

将来(1ペタフロップス超・2010年代前半頃を想定)の 超高速計算機に必要な要素技術の研究開発について

平成16年8月20日

富士通株式会社



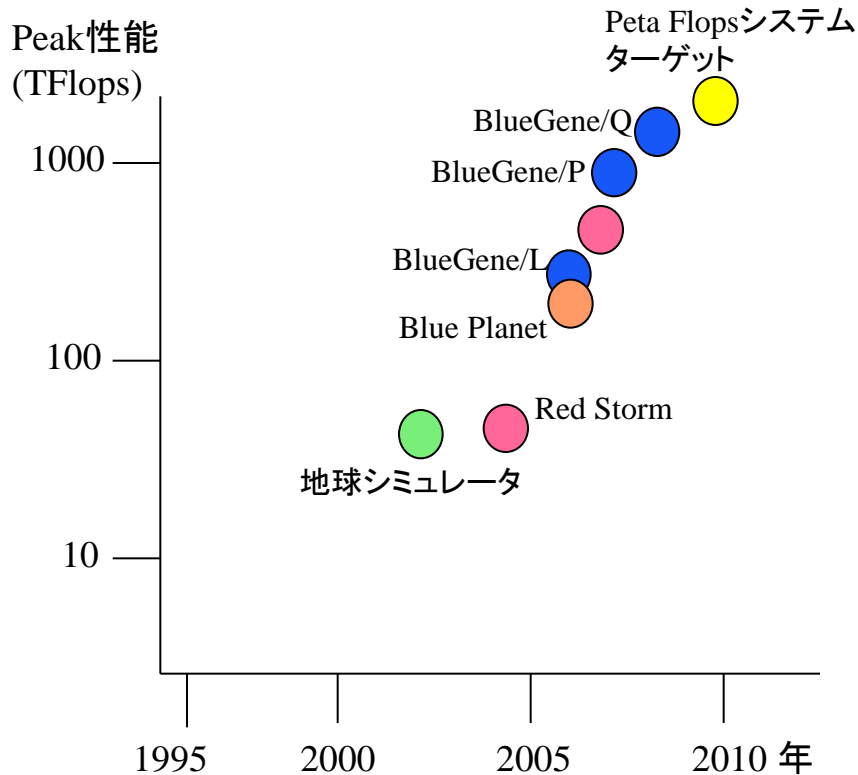
1. 将来の超高速計算機システムについて

1PFlops超高速計算機(ペタフロップス)システム実現へ向けての挑戦

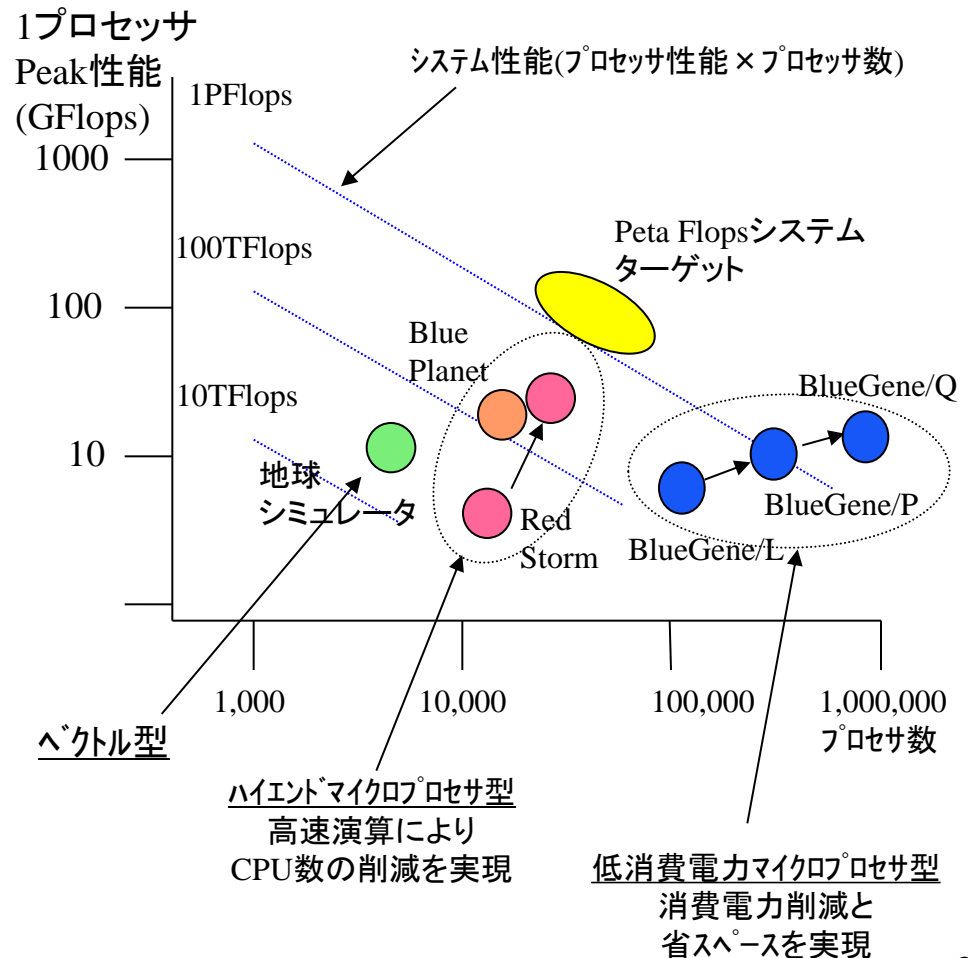
- 実効1PFlopsを引き出すアーキテクチャの開発
 - **性能=プロセッサ+インターコネクト+ソフトウェア開発環境**
 - 実効1PFlopsを実現するには、ピーク性能で3PFlops以上が必要
 - 数千ノード以上の並列実行には、高バンド幅、低遅延のインターコネクトが必要
 - アプリケーションから並列性を引き出すコンパイラ、ソフトウェア開発環境が必要
- 大規模システムの実現を可能とする技術の開発
 - **高密度化**: コンパクトなシステム
 - **省電力**: プロセッサ/システム全体としての省電力化
 - **半導体技術**: 微細加工技術, 低消費電力技術
 - **高信頼**: 障害検知, 切り離し, 交換, 自動復旧(ホットスワップ)
 - **運用管理システム**: 超高速計算機システムを効率よく活用, 運用する技術

