



スーパーコンピュータRICCの現状と取組みについて

理化学研究所 情報基盤センター

野田 茂穂





RICCの紹介と運用状況





RICC設計コンセプト

- I. 次世代スーパーコン ピュータを活用するアプ リケーションの開発環境
- II.新しい計算機技術の利用
- III.加速器・次世代シーケンサ などの実験データのストレー ジ・解析



数千並列規模の並列ジョブ実行 が可能



100ノードにグラフィックボードを 利用したアクセラレータを搭載



4PBのテープアーカイブ装置と 500TBのディスク装置が実験結果 の大規模データに対応

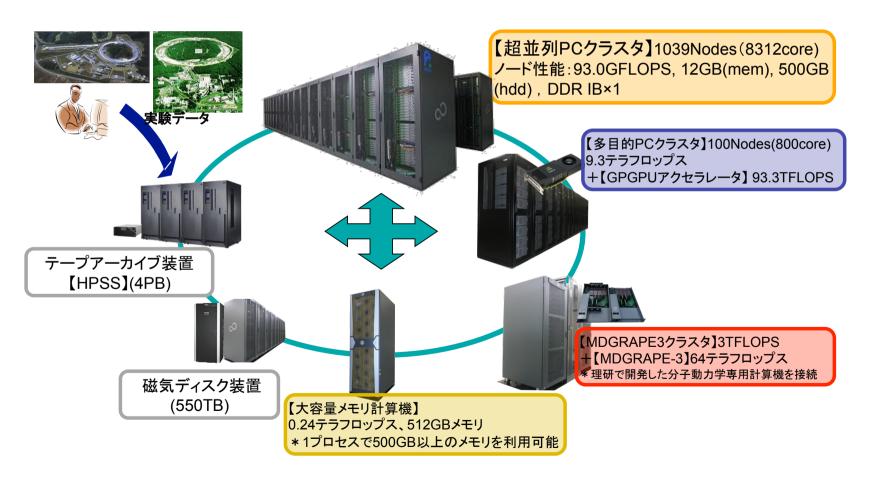




RICCの概要

【システム構成】

超並列PCクラスタ+GPUクラスタ+専用機クラスタ+大容量メモリ計算機







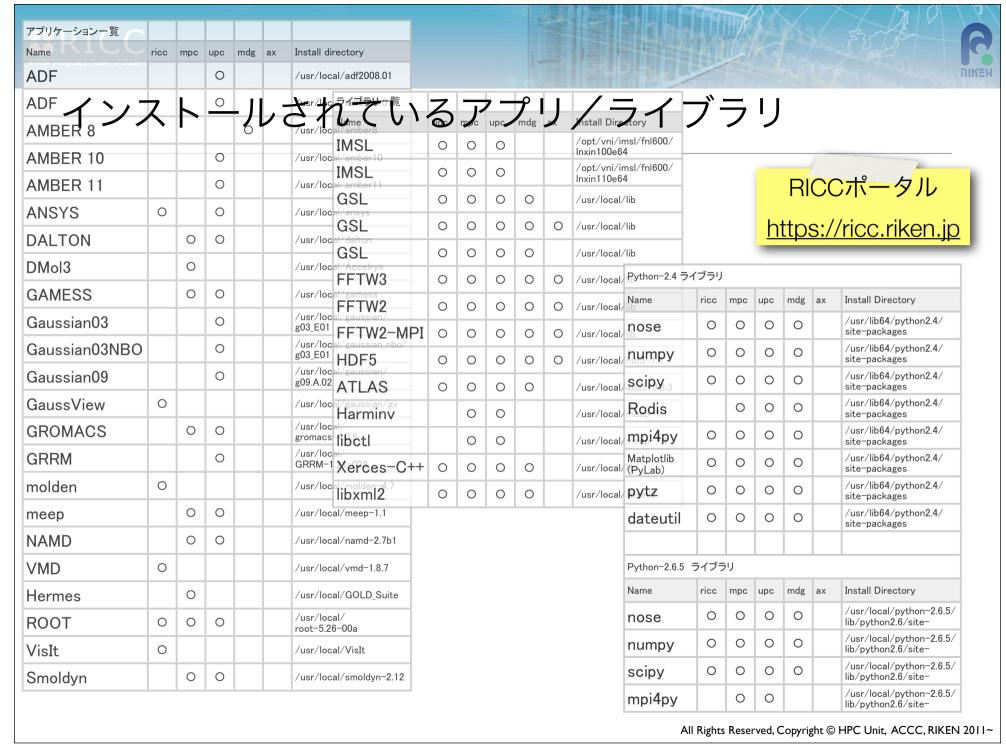
各コンポーネントの特長 (利用手引書(簡易版): P2~3)

