

# レイテンシコアの高度化・高効率化による 将来のHPCIシステムに関する調査研究

## Feasibility study on advanced and efficient latency core-based architecture for future HPCI R&D

---

石川 裕 東京大学情報基盤センター  
平木 敬 東京大学情報理工学系研究科  
青柳 睦 九州大学情報基盤研究開発センター  
新庄 直樹 富士通  
飯田 恒雄 日立製作所  
中村 祐一 日本電気

レイテンシコアの高度化・高効率化……

2012//08/10

# FS (Feasibility Study)概要

2018年頃設置可能並列システムを汎用型プロセッサからのアプローチでFS  
アプリケーション、システムソフトウェア、アーキテクチャのco-design  
システムソフトウェアスタック共通化 (From PC cluster to high-end machines)

ターゲットアプリケーション群

電力: 20MW~30MW  
設置面積: 2000m<sup>2</sup>

ポータビリティ

スケーラビリティ

メモリ階層  
と性能

低消費電力機  
構

耐故障性

並列マシンアーキテクチャ  
CPUアーキテクチャ、ノードアーキテ  
クチャ、インターコネクトアーキテ  
クチャ

階層ストレージアーキテクチャ  
記憶デバイス、ネットワーク、ファイ  
ルシステム

システムソフトウェア  
オペレーティングシステム、通信ライ  
ブラリ、ファイルI/Oライブラリ、数値  
計算ライブラリ

- プロジェクトの総合的推進及びシステムソフトウェア及びアプリケーション性能予測に関する検討 (東京大学情報基盤センター・大気海洋研・物性研・工学系研究科、協力機関: 理研AICS)
- アーキテクチャ評価およびコンパイラ技術と省電力機構の検討 (東京大学情報理工学系研究科)
- インターコネクト性能推定環境の検討 (九州大学情報基盤研究開発センター、協力機関: 九州先端科学研究所)
- アーキテクチャ概念設計およびシステムソフトウェアに関する検討 (富士通)
- 階層ストレージとコモディティ向けシステムソフトウェアに関する検討 (日立製作所)
- 低遅延通信機構の検討 (日本電気、協力機関: NECドイツ)

# 体制

## アーキテクチャ

東大情報理工  
業務主任者: 平木  
アーキ評価&コンパ  
イラ: 平木、須田  
省電力: 中村

九州大学  
業務主任者: 青柳  
インターコネク  
ト: 井上、  
稲富  
協力機関: 九州先端研  
柴村、眞木

東大情報基盤センター  
業務主任者: 石川  
アプリケーション: 藤堂、岩田、  
内田、羽角  
システムソフト: 石川、實本  
性能予測: 片桐、中島、大島  
協力機関: 理化学研究所  
堀、西澤、他