HPCシステムのトレンドと開発ターゲット

＜HPCシステムの性能トレンド＞



＜プロセッサ数 vs. 単体性能＞

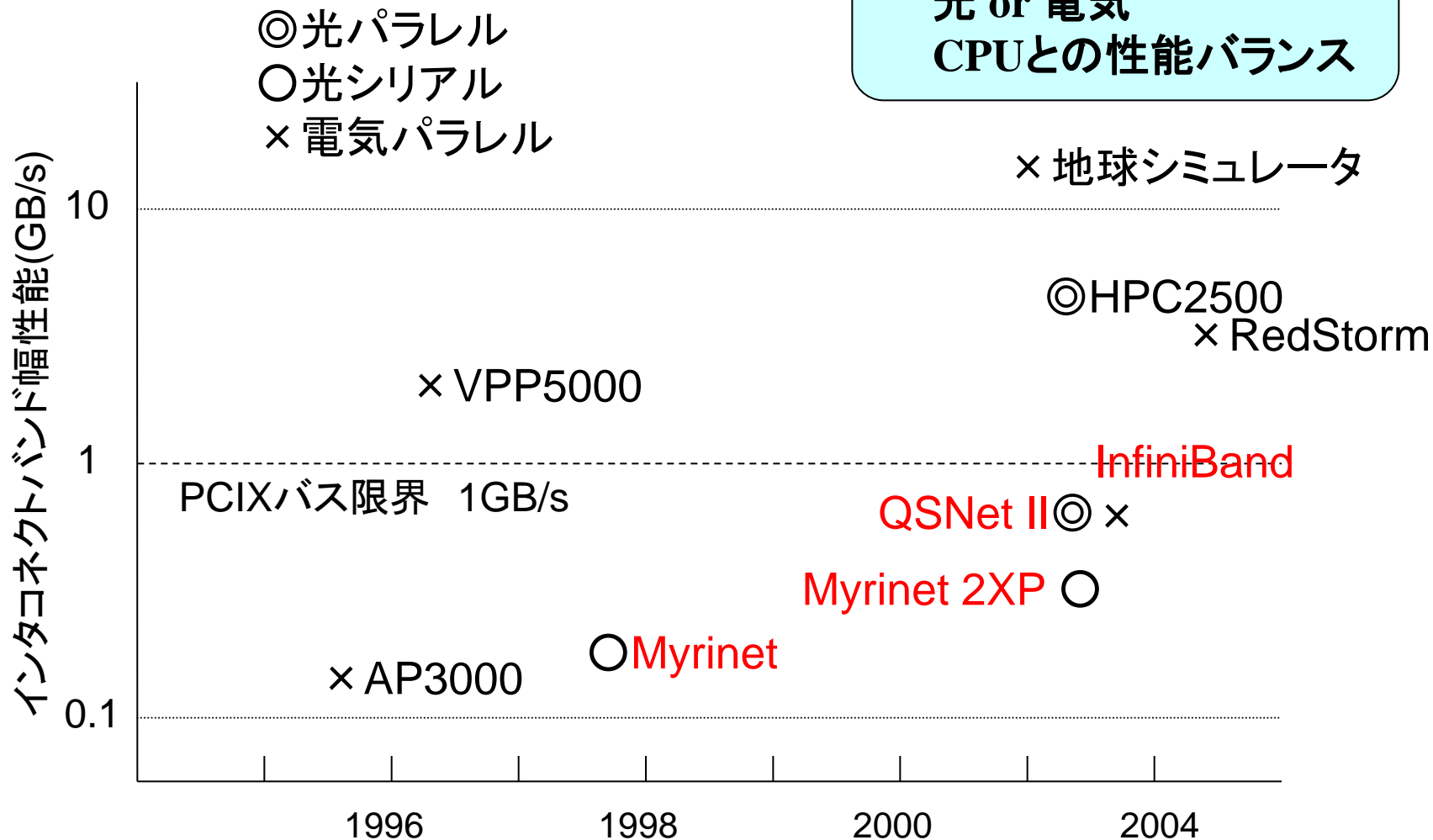


2010年のターゲット
ピーク性能3PFlops, 実効1PFlops

HPC向けインタコネクットの性能トレンド

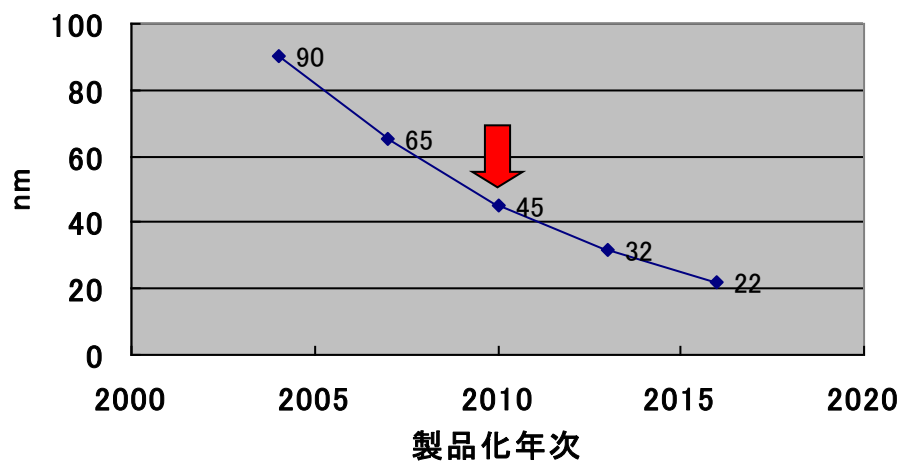
インタコネクットのチャレンジ

光 or 電気
CPUとの性能バランス

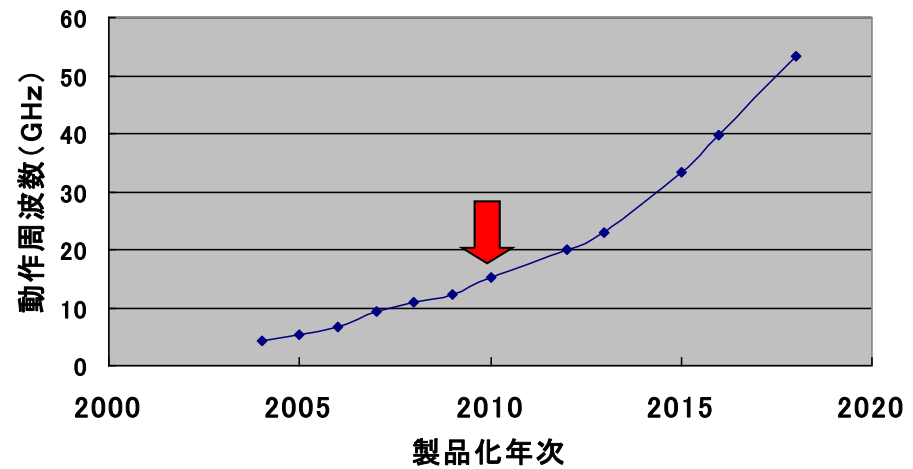


半導体技術のトレンド (ITRS2003資料から) FUJITSU

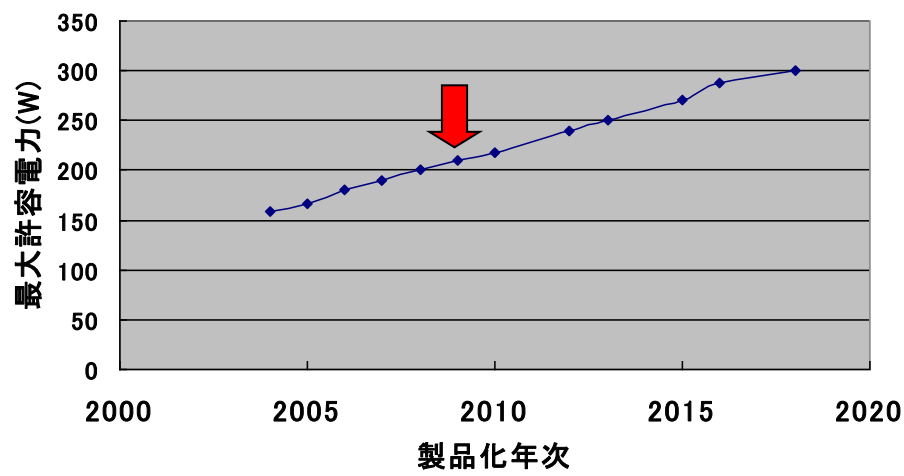
プロセス



On Chip周波数(GHz)



最大許容電力(W)