*****計算ノード

- 超並列PCクラスタ MPC (8000->8192)
 - →大規模並列計算用途
- 多目的PCクラスタ UPC (800)
 - ➡商用・フリーアプリケーション実行
 - ➡GPUアクセラレータを搭載
- インタラクティブジョブクラスタ (32)
- 占有利用課題クラスタ (160)
- MDGRAPE-3クラスタ MDG
 - →分子動力学専用計算機
- 大容量メモリ計算機 AX
 - ➡1プロセスで大量メモリを必要とする 計算向け

*****フロントエンドシステム

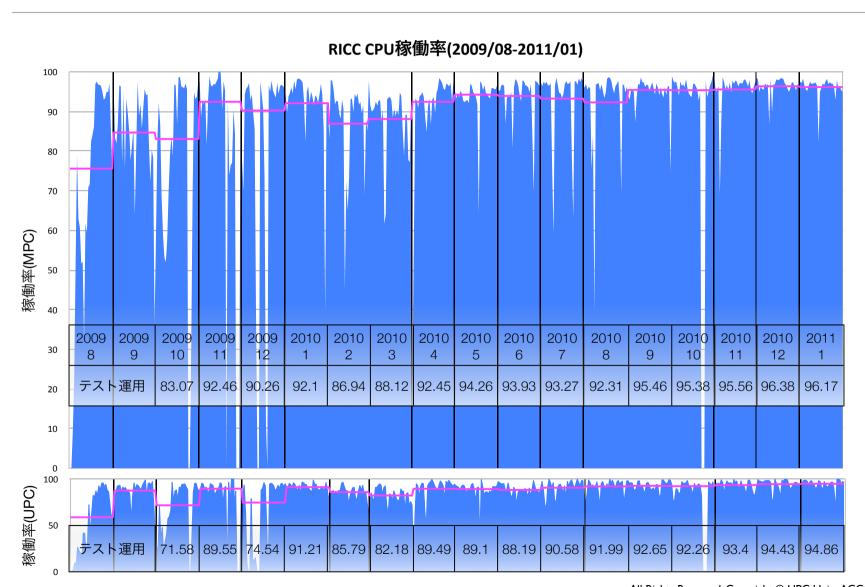
- 磁気ディスク装置
 - ➡ホーム領域(/home): 30TB (1ユーザあたり最大500GB)
 - ➡データ領域(/data) : 300TB
- テープアーカイブ装置(2PB)