富士通  
業務主任者: 新庄  
概念設計: 清水  
安里、森田、青木、安島  
システムソフト: 住元、  
小田和

## アプリケーション

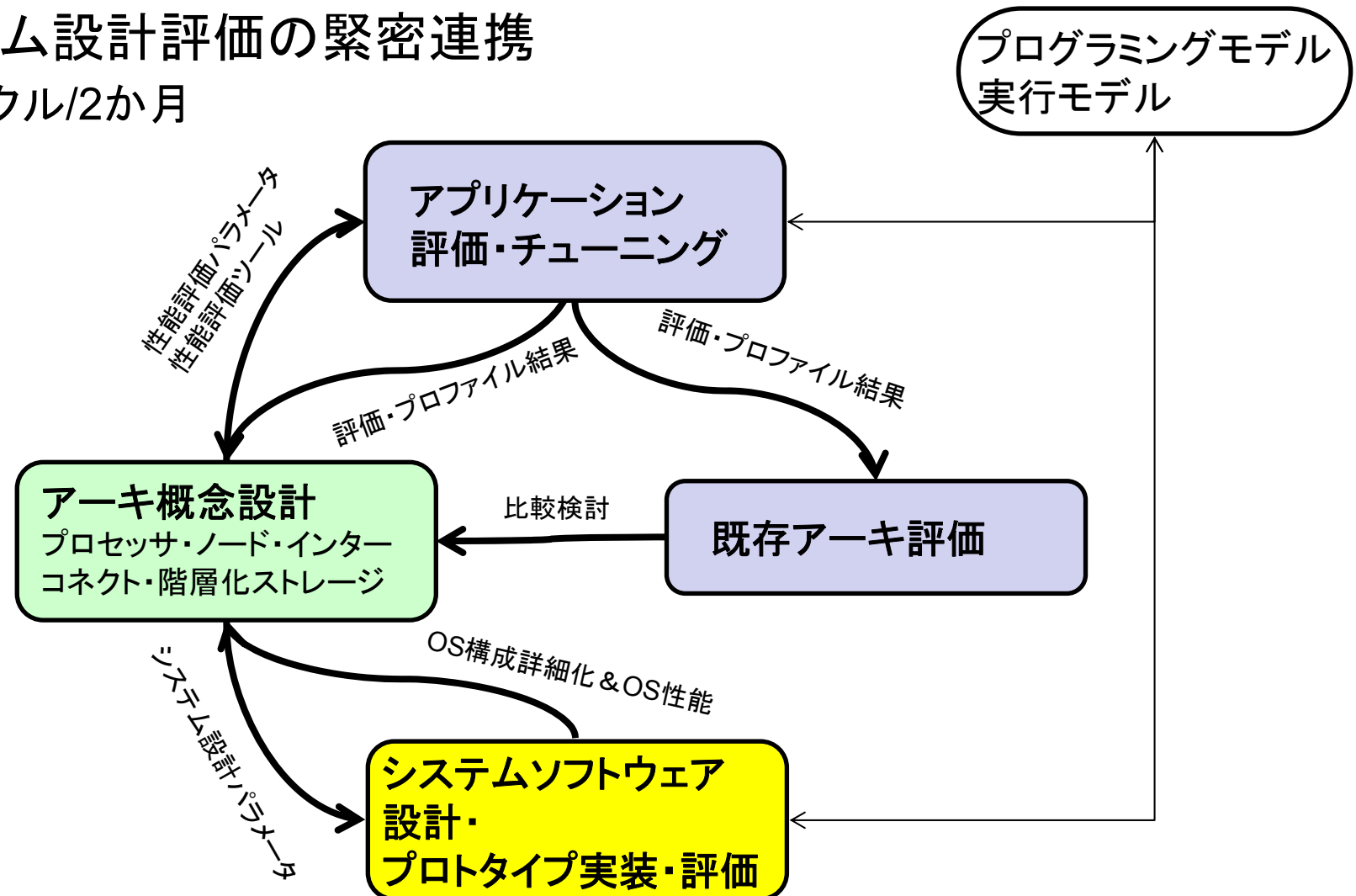
日立  
業務主任者: 飯田  
階層ストレージ: 飯田、清水  
アプリケーション評価: 米村  
佐伯