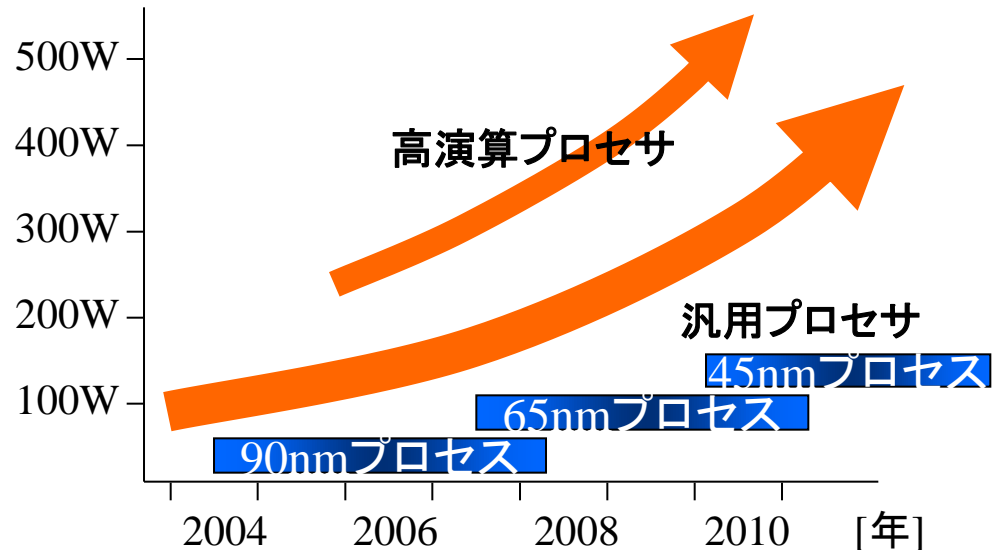
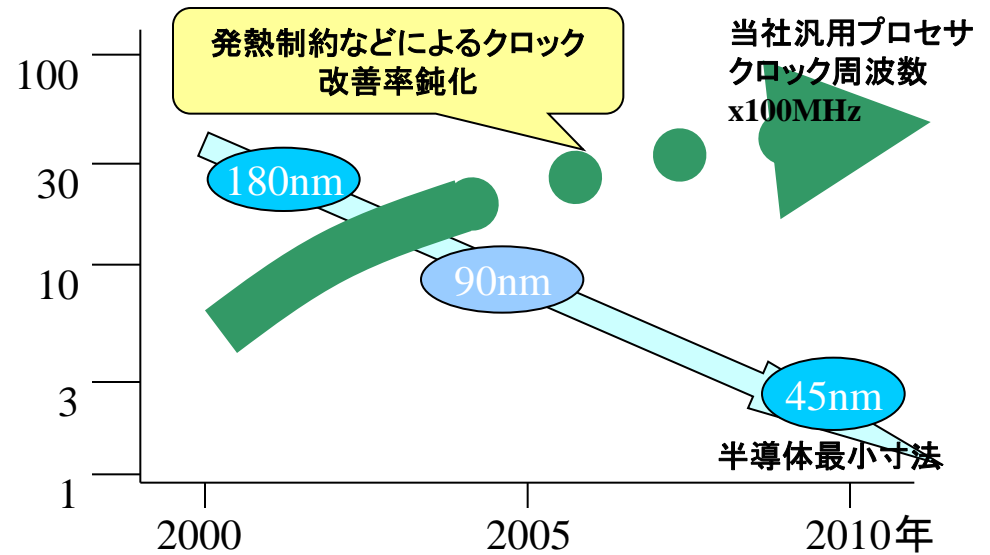


2010年頃
プロセス: 45nm
オンチップ周波数: 15GHz
最大許容電力: 210W

半導体技術と高速化のトレンド

- 微細化は進展するが、以下の問題が顕著化
 - 製造バラツキの増大
 - リーク電流の増大
 - 消費電力の増大

- 要素技術上のチャレンジ
 - 半導体製造技術
 - 低リークゲート絶縁膜
 - 移動度改善(歪シリコン等)
 - 回路設計
 - 製造ばらつき, 雑音に強い回路
 - 低電力回路



2. 主なアプリケーションソフトウェアについて

対象分野	アプリケーション	概要	目的	期待されるブレイクスルー	経済的波及効果
ナノテクノロジー	Car-Parrinello法第一原理分子動力学	<ul style="list-style-type: none"> MOSデバイス絶縁層・半導体・金属層界面の10,000原子超の系の、量子論に基づいた原子拡散、欠陥生成予測を行う。 CNTデバイス10,000原子超の系の、量子論に基づいたCNT成長、電極接合成長の予測を行う。 	<ul style="list-style-type: none"> シリコンデバイスの極限微細化設計 CNTデバイス構造の自己組織化制御 	超微細デバイス、CNTデバイスの実現	我が国の半導体技術、 しいてはIT技術が国際的に先行でき、高い競争力につながる。
バイオ・インフォマティクス	分子軌道法第一原理分子動力学 (or 半古典MD, 溶媒ポテンシャルモデルとの連成解析)	<ul style="list-style-type: none"> タンパク質や水、他分子の数千～数百万原子の系の量子論に基づいた構造解析や他分子との相互作用の予測を行う。 	<ul style="list-style-type: none"> タンパク質構造 動的機能予測 	創薬候補の計算機スクリーニング	我が国の創薬・医療技術の高度化と米国寡占からの脱却
気候予測	大気・海洋・海水結合モデル、他	<ul style="list-style-type: none"> 水平解像度1KM (現在～100KM) で大気、海洋、海水の振舞いを100年レンジでシミュレーション。 大気中化学反応、海中生物化学過程を含む。 	<ul style="list-style-type: none"> 大気変動、水循環、大気汚染の影響の予測 	長期気候変動の予測と社会経済動向の予測	エネルギー需要、環境コストの長期的予測につながる
CAE	衝突解析、流体解析、及びこれらの連成解析	<ul style="list-style-type: none"> 数千万メッシュレベル大規模構造解析、数百万メッシュレベルの衝突解析の大量計算 現行メッシュサイズの1000倍クラスの大規模メッシュによる精密又は大規模計算 高度な流体、構造連成計算 	<ul style="list-style-type: none"> 自動車、航空機、電子機器等の設計 	最適化設計による開発期間短縮、コスト削減	基幹の産業競争力確保

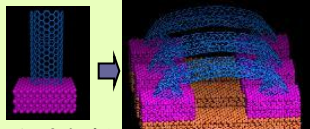
ペタスケールコンピューティングの拓く世界



大規模処理 計算規模の拡大・計算精度の向上

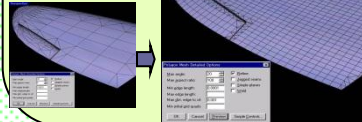
材料解析

1000倍以上の計算規模拡大による実際の設計に直結するシミュレーションへ



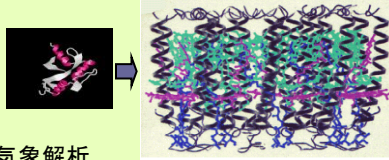
構造・衝突解析

1000倍以上の計算メッシュ拡大による解析精度向上。より正確かつ安全な設計へ



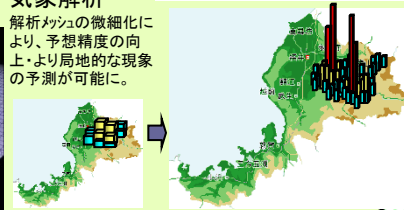
タンパク解析

厳密解法による精度向上と解析規模拡大により現実に近い複雑系の挙動解析を実現。



気象解析

解析メッシュの微細化により、予想精度の向上。より局地的な現象の予測が可能に。



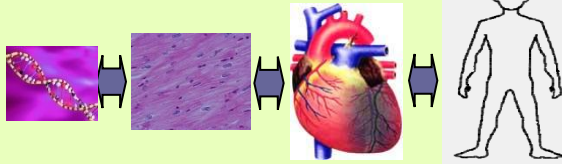
大規模処理 Scale up
複雑処理 Super-Scale
大量処理 Scale out