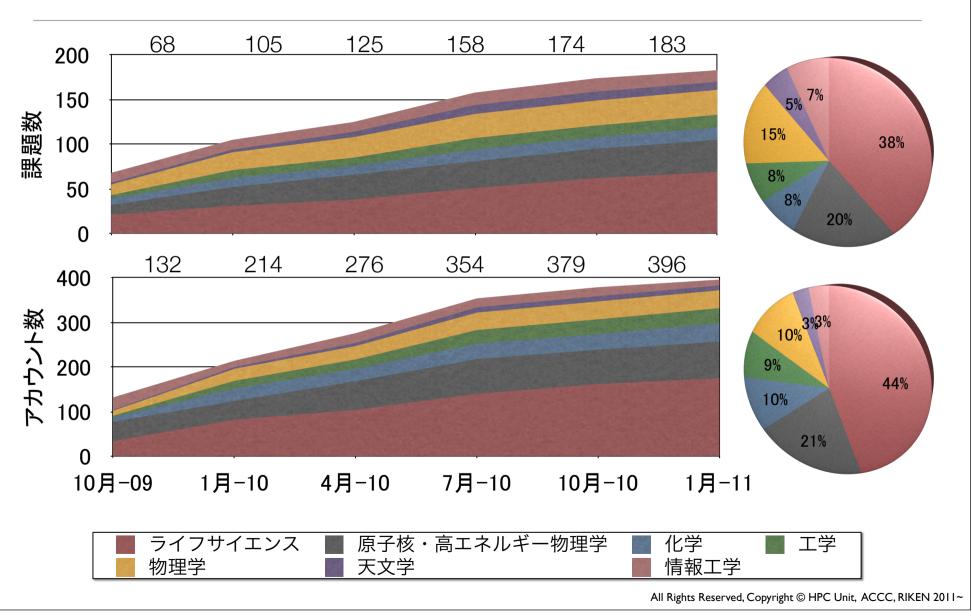
稼働率





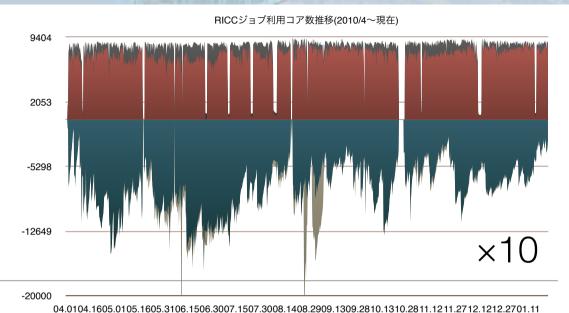


課題数 アカウント数









運用状況まとめ

PCクラスタが95%を超える稼働率

→常時約50000コア程度の待ちジョブがある

課題審査の結果を元にジョブ実行を制御

- →独自に開発したMJS(メタジョブスケジューラ)
- →審査で高い評価が得られている課題は希望通りの演算が実行できるように制御 高並列のジョブは週末運用(※後述)を用意





情報基盤センター取り組みについて





RICC設計コンセプト

I. 次世代スーパーコン ピュータを活用するアプ リケーションの開発環境

| II.新しい計算機技術の利用

III.加速器・次世代シーケンサ などの実験データのストレージ・解析



数千並列規模の並列ジョブ実行 が可能



100ノードにグラフィックボードを 利用したアクセラレータを搭載

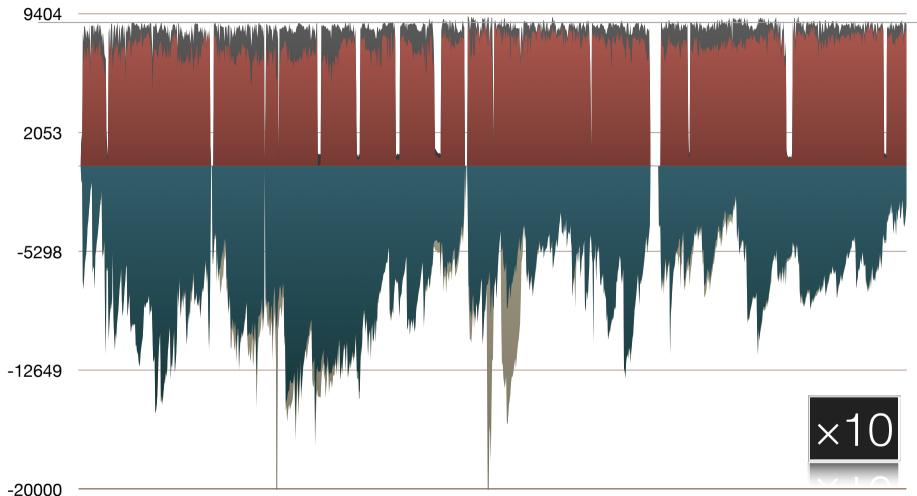


4PBのテープアーカイブ装置と 500TBのディスク装置が実験結果 の大規模データに対応





実行と待ちJobの推移



04.01 04.16 05.01 05.16 05.31 06.15 06.30 07.15 07.30 08.14 08.29 09.13 09.28 10.13 10.28 11.12 11.27 12.12 12.27 01.11