NEC  
業務主任者: 中村  
低遅延通信: 中村、高木  
協力機関: NECドイツ Ritzdorf

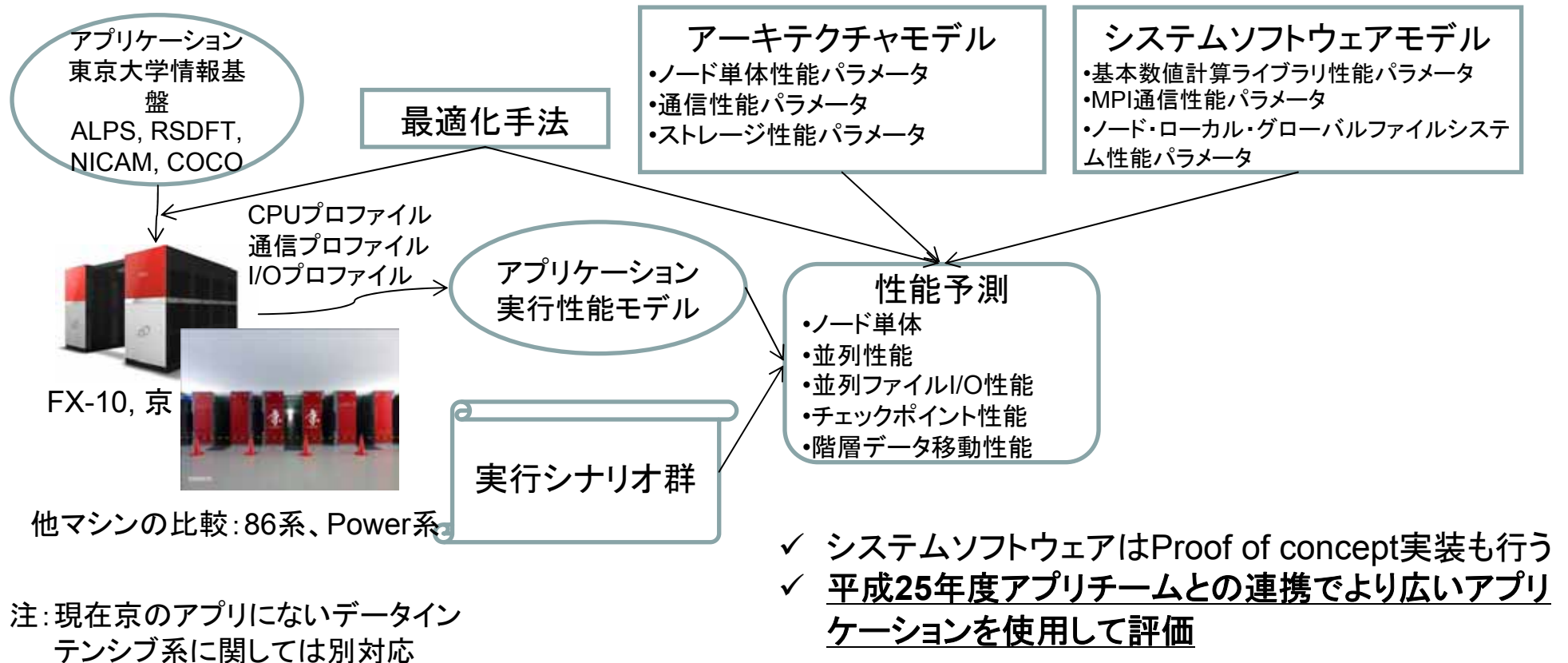
## システムソフト

# Co-designの進め方

- アーキテクチャとアプリ評価チューニングおよびシステム設計評価の緊密連携
  - 1サイクル/2か月



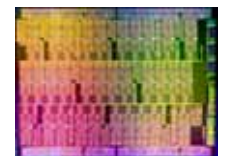
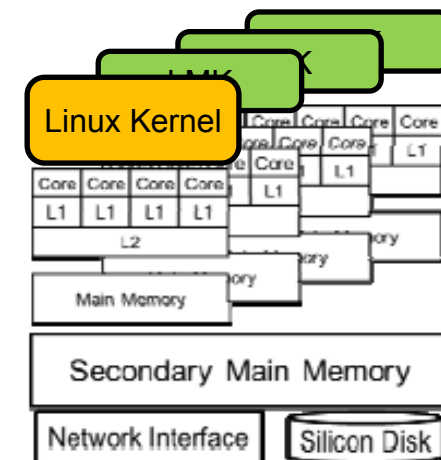
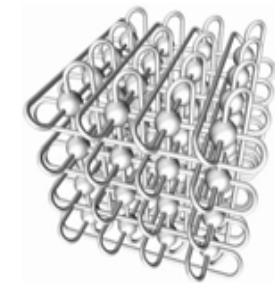
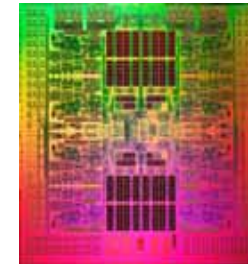
# FS概要:進め方



「今後のHPCI技術開発に関する報告書」を尊重し、京およびFX10におけるアプリケーション並列性能およびI/O性能、耐故障性および運用・保守の観点で課題を精査、概念設計に反映

# FS概要：特徴と公開成果物

- 汎用型プロセッサからのアプローチ
  - 京コンピュータ、FX10CPUを基に、メニコア化、SIMD化、通信機構の高度化を検討
  - ハードウェア概念設計時からシステムソフトウェアスタックの設計開発の一部を先行し、コモディティベースクラスタでも利用できるよう設計。国際協力。システムソフトウェアスタックのガラパゴス化を回避
    - ヘテロOSカーネル
    - 低遅延通信ライブラリとMPI
    - スケーラブルファイルI/Oと階層化ストレージ
- 運用を考えたシステム設計
  - ユーザアプリケーションの利用シナリオ
    - 課題解決に必要とされる資源量と必要時間
      - 大規模、中規模、小規模実行
  - 京コンピュータ、FX10@東大における運用経験
    - スケジューリング、耐故障性
- 開発計画、必要経費(開発経費、設置経費)、保守・運用経費の見積もり