次元を越える
視点が変わる
見えないものが見える
価値観が変わる

複雑処理 マルチスケールコンピューティング・マルチフィジクスコンピューティング

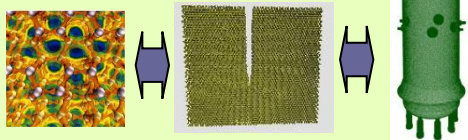
DNAから人体までの一貫した解析

DNAから発現するタンパク、細胞内のタンパクの挙動、それらが構成する器官、さらに器官から構成される人体までをモデル化し、薬剤の効果、遺伝疾患を遺伝レベルで対応可能



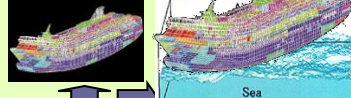
原子レベルからの破壊の一貫した解析

原子レベルの結合切断から目に見える破壊までを解析し、原理の理解、材料の改良へ。



構造流体連成解析

船の構造解析



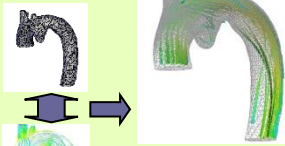
波の解析



関連する全ての要素を計算する動的解析によるより安全な、より効率的な製品作りへ

生体の動的解析

血管の構造解析



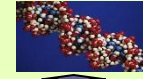
流れ解析

流れによる影響を考慮した血管の動的な解析により、より安全な診断・手術へ

大量処理 大量データの検索

遺伝情報検索

大量の遺伝子情報の中から目的遺伝子との類似情報を超高速に検索。未知ウイルスの探索、疾病診断などに活用。



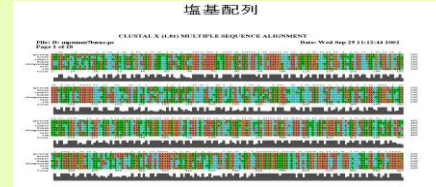
暗号解析

コンピュータウイルス検出

大量のインターネット情報からのコンピュータウイルス検出、暗号解析を高速に行い、インターネットのセキュリティを確保

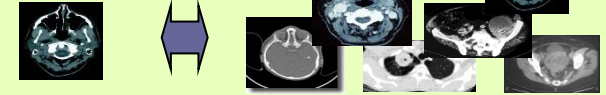


塩基配列



多次元類似度検索

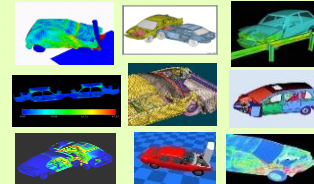
大量の4次元画像からの類似画像を高速に検索。医療診断等に活用



大量処理 処理回数の増大

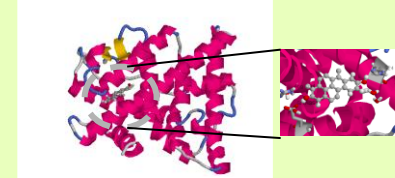
衝突解析

パラメータを変化させた大量解析による、計算モデルの検証、コード検証統計解析を行い、より安全かつ効率的な製品作りへ



タンパク解析

目的タンパクと多種類の薬候補化合物とのドッキング計算を高速に実行し、候補薬剤を絞り込み、疾患対策を迅速に。



安全・安心な社会の実現に向けて

創薬支援
診断・治療支援
未知疾病対策支援

医療・健康

災害シミュレーション
災害予知・予測
災害避難支援

防災

暗号処理
コンピュータウイルス検出
金融工学

セキュリティ

新材料設計・開発支援
製品設計・開発支援

物作りの革新

.....

3-1. 必要な要素技術について(ハードウェア) FUJITSU

- プロセッサアーキテクチャ

- 実効1PFlopsを実現する演算器, キャッシュ
 - 数百GFlops/チップの実現方式
- 演算エネルギー消費の低減
 - 現状の延長で性能を追求すると500W/chip → 200W以下
- 高信頼度の実現(演算エラー等の検出)
 - 従来の汎用マイクロプロセッサのチェックサム
 - 20,000CPUシステムでは, 500時間に1回エラー発生
 - 演算器二重化は電力, チップ面積の負担大 → エネルギー効率の高い検出手法の開発

SIMD型CPU vs Scalar型CPU

高信頼化技術

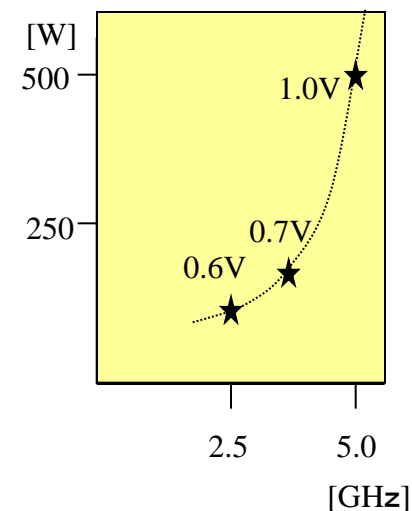
- システムアーキテクチャ

- 対象アプリケーション群に対して実効1PFlopsを実現できるシステムアーキテクチャ
 - 現状では数百~数千プロセッサが限界 → 数万~数十万プロセッサを効率的に動かす技術
 - 分散メモリ+インタコネク, 階層メモリ構造(バンド幅, レーテンシ, メモリ量など)の最適化
- 高可用性(エラーリカバリ)の実現方式
 - ホットスワップの実現

両方式の比較

	汎用Scalar型CPU	SIMD型CPU
ノード	2CPU+Memory	1CPU+L3\$+Memory
CPUチップ	5GHz 通常電圧 160GFlops ~200W	3GHz 低電圧駆動 768GFlops ~200W
メモリBW	128GB/s/CPU	400GB/s/CPU L3\$ 102GB/s/CPU
システム	10240ノード ~8MW	4096ノード ~1.6MW
アプリケーション	スカラ並列の高度化	スカラ並列の高度化 +SIMD型CPU向き アルゴリズム

低電圧駆動の効果
(SIMD型CPUの例)