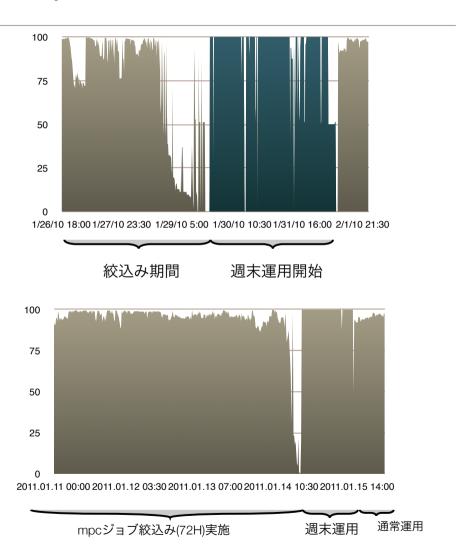
- ●PCクラスタが95%を超える稼働率
- ●→常時約50000コア程度の待ちジョブがある





大規模並列ジョブのための週末運用

- ◆常時高並列ジョブがサブミットできる が高い稼働率のため現実的に通常運用 中の実行は困難
 - 大規模並列ジョブをまとめて実行す る日を設ける
- 利用区分に関わらずどの課題でも参加 できる
 - 目的:高並列なプログラム開発や チューニング
- 2010年1月から開始(1~2回/ 月のペース)







参加状況

2010年度は12回実施

参加課題数:47件(13課題) (1月まで)

各回の平均演算時間 179,560コア時間(1月まで)

→平均すると30時間程度

主な参加目的

- ・スケーラビリティの調査
- ・プログラムの高並列実行時の動作確認





参加課題の例

- タンパク質・核酸など生体高分子の分子シミュレーション(木寺 詔紀、松永 康佑、森次 圭、他)
 - タンパク質の構造変化パスをサンプリングするプログラムのベンチマーク
- 大規模遺伝子ネットワーク推定プログラムの研究開発(宮野悟、玉田 嘉紀)
 - マイクロアレイデータなどの遺伝子発現データから遺伝子間の 発現制御ネットワークを推定・予測するプログラム.
- 生体高分子生化学的機能解析のための分子計算技術の開発(木寺 詔紀、中村春木、他)
 - 開発中の量子化学計算プログラム Platypus-QM のRDFT計算機能の並列性能測定
 - 異なるスケールモデル(今回は生体高分子の全原子モデルとより粗視化した残基レベルモデル)の分子動力学シミュレーションを連成した計算(MSES法)





RICC設計コンセプト

I. 次世代スーパーコン ピュータを活用するアプ リケーションの開発環境

II.新しい計算機技術の利用

III.加速器・次世代シーケンサ などの実験データのストレージ・解析



数千並列規模の並列ジョブ実行 が可能



100ノードにグラフィックボードを 利用したアクセラレータを搭載



4PBのテープアーカイブ装置と 500TBのディスク装置が実験結果 の大規模データに対応

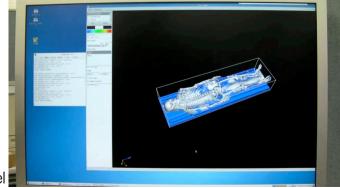




アクセラレータの利用状況と利用推進

- 状況
 - アクセラレータで10倍程度の加速が可能
 - 既にプログラム移植が済んでいる研究
 - 第一原理計算による機能性物質としての分子結晶の理論解析
 - 分子動力学計算商用アプリ (Amber)
 - アクセラレータ向けアルゴリズムの研究
 - LSV: 次世代スパコンの計算結果可 視化のための大規模可視化ソフトの

開発



- 現在はアクセラレータの稼働率は低い
- アクセラレータを活用できるアプリケー ションの整備が課題
 - プログラムを容易にGPGPU化出来る環 境の整備
 - PGI等のコンパイラの導入
 - GPGPUライブラリの開発と利用ツー ル(RIVER)の整備
 - GPGPUプログラミング教育の実施
 - プログラムのGPGPU移植作業の実施



Test LSV:107GB Voxel





RIVER (RIKEN IBM Visual Programing Environment)





メニーコア化・ハイブリッド化の流れ

- CPUの多コア化、および多量のコアを備えたアクセラレータの混載により、HPCのピーク性能は大幅に向上した
 - すでにハイブリッド・システムが、Top 500の上位を多数占めている
- さらなるメニーコア化、ハイブリッド化の進行が予想される

RICC



Xeon 4 コア + GPGPU 240コア

Roadrunner (IBM)



