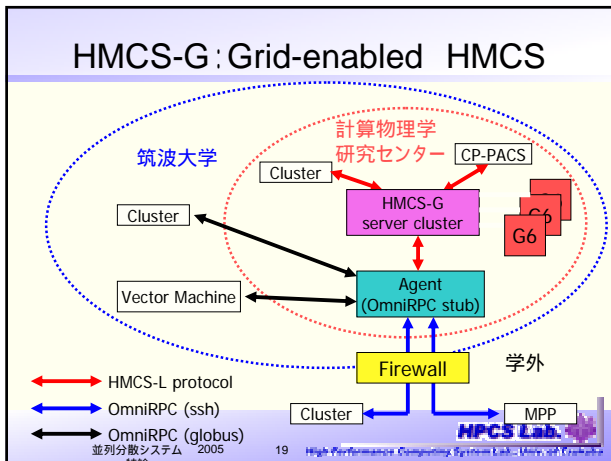
並列分散システム 2005 17 High Performance Computing System Lab., Univ. of Tsukuba

Heterogeneous Multi-Computer System

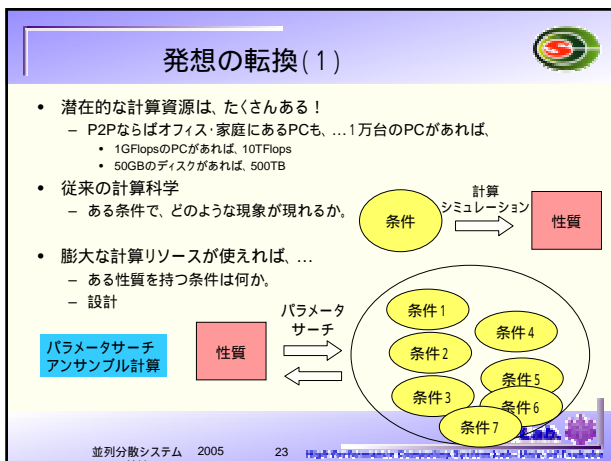
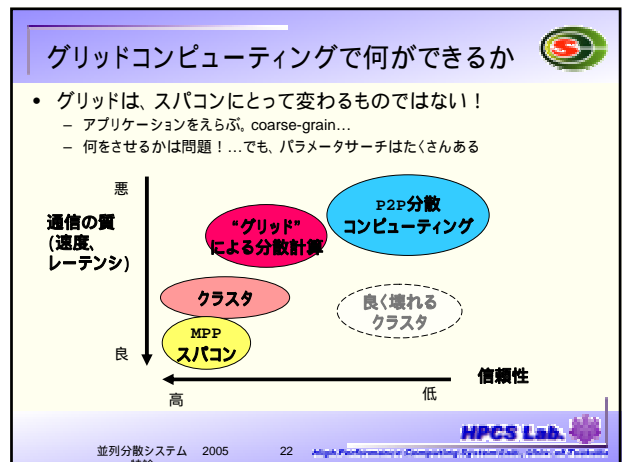
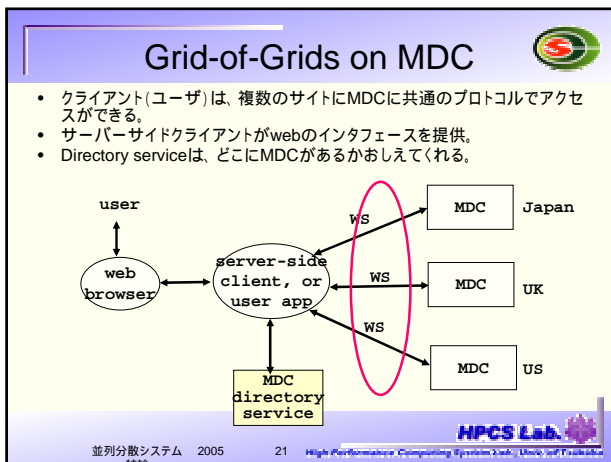
- 計算物理学のための計算性能は...
 - 最先端の計算物理学においては、詳細で精密なシミュレーションのために各種物理現象を達成困難として扱う必要がある
 - それらの中には大規模計算のために膨大な演算量を必要とするものがある
 - FFT: $O(N \log N)$
 - 重力、分子動力学法: $O(N^2)$
 - ナノスケール物性物理: $O(N^3)$, $O(N^4)$, ...
 - 通常の汎用計算機は多くの場合十分な計算性能を提供しない
 - 何らかの形で専用計算機を導入する必要がある
- 粒子系シミュレーション (例: 重力計算) と 連続系シミュレーション (ex: 場・流体力学) を一つのシステムで統合化
- 汎用プロセッサ (柔軟性) と専用プロセッサ (高速性) を統合
- HMCS: GRAPE-6 (重力計算専用機) + 汎用並列システム
- HMCS-G: Grid-enabled version!