# H24年度全体スケジュール



# アーキテクチャ概念設計の調査研究の進め方

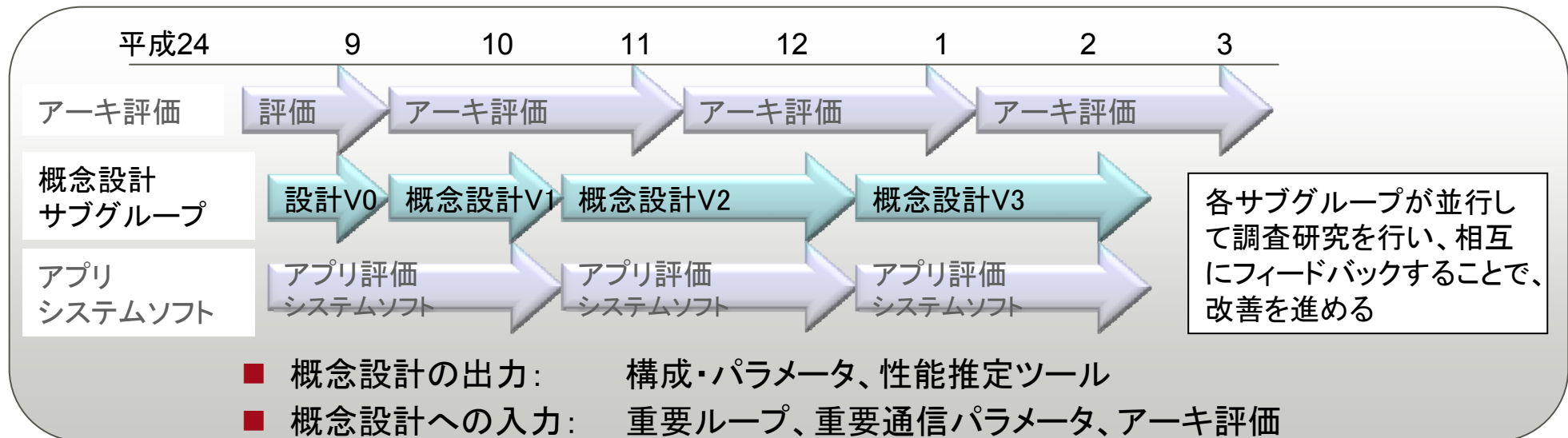
## 「今後のHPCI技術開発に関する報告書」(白書)の汎用型スペックとアプリ要求スペック

- 白書の汎用(従来)型のスペック
  - 設置面積 2,000m<sup>2</sup>
  - 消費電力 20~30MW
  - 総演算性能 200~400 PFlops
  - 総メモリ帯域 20~40 PB/s
  - 総メモリ容量 20~40 PB
  - インタコネク 16 GB/s × 8/node (4D-Torus)  
or 32 GB/s × 64 port/32 nodes (Dragonfly)