Opteron 2 コア + PowerXCell 9コア

TOP 10 Systems - 11/2010

- Tianhe-1A NUDT TH MPP, X5670 2.93Ghz 6C, NVIDIA GPU, FT-1000 8C
- 2 Jaguar Cray XT5-HE Opteron 6-core 2.6 GHz
- Nebulae Dawning TC3600 Blade, Intel X5650, NVidia Tesla C2050 GPU
- TSUBAME 2.0 HP ProLiant SL390s G7 Xeon 6C X5670, Nvidia GPU, Linux/Windows
- Hopper Cray XE6 12-core 2.1 GHz
- Tera-100 Bull bullx super-node S6010/S6030
- Roadrunner BladeCenter QS22/LS21 Cluster, PowerXCell 8i 3.2 Ghz / Opteron DC 1.8 GHz, Voltaire Infiniband
- 8 Kraken XT5 Cray XT5-HE Opteron 6-core 2.6 GHz
- JUGENE Blue Gene/P Solution
- Cielo Cray XE6 8-core 2.4





直面するプログラミングの困難

● メニーコア化、ハイブリッド化が進むとプログラミングが課題になる

クスク並列化 〉

処理AはCPUで、処理BはGPUで・・・

キャッシュサイズ、命令セットの違いを意 識した最適化・・・

7 7

通信の実装

ごとの最適化

• CPU同士はMPIで、GPUとは専用の転送で・・・

CPUの演算とGPUの演算と通信は同時に実行・・・

計算リソースへの割り当て

・処理AはCPUでN並列、処理BはGPU だからM並列で・・・



性能と生産性を両立させる新しい枠組みが必要





ハイブリッド・システムプログラミング環境 "RIVER" 2レベルプログラミング 自動並列化

コンポーネント・プログラミング

A.c D.cpp C.c

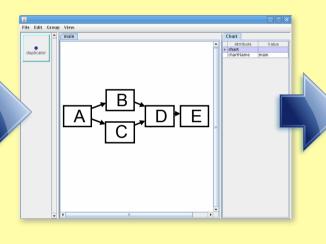
E.c

B.f95 C.cu

タスク並列化

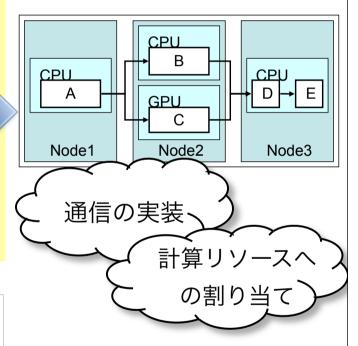
アーキテクチャでとの最適化

熟練プログラマーにより作成したものを共有 たものを共有 シングルスレッド用に書かれた 既存の言語・資産を流用可能 追加はわずかな指示文だけ データ・フロー・プログラミング



一般プログラマーは、アルゴリ ズムのデータフローを記述する だけ

最適化された部品ライブラリと (左)、利用するシステムに最 適化された並列化(右)を自動 的に利用できる コンパイル



コンパイラーが、プロセッサーご との処理能力とシステム構成を 考慮して処理の割り当てと並列 度を判断

通信コードも自動的に挿入

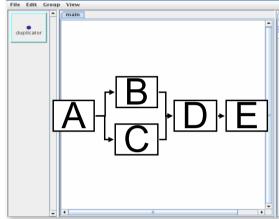




RIVERを利用したプログラミングのイメージ

1. アルゴリズムのデータフローを記述する



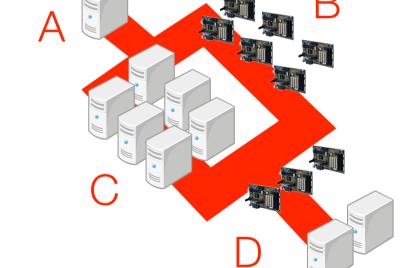


計算機利用者

3. コンパイラーが、プログラムを生成

処理内容とコードライブラリを参照して、最適 な計算ノード、並列数を決定。

必要なプログラム部品、通信コードを含んだ 実行形式を、ノードごとに生成。(MPMD)