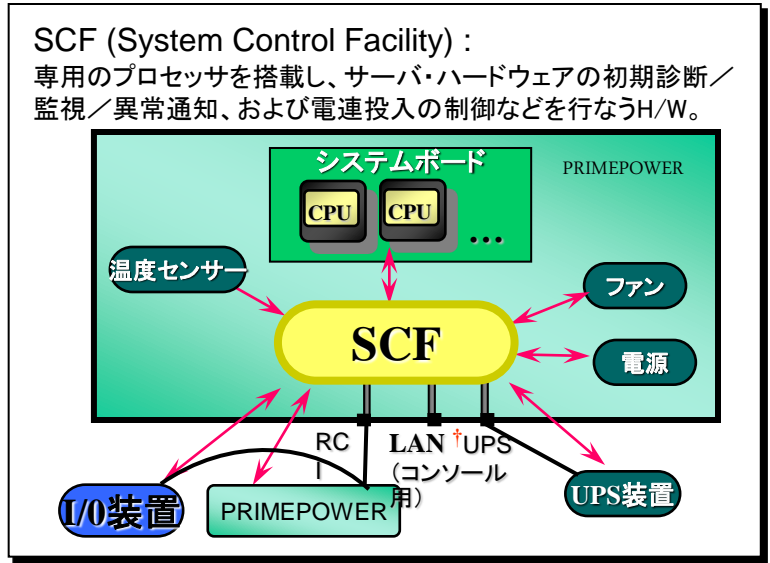
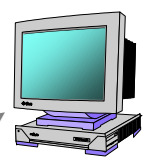
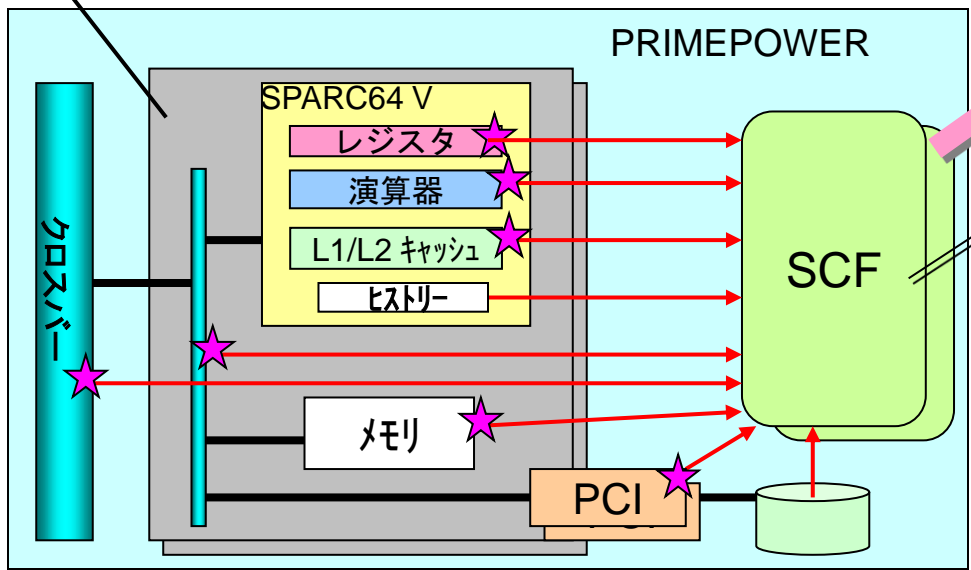
青字:チャレンジ項目

高信頼化技術：高精度な予兆監視・エラー原因特定

～ 万一のエラー発生時にも迅速に対処・復旧 ～

- エラー検出回路(チェッカー)を随所に内蔵
- 障害部品の特定、ヒストリデータ等による解析
- e-mailにて遠隔監視センターへ通知

▶ チェッカー
 • SPARC64™ V プロセッサ当たり **679** 個
 • システム当たり **90,000** 個
 (PRIMEPOWER2500/32CPUの場合)