- 白書のアプリケーション要求スペック
  - 総演算性能 800~2500PFlops

## 調査研究の進め方

- 技術動向の予測に基づき、複数のアーキテクチャ設計方針を策定
- 白書のアプリケーション要求性能に近づけるために、新技術開発を概念設計に加味

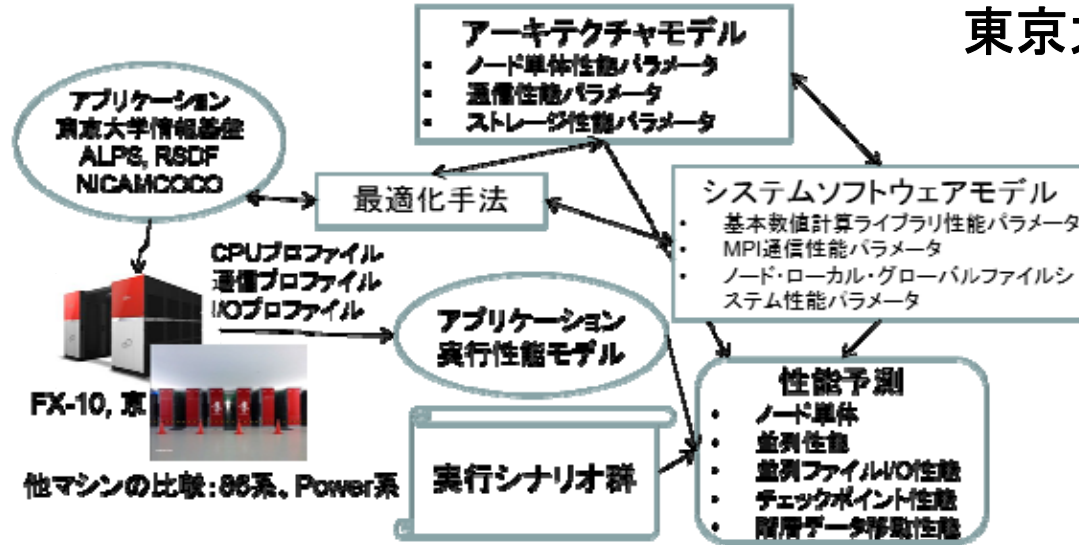






# アーキテクチャ評価およびコンパイラ技術と省電力機構の進め方

## 東京大学情報理工学系研究科



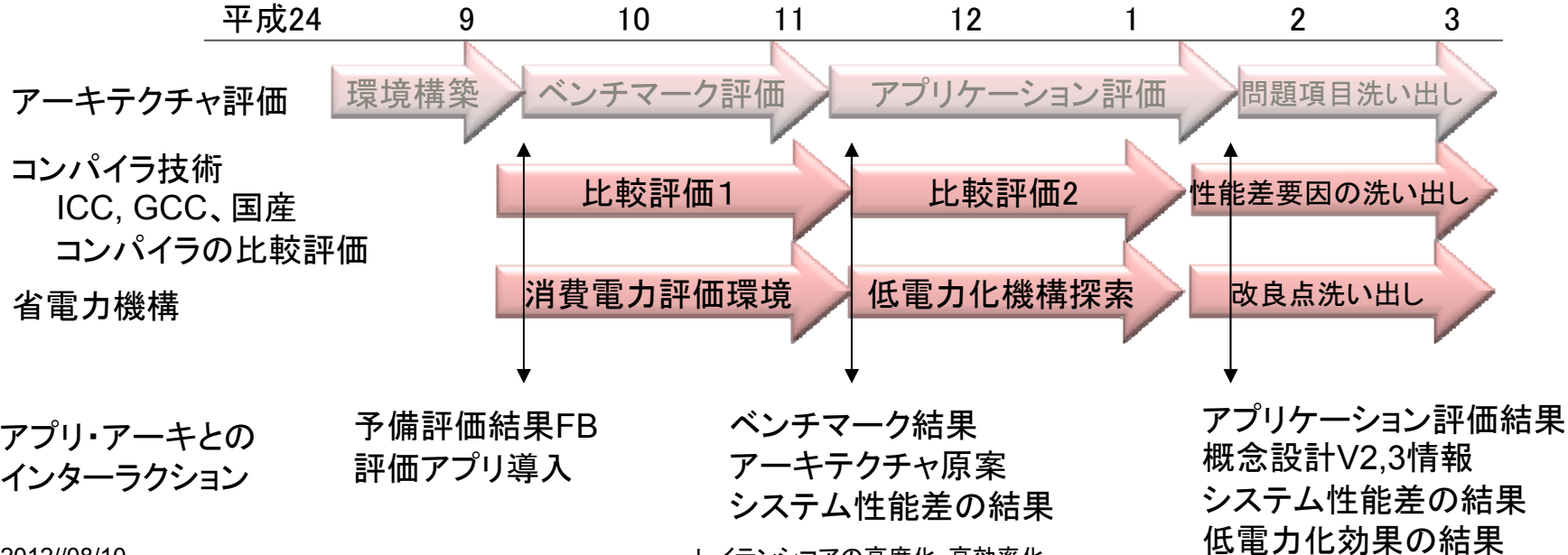
### アーキテクチャ評価およびコンパイラ技術

- 評価環境の構築: FX10, IvyBridge 環境整備、ソフトウェアの開発
- アプリケーションによる評価: 標準的ベンチマークソフトと第一原理計算等アプリケーションによる評価