2. システム構成を与えてコンパイル

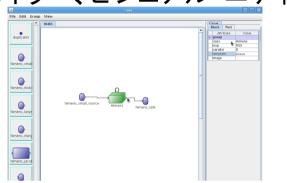




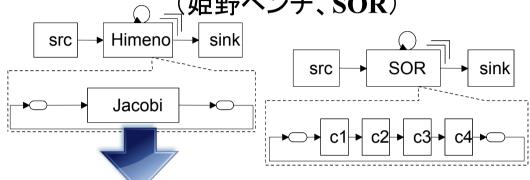


2010年度の成果

ープログラミング環境の開発 (コンパイラー、ビジュアル・エディタ)



部品ライブラリ開発、性能実験 (姫野ベンチ、SOR)

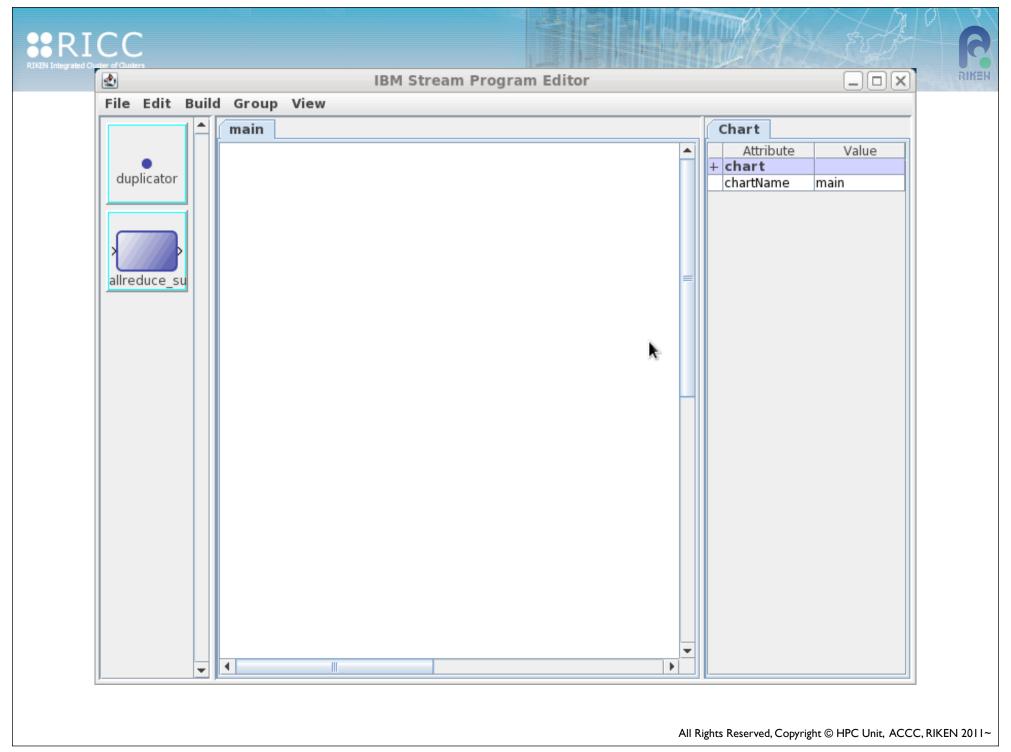


HimenoBMT (M size = 128 x 128 x 256) での割り当て例

実験環境

	x86	GPGPU
言語	С	С
コンパイラー	icc	nvcc
並列度	1-80	1-80
コア数またはカード数	1-80コア	1-80カード
使用ノード数	10	80
特記事項	8コア/ノー	CUDA for C

	スケジューリング結果	ホスト数
っ 2ホストは通信コスト	GPGPU x 1	1
が高く1並列で実行	GPGPU x 1	2
	GPGPU x 4	4
	GPGPU x 8	8
	GPGPU x 16	16
	GPGPU x 16	32
し 16ホスト以上は頭 うちと判断	GPGPU x 16	64
Copylight © HPC Unit, ACCC, RIKEN 2011~	GPGPU x 16	80





自動並列化・スケジューリング結果



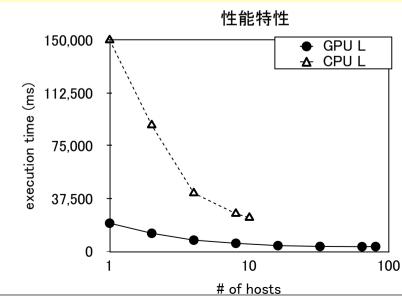
実験結果 HimenoBMT (L size = 256 x 256 x 512)