RCI (Remote Cabinet Interface): ETERNUS等との電源連動、クラスタでのノード監視

UPS: 無停電電源装置との信号ケーブル接続

†: PRIMEPOWER250/450 に標準装備するXSCF (eXtended System Control Facility)でサポート

高信頼化技術：連続運転へのこだわり

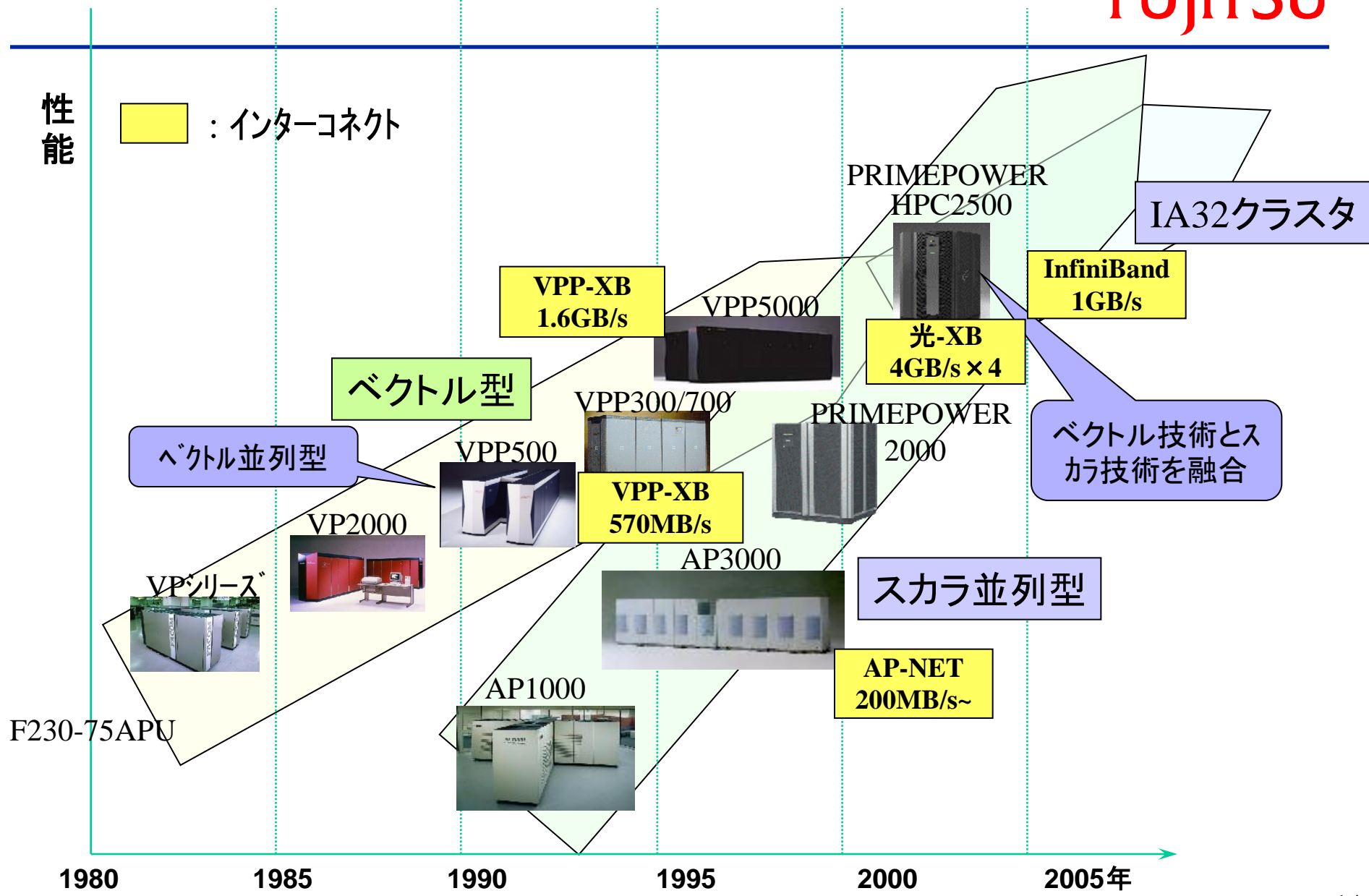
当社の 連続運転技術例

- 自己診断による予兆監視機能
- 局所的電源切断機能
- 故障箇所の自動切離し機能

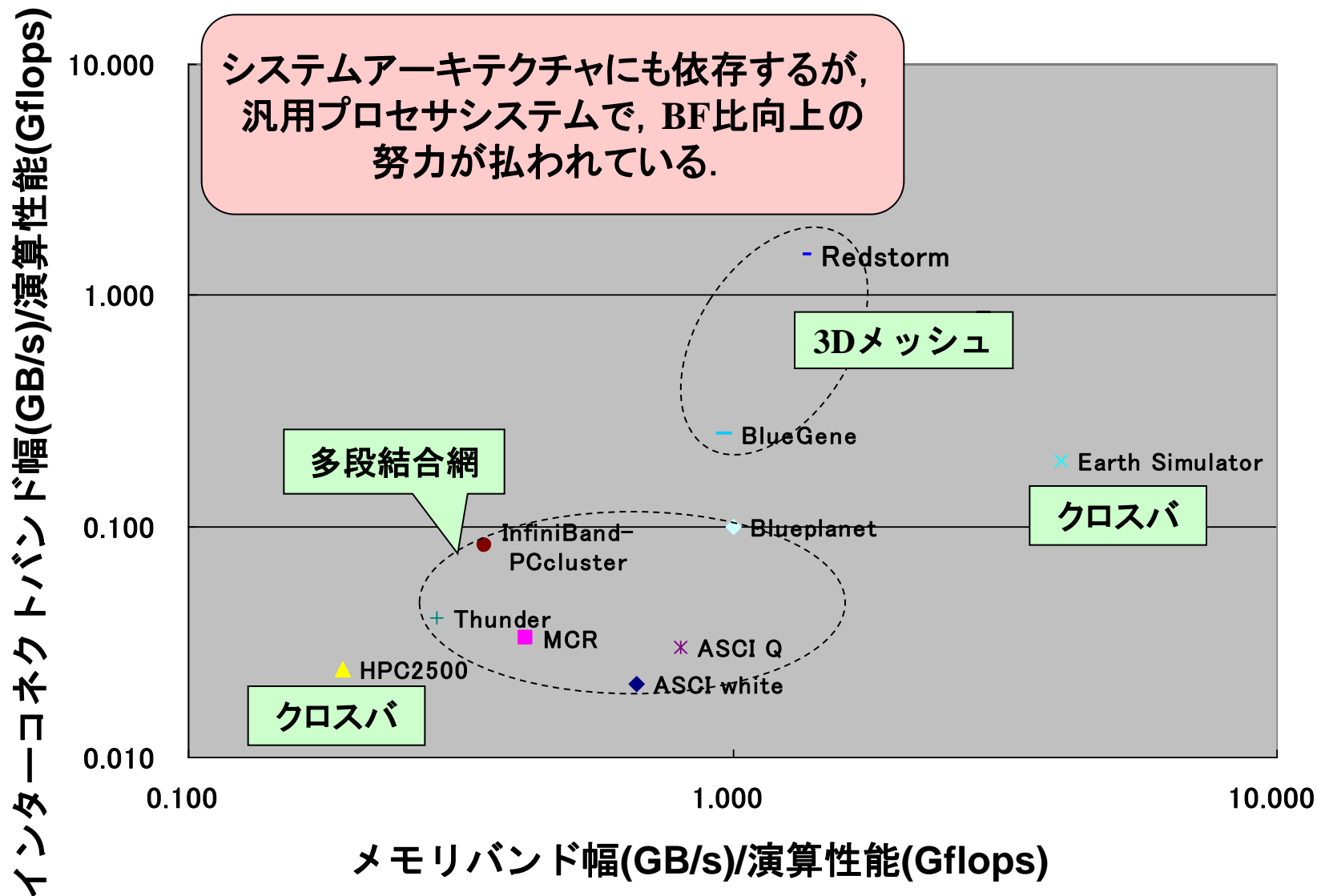
	A社マシン		B社マシン(推定)		
	動的縮退 冗長化	活性交換	動的縮退 冗長化	活性交換	
プロセッサ メモリ	[Solid Light Blue Box]	[Solid Light Blue Box]	[Solid Light Blue Box]	[Purple Box: 計画停止]	
バックプレーンクロスバ		[Purple Box: 計画停止]	[Dotted Box]	[Dotted Box]	
システムボード		[Solid Light Blue Box]	[Solid Light Blue Box]	[Yellow Box: 未対応]	[Orange Box: システム停止]
サービスプロセッサ			[Solid Light Blue Box]	[Solid Light Blue Box]	[Solid Light Blue Box]
電源・ファン			[Solid Light Blue Box]	[Solid Light Blue Box]	[Solid Light Blue Box]
I/Oカード			[Solid Light Blue Box]	[Solid Light Blue Box]	[Solid Light Blue Box]
ディスク			[Solid Light Blue Box]	[Solid Light Blue Box]	[Solid Light Blue Box]

- 低電力論理回路
 - 200W/チップで実現可能なロジックの開発と検討
 - 低電圧動作と動作マージンを確保する回路設計技術
- 高速伝送回路
 - プロセサチップとメモリ, インタコネクタ間の高速信号伝送
 - 現状3~6Gbps, 20~100mW/Gbit/s → 6~12Gbit/s, 10mW/Gbit/s以下の消費電力
- 高速ノード間インタコネクタ
 - 最適トポロジー, 必要性能の検討
 - アプリケーションが要求するデータ転送パターンを低オーバーヘッドで実現するホストインタフェース
 - スイッチチップ(あるいはプロセサ内蔵のスイッチ機能)の検討
 - 要求性能を実現するアーキテクチャ, ハードウェア規模, 消費電力等の検討
 - 現状, 10Gbps x 12port → 20Gbps x 16port のインテリジェントスイッチ
 - インタコネクタを実現する素材の検討