### 省電力機構

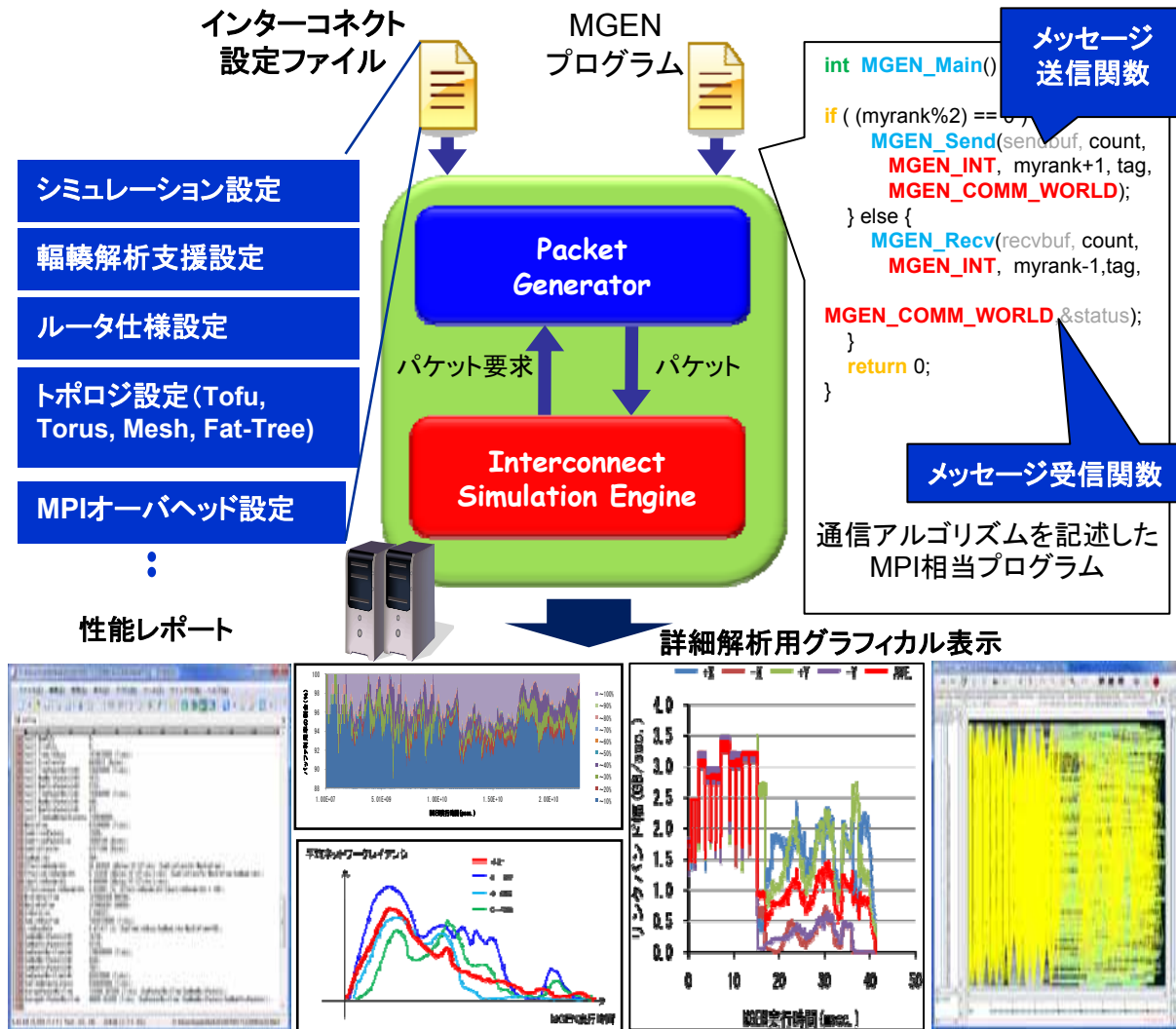
- インターコネクタ部の消費電力分析と省電力機構の検討

### スケジュールとインターラクション

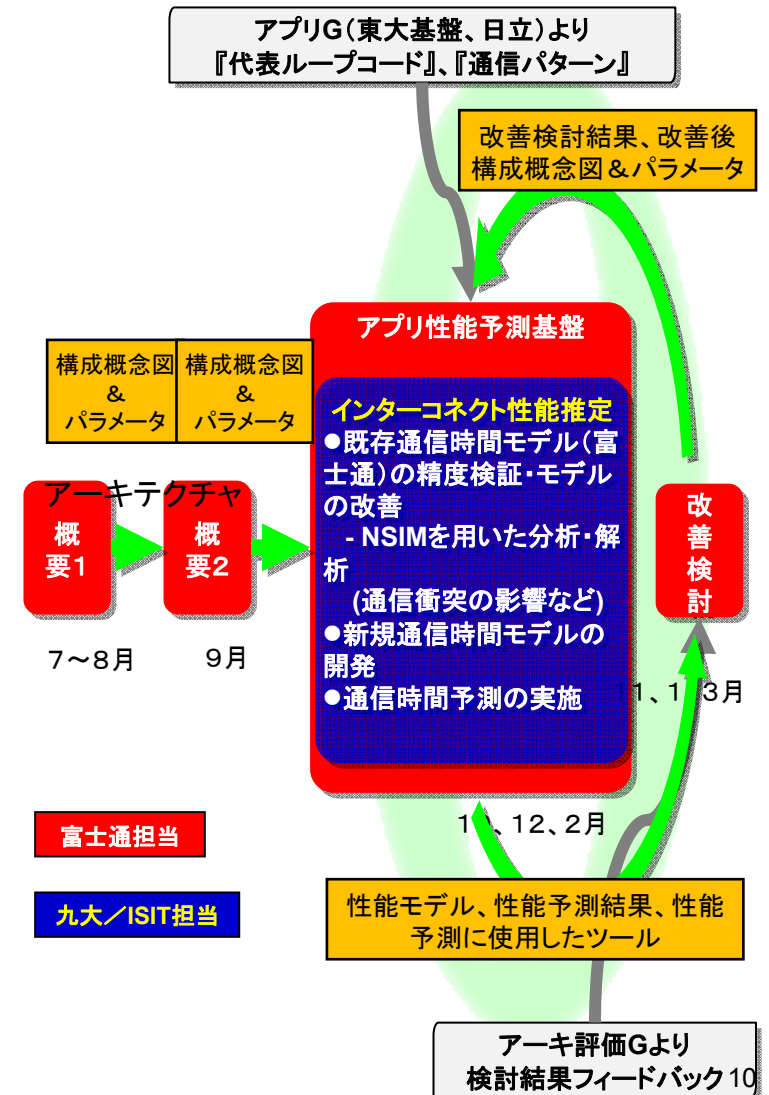


# インターコネクト性能評価に関する検討の進め方

## ■ 性能予測手法



## ■ スケジュール(進め方)



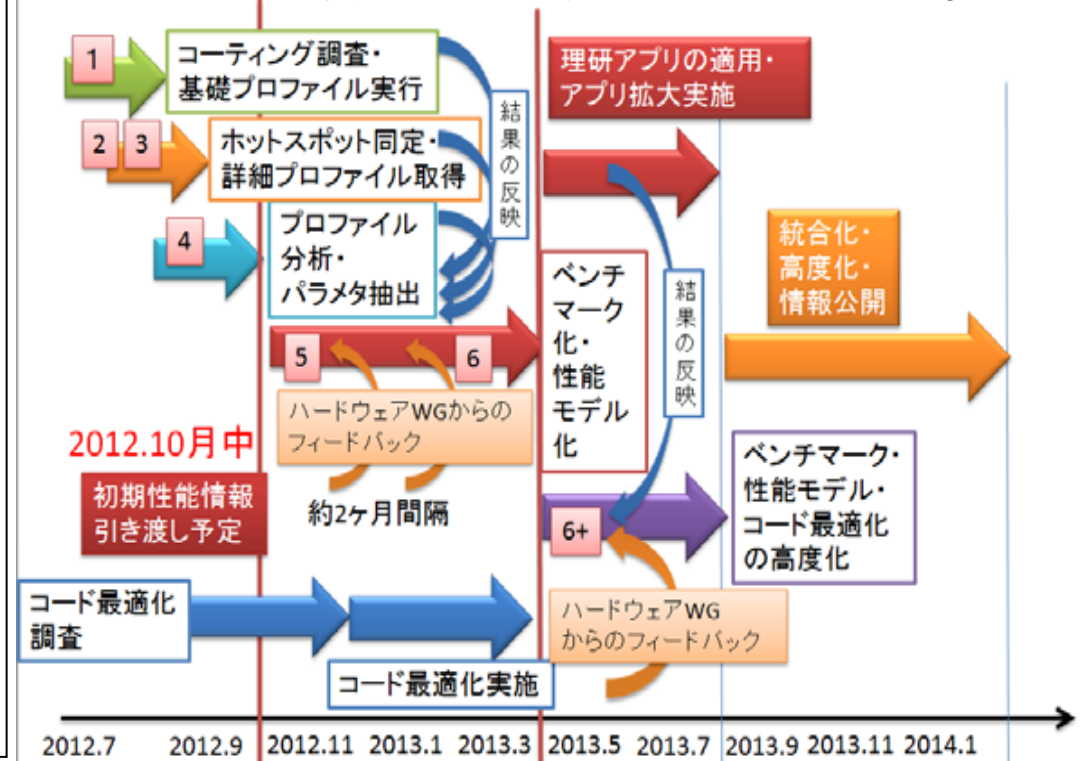
# アプリケーション性能予測に関する検討の進め方

## ■ 性能予測手法

1. **ホットスポット同定**: 基本プロファイラで複数のホットスポット(ループレベル)を同定、全体性能の近似
2. **カーネル分離**: (目視により) 計算部分、通信部分、I/O部分の分離
  - 計算部分: 演算カーネル
  - 通信部分: 通信カーネル
  - I/O部分: I/Oカーネル
3. **通信パターン確認**
4. **詳細プロファイルと分析**: 詳細プロファイラを用い、ホットスポットごとにハードウェア性能情報を取得し分析
  - 演算カーネルの演算効率/命令発行量/キャッシュ利用効率
  - 通信カーネルの通信回数/量/通信待ち時間
  - I/Oカーネルのデータ読み書き量/頻度
5. **ベンチマーク化**: ホットスポットのみで動作するようにコードを再構成
  - マシン特化の書き方、および、汎用的な書き方の2種を区別
  - 演算カーネル、通信カーネル、I/Oカーネル分類
6. **詳細モデル化**: ハードウェア因子を用いた数式による実行時間の近似