HPCS Lab.

並列分散システム 2005 18 High Performance Computing System Lab., Univ. of Tsukuba



- ### データグリッドへの取り組み: ILDG
- ILDG (International Lattice Data Grid)
 - 素粒子物理学分野でのデータグリッド
 - 計算されたconfigurationのデータの統一
 - configurationのデータを共有
 - (Lattice theoryの計算のための資源の共有)
 - Working Group
 - Metadata WG: 計算されたconfigurationの記述フォーマット(QCDML)の策定
 - Middleware WG: Webserviceを利用して、configurationのデータの共有のための仕組みを提供。
 - Metadata Catalog (MDC)
 - Replica Management
 - Storage Resource Management (SRM)
- HPCS Lab. 並列分散システム 2005 20



- ### 発想の転換 (2)
- グリッドコンピューティングの計算の性質に適した新しいアルゴリズム、発想が必要！
 - 例: Folding@home
 - スタンフォード大学のビジャ・パンデ助教授 (<http://www.stanford.edu/dept/chemistry/faculty/pande/>) らが開発した“確率分割法”
 - 一般的なMD計算では、たん白質を周期境界条件で区切って“空間分割”によって並列化、この手法ではCPU間通信がボトルネックとなり、並列度をあまり上げられない
 - “確率分割法”では、各CPUに薬物分子とたん白質とのクラスター構造を与え、周囲の水分子の状態などの計算条件を変えたシミュレーションを並列で実行させる。(パラメータサーチ！)
 - folding問題も、グリッド(P2P)で可能になった！？
- HPCS Lab. 並列分散システム 2005 24

データグリッドの動向



- データグリッド
 - 様々な電子的なデータが蓄積しつつある
 - 加速器
 - 衛星データ
 - 天体観測
 - センサーデータ
 - ゲノム
 - 例: Gfarmによる天文データの解析
 - Gfarm: 産業総合研究所で開発されたデータグリッドのためのミドルウェア
 - 膨大な天体観測の画像データを処理し、解析
 - データマイニング技術による発見
 - 分散されたデータの統合
- データグリッドの成功のためには、分野のコミュニティのデータの共有利用のための国際的な協力が必要
 - ILDG

HPCS Lab.

並列分散システム 2005

25

High Performance Computing System Lab., Univ. of Tsukuba

グリッド技術の最近の動向



- Gridの本質は「安全なリソースの共有」
 - Globusの重要なところは、PKI (X509ベースの認証)
- GGF (Global Grid Forum)の変質(?)
 - Scientific Grid (VO)からビジネスグリッドへの変質
 - Web Serviceベースに移行
 - OGSAからWSRF(WS Resource Framework)へ
- アプリケーションは確実に増えつつある
 - データグリッドは重要
 - パラメータサーチは重要
- エンタープライズ・グリッド
 - GridMP(United Devices), ...
 - 社内需要はある。
- P2Pは?
 - CS的にはこちらのほうが面白い
- これからの興味のある課題は、resource finding

HPCS Lab.

並列分散システム 2005

26

High Performance Computing System Lab., Univ. of Tsukuba

グリッドとは



- 認証と利用
 - PKI (Public Key Infrastructure)によるセキュリティ
 - Single-sign ON
- 資源共有
 - コミュニティで、データを共有するための標準化(XML)
 - Web Serviceなどの共通のインタフェースを作る
- 資源発見
 - これは、期待されているがあまりない。
 - UDDI (Universal Description, Discovery and Integration) ?
- 大規模
 - 大規模な計算資源の統合が可能
 - 大規模なデータが蓄積が可能

HPCS Lab.

並列分散システム 2005

27

High Performance Computing System Lab., Univ. of Tsukuba