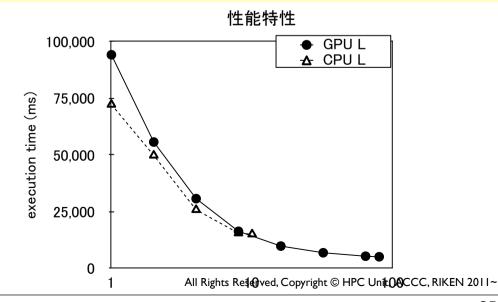
SORBMT (L size = $512 \times 512 \times 512$)

使用可能ホスト数	スケジューリング結果	
1	GPGPU x 1	
2	GPGPU x 2	
4	GPGPU x 4	
8	GPGPU x 8	
16	GPGPU x 16	
32	GPGPU x 32	
64	GPGPU x 64	
80	GPGPU x 80	

スケジューリングミ ス;最速構成は GPGPU x 64

スケジューリング結果	
x86 x 8 コア	
x86 x 16コア	
x86 x 32コア	
x86 x 64コア	
GPGPU x 16	
GPGPU x 32	
GPGPU x 64	
GPGPU x 80	









ロードマップ

2010年度

基幹ツールの開発

2011年度

RICCでのパイロット運用

言語仕様

目標

部品化指示構文の定義

ビジュアル エディター

プロトタイプ

コンポーネント ライブラリ

姫野ベンチ

SORベンチ

コンパイラー (並列化スケジュー ラー)

CPUおよびGPUに対応し た自動スケジューリング

通信コードの自動生成 (MPI)

2011/04公開予定

Eclipseプラグイン化

デバッグ支援機能

実用的なアプリケーション

通信隠蔽

計算能力にあわせ た自動領域分割

MPIとスレッドの 通信ハイブリッド化





RICC設計コンセプト

I. 次世代スーパーコン ピュータを活用するアプ リケーションの開発環境

|||.新しい計算機技術の利用

III.加速器・次世代シーケンサ などの実験データのストレー ジ・解析



数千並列規模の並列ジョブ実行 が可能



100ノードにグラフィックボードを 利用したアクセラレータを搭載



4PBのテープアーカイブ装置と 500TBのディスク装置が実験結果 の大規模データに対応





開発中:データデポジトリシステム(7月開始予定)

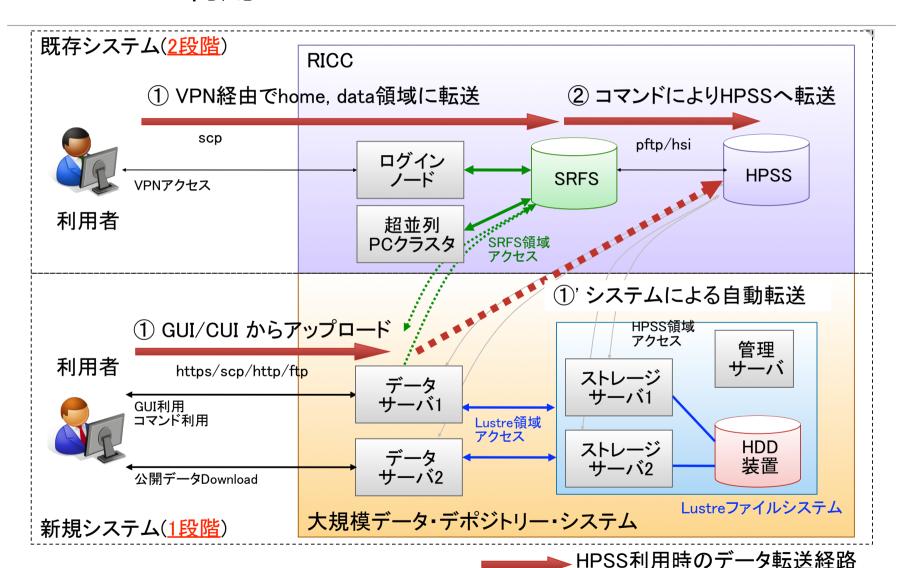
概要:

- ・大量の実験データが生成され、ストレージ容量が不足
 - →演算資源を使わないがHPSSテープアーカイブのみを利用
- ・HPSSを利用するには複数のコマンド実行が必要
 - →GUIを開発する事により、利便性の高いシステムを開発