## ■ スケジュール予定

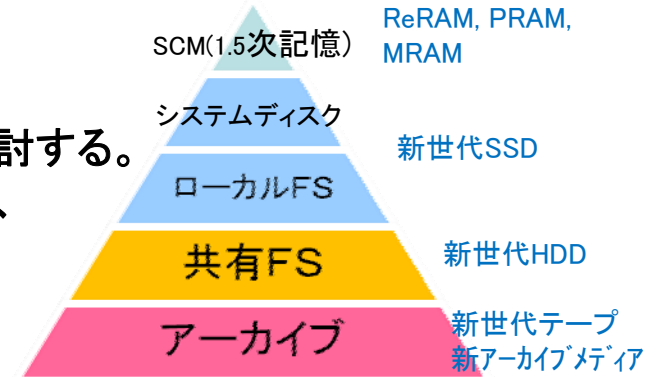
- 平成24年度: 4種アプリの性能プロファイリングを行い、ハードウェア設計に反映
  - エクサ環境で想定される実行モデルをアプリケーションごとに作成
  - 2か月間隔でハードウェア設計をフィードバック
  - コード最適化も同時に行い適切な性能プロファイリングを実現
- 平成25年度: 理研と連携しアプリケーション拡大



# 階層ストレージに関する検討の進め方

## 【目的】

- ・計算ノードローカルストレージ(システムディスク)、  
計算ノードローカルストレージ(ユーザ、ステー징、チェックポイント)、  
グローバルストレージ(共有ファイルシステム)、  
アーカイブをトランスペアレントに扱えるファイルI/O機構(階層間API)を検討する。
- ・2018年に利用可能な、新メモリ、新アーカイブメディアの技術動向も調査し、  
ストレージ階層への取り込みを検討する。



## 【調査研究の進め方】

- ・「京」、FX10の現状アプリからストレージ要求パラメータを取得、将来性能導出
- ・ゲノム系アプリからデータインテンシブ系ベンチマークを策定(ランダムI/O系?)
- ・トレンド・動向から利用可能なストレージパラメータを抽出
- ・ローカル・グローバル・アーカイブ階層の性能設計、機能設計

平成24 9 10 11 12 1 2 3



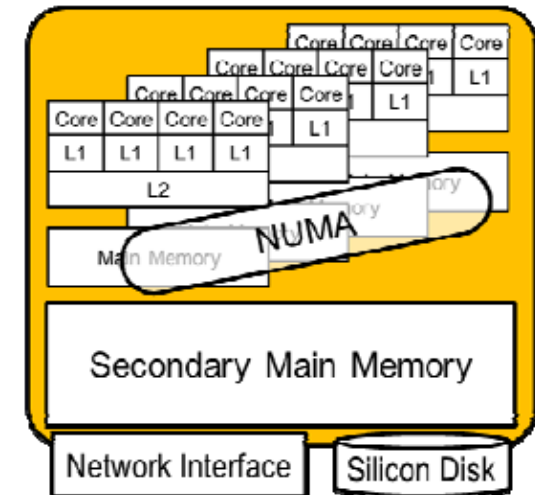
# システムソフトウェアに関する検討の進め方

- OS設計に必要なハードウェア構成を決めて、proof of concept実装を行いアーキテクチャ概念設計に反映
- 下表の検討すべき性能阻害要因に基づくシステムソフトウェア構成法および新OS機構の設計・試作・評価
- システムソフトウェア研究開発の国際協力(国際標準を含む)実現に向けて関係機関と調整

## 性能阻害要因

| Core Capability   | O(100) Core / Node   | O(1M+) Node, O(100M+) Core   |
|---|--|--|
|   | <ul style="list-style-type: none"> <li>Architectural support for scalability</li> <li>Programming Model &amp; Execution Model</li> </ul>   |  |
| <ul style="list-style-type: none"> <li>Reducing cache pollution</li> <li>Localizing data</li> <li>Managing memory hierarchy</li> <li>Utilizing SIMD resource</li> </ul> | <ul style="list-style-type: none"> <li>Reducing memory contention</li> <li>Reducing data movement among cores</li> <li>Providing fast communication</li> <li>Parallelizing functions achieving less data movement</li> <li>Sharing hardware resources</li> </ul> | <ul style="list-style-type: none"> <li>Reducing information whose size grows by # of cores &amp; nodes, e.g., comm. connection.</li> <li>Supporting strong scale in terms of communication</li> <li>Handling huge number of files, access and store</li> <li>Providing fault resilience</li> </ul> |

## ノードアーキテクチャの一例



## OS構成の一例