富士通のHPCプラットフォームとインターコネクト FUJITSU



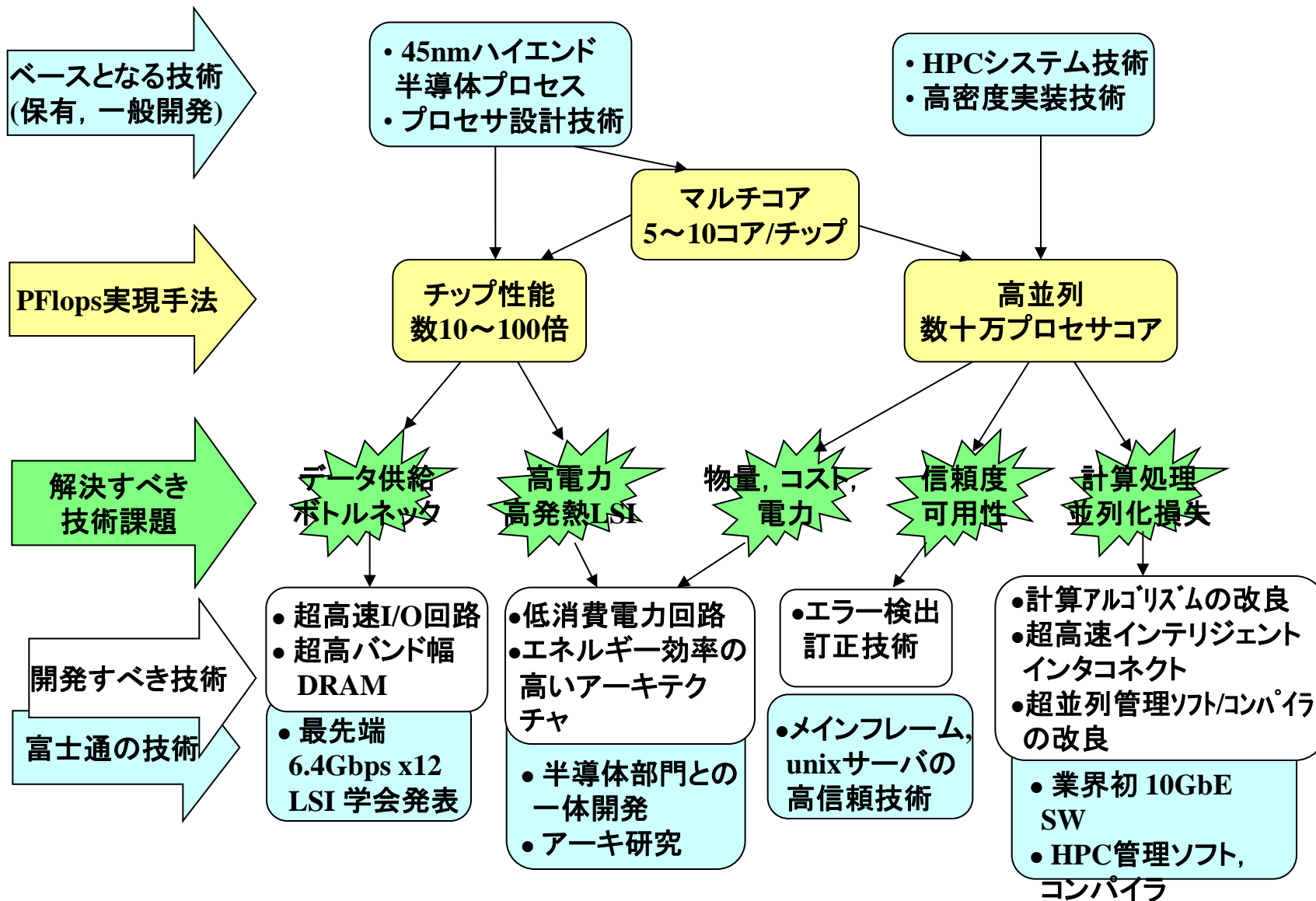
高速インタコネク性能バランス



3-2 必要な要素技術について(基本ソフトウェア)

- OS・ミドルウェア
 - ～数万ノードの効率的運用管理機能の検討
資源管理, ジョブスケジューリング, 操作・監視, 稼動情報など
 - 大容量高速分散並列ファイルの検討
ローカルファイル, ネットワークファイルなど
 - 超並列向けカーネル拡張の検討
- コンパイラ
 - ベクトル並みの実効性能を実現する自動並列コンパイラとハードアシスト機能の検討
 - ～数万ノードの並列最適化と開発環境の検討
 - 超並列ライブラリ(MPI、数学ライブラリなど)の検討
- 効率的な超並列実行の実現
 - 15万～50万演算器に処理を分割する方法は自明ではない.
 - バッファ, キャッシュ容量, メモリバンド幅, インタコネクティブバンド幅等を考慮して, 処理の並列化を支援するツールが必要
 - 実システムでのチューニングを可能にする豊富な動作状態情報の提供

必要な要素技術のまとめ



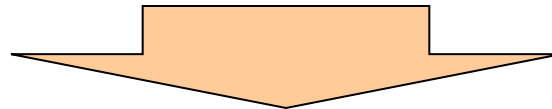
3-3. 必要な要素技術について(アプリケーション開発) — 利用技術の研究開発と普及 —

超PFlops・シミュレーションの実現には、ハードウェア・ミドルウェアの研究開発のみならず、既存手法の延長線上にない新しいシミュレーション手法・アルゴリズム等の研究開発が重要。

- 超大規模シミュレーションを実現するアルゴリズム, 最適化設計問題・連成解析などの先端シミュレーション手法の研究開発
- 超高並列処理技術・アルゴリズムの研究開発

シミュレーション技術の発展には、理論研究者, 利用技術開発者等が連携して、統合的にシミュレーション技術に関する研究開発を目的として、これを実施する環境・仕組みの整備が必要。

(米国では、ソフト開発者, 利用者, 理論研究者が連携してシミュレーションの研究開発を実施)

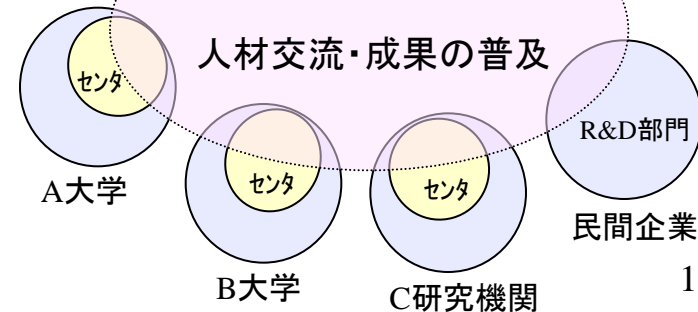


シミュレーション技術センターの設立

- 産学官連携体制による計算科学に関するCOEを設立し、人材育成, アプリケーションに関する研究開発を実施。成果を全国に普及。
- 必要に応じて、各機関に人材を派遣。