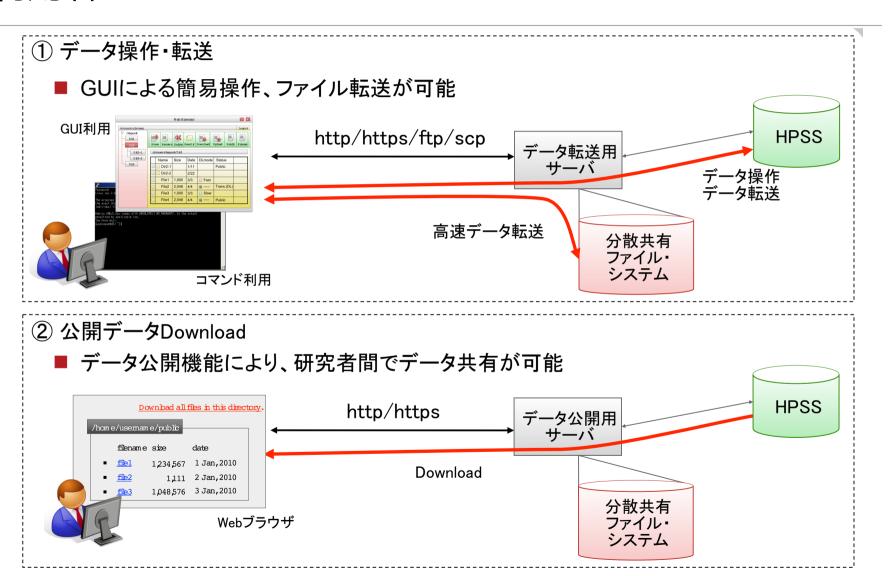
システムの利用イメージ







利用者アクセスイメージ







その他の取り組み

チューニング

单体性能、並列性能

講習会(RICCユーザ向けGPGPU講習会 2/18、一般向けも計画中)

研究の状況に応じた柔軟な対応

緊急の場合、一時的にJobの優先度をコントロール。

簡易利用でも一時的に使用可能core数制限を解除

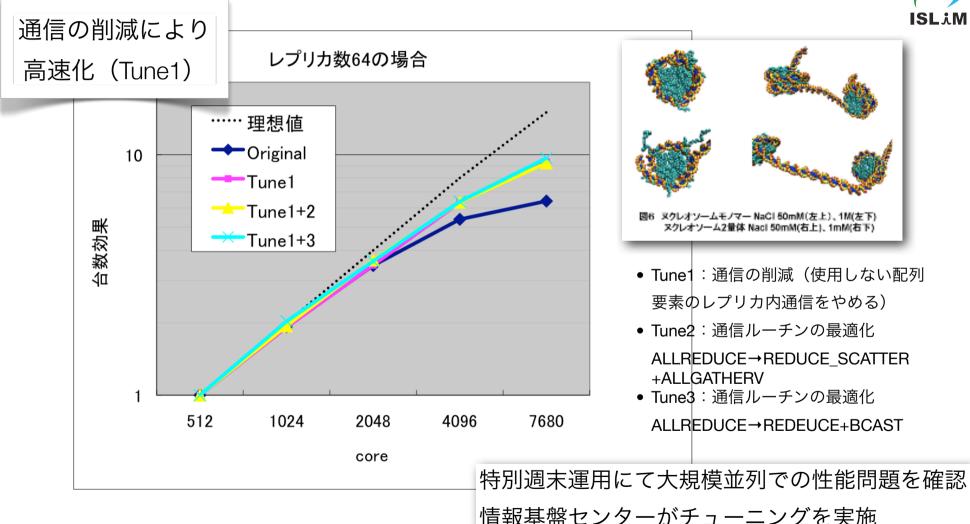
アプリケーション/ライブラリのインストール





チューニング支援の例(粗視化モデル計算CafeMol)









まとめ

RICCは非常に高い稼働率

- →より多くのジョブが実行できるよう運用に努める
- →プログラムのチューニング支援により効率的な実行を支援
- →アクセラレータの利用を進めて行く

ご意見、ご要望等:<u>hpc@riken.jp</u>