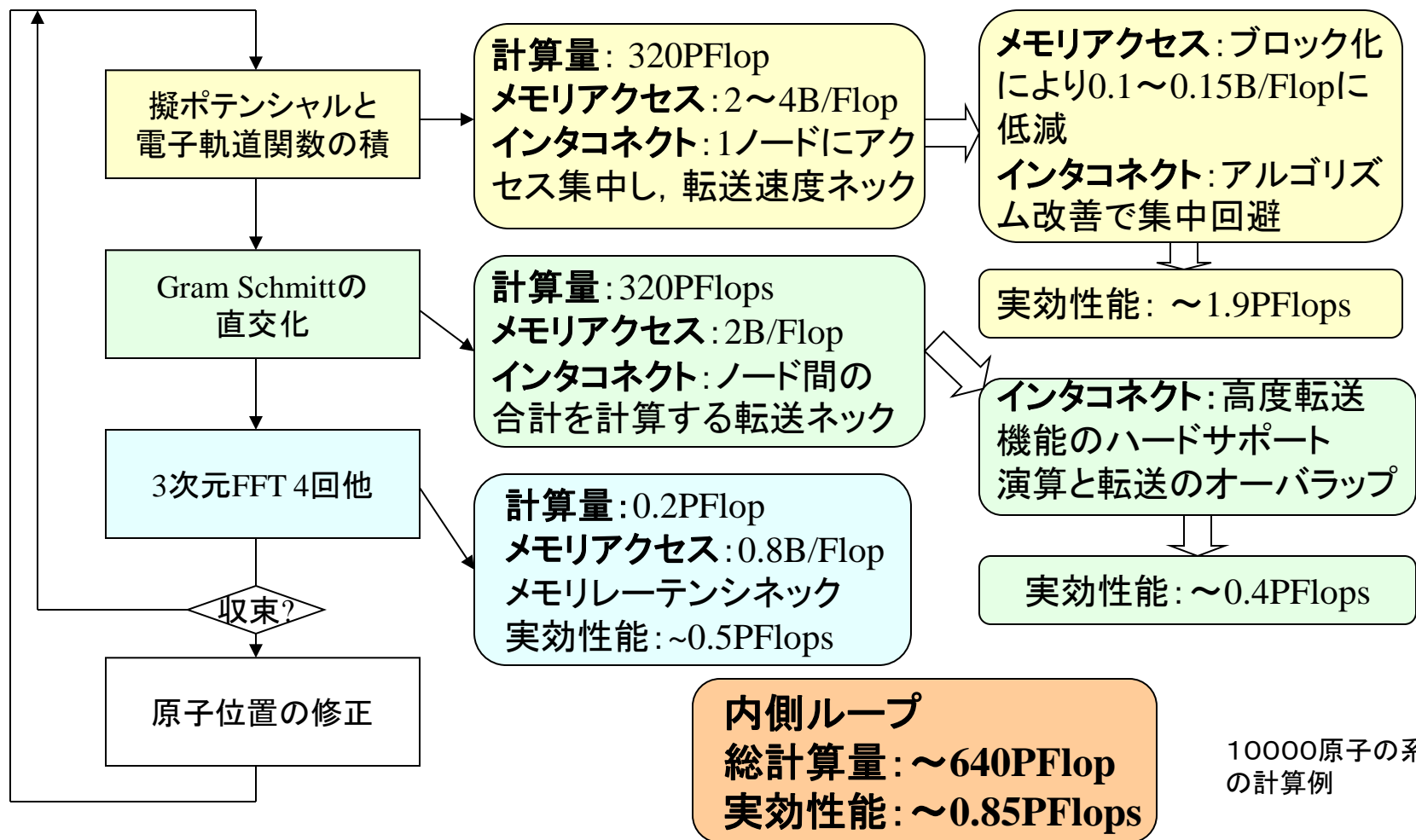
シミュレーション技術センター

- ・理論研究者を交えたシミュレーション技術の研究開発
- ・利用技術の研究開発



SIMD高演算システムによるCar-Parrinello法の処理分析例

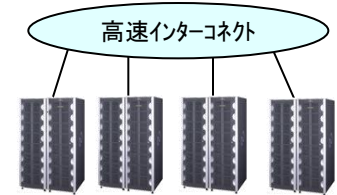
シリコン等の固体の原子レベルの第一原理シミュレーション法



実現する技術の波及分野

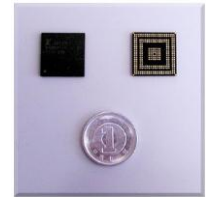
- クラスタ

- PFlopsシステムより小規模であるが、PFlopsシステムで開発したシステムアーキテクチャ、インタコネク、並列実行技術が有効に利用できる。



- メディア処理LSI等

- 画像、ビデオ処理、音声認識等の分野では大量の演算を低電力で実行する必要があり、演算エネルギーの低減技術は汎用的に利用可能である。また、メモリとの間の高速伝送技術も適用可能である。



- LAN, 機器間接続(次世代インタコネク等)

- 現状は2.5~3Gbpsの伝送であるが、低電力の6~12Gbps伝送技術は、これらの性能向上を可能にする。



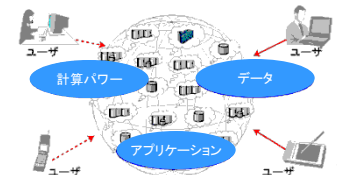
- 組込機器等

- 半導体の低消費電力技術は組込機器への応用が可能である。



- グリッドコンピューティング

- 超高速計算機の運用管理システムは、グリッドコンピューティングへの展開が可能。



ターゲットシステムの実現イメージ

既存のスーパーコンピュータシステム



既存スーパーコンピュータの
一つのロッカー並の性能

既存スーパーコンピュータ
システム全体並の性能

Peta Flopsシステム



プロセッサチップ
数百GFlops

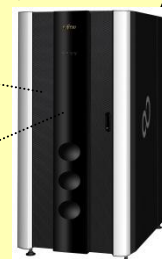
単体プロセッサ
8個程度を
1チップに集積



ボード

1~4プロセッサ
0.5~1TFlops

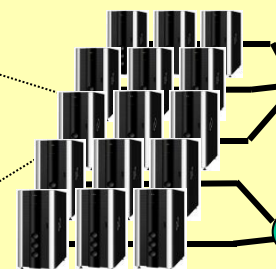
DRAM



ロッカー

数十ボード

10~数10TFlops



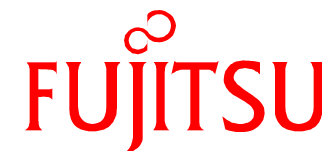
システム

100~500ロッカー

数PFlops

超高速
インテリジェント
インターコネクト

4-1. 要素技術(ハードウェア)の研究開発について



○研究開発実施体制

大学、研究機関

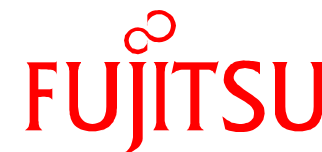
○スケジュールおよび費用

要素技術の実現: 20億円

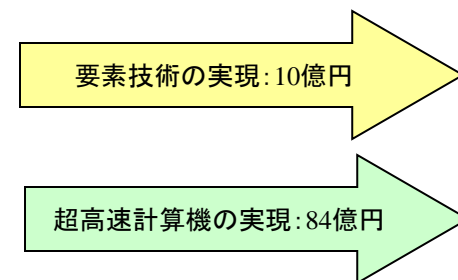
超高速計算機の実現: 460億円

要素技術	初年度	第2年度	第3年度	第4年度	第5年度	第6年度
システムアーキテクチャ, およびプロセッサアーキテクチャ	方式検討 : 1億円	シミュレータ開発 : 2億円	シミュレーションによる実証 : 1億円 商用システム方式設計 : 15億円	商用システム設計 : 120億円		チップ, システム製造, 試験 : 250億円
低電力論理回路, および高速伝送回路	回路設計 : 2億円	チップ試作(65nm) : 6億円	評価 : 1億円			
高速ノード間インタコネク	方式検討 : 1億円	論理設計 : 1.5億円 FPGA実装 : 1.5億円	商用方式設計 : 5億円 実証 : 3億円	商用インタコネク設計 : 40億円		製造, 試験 : 80億円

4-2. 要素技術(ソフトウェア)の研究開発について



- 研究開発実施体制
大学、研究機関
- スケジュールおよび費用



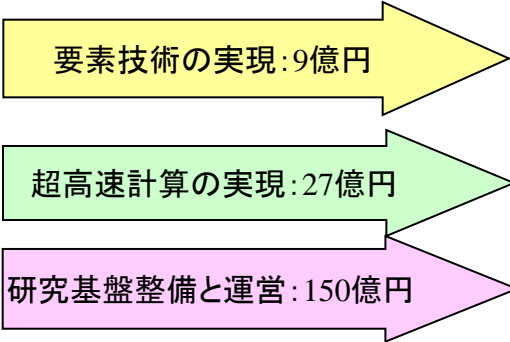
要素技術	初年度	第2年度	第3年度	第4年度	第5年度	第6年度
OS・ミドルウェア ・運用管理機能 ・高速分散並列 ファイル	方式検討 :1.5億円	プロトタイプ開発 :3億円	プロトタイプ実証 :3億円	商用方式設計 :8億円	商用システム開発:32億円	商用システム実証 :16億円
コンパイラ ・自動並列 ・最適化/ 開発環境 ・ライブラリ	方式検討 :0.5億円	プロトタイプ開発 :1億円	プロトタイプ実証 :1億円	商用方式設計 :4億円	商用システム開発:16億円	商用システム実証 :8億円

4-3. 要素技術(アプリケーション)の研究開発について

○研究開発実施体制

大学、研究所、民間企業が積極的に連携

○スケジュールおよび費用



要素技術	初年度	第2年度	第3年度	第4年度	第5年度	第6年度
超大規模・超高並列シミュレーション技術開発	アルゴリズム・手法研究開発 :1億円/年					
		超並列アプリプロトタイプ開発 :3億円/年		実用シミュレーション・ソフトの開発 :9億円/年		
研究基盤整備技術の普及	基盤整備 10億円		基盤整備 40億円			
		産官学連携シミュレーション技術センターの運営 :20億円/年				

まとめ

- プロセッサ
 - コスト性能比の良いアーキテクチャの採用
 - SIMD型 vs 汎用Scalar型
- インタコネクト
 - 光 vs 電気
 - 性能バランス(メモリスループット:BF比)の改善
- 半導体
 - 微細化に伴う課題の解決:製造ばらつき, リーク電流, 消費電力増大
- ソフトウェア
 - 開発環境の整備:コンパイラ, 性能モニタ, etc
- 運用管理
 - 高信頼の追及
 - 使いやすさの追求
- アプリケーション
 - 適用分野の拡大: Compute-intensive から探索, マイニング等へ