

# 地球シミュレータの運用状況

大倉 悟, 中川 剛史, 甲斐 恭, 板倉 憲一

独立行政法人海洋研究開発機構 地球情報基盤センター

情報システム部 基盤システムグループ

okura@jamstec.go.jp

概要：独立行政法人海洋研究開発機構(JAMSTEC)では、2002年3月より世界最大規模のベクトル並列型スーパーコンピュータ「地球シミュレータ」を運用している。2009年3月の地球シミュレータ(ES2)へのシステム更新を経て、2015年3月には、後継システムの運用を開始する予定である。本稿では、地球シミュレータ(ES2)の概要及び運用状況と、後継システムの概要について述べる。

## 1 はじめに

初代地球シミュレータは、故・三好甫プロジェクトリーダーの元、旧・宇宙開発事業団(現・宇宙航空研究開発機構)、旧・日本原子力研究所(現・日本原子力研究開発機構)、旧・海洋科学技術センター(現・海洋研究開発機構)の三法人によって開発された、世界最大規模のベクトル並列型スーパーコンピュータである。大気大循環モデルの実行性能で5TFLOPS以上という開発目標を大きく上回る26.58TFLOPSを達成するとともに、LINPACKベンチマークにおいても

35.86TFLOPSを達成した[1]。2009年3月にはシステムを更新し、現在はNEC製SX-9Eを中心としたシステム(ES2)が運用されている。

本稿では、ES2の概要と運用状況について紹介する。また、2015年3月に稼働が予定されている、後継システムの概要について紹介する。

## 2 ES2の概要

ES2の概要を、図1に示す[2]。160台の計算ノードは、Fat-Tree構成のノード間ネットワークにより接続される。計算ノードは、インタラクティブノード(1台)、プリ・ポスト処理やデバッグ処理用のS系バッチノード(3台)、大規模並列処理用のL系バッチノード(156台)の、三種類の用途に使い分けられる。各計算ノードには、バッチジョブ実行中のテンポラリ領域としてのワークディスクが接続されており、パーマネント領域であるユーザディスクとの間で、バッチジョブの実行前後にファイルのステージング処理が行われる。

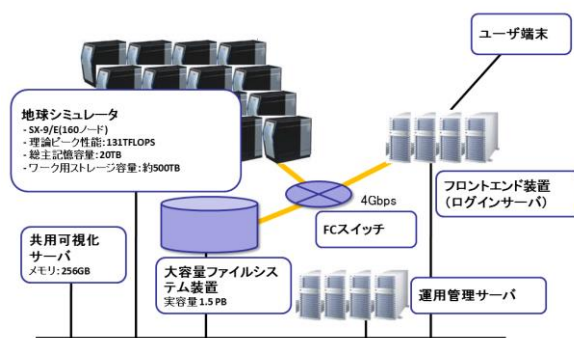


図1 ES2システムの概要

## 2.1 計算資源

ES2では、計算資源の割当てを、プログラムが使用する計算ノードの数と、それを実行した経過時間の積を用いて行う。例えば、計算ノード1台を1時間使用した場合、「1ノード時間」を使用したものとする。L系バッチノード数に、各年度において最大利用可能な時間数を乗じた値を、配分可能な総ノード時間とする。図2に2014年度の資源配分割合を示す。

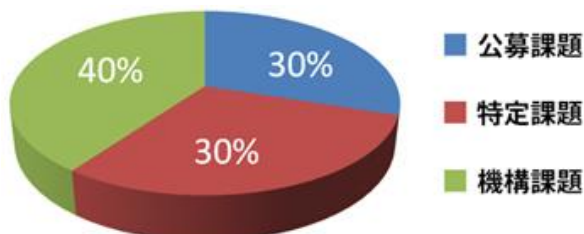


図2 ES2の資源配分割合

現在、公募による利用(公募課題)に30%、国等からの委託・補助等による利用(特定課題)に30%、

JAMSTEC が主導する研究課題による利用(機構課題)に 40%が充てられている。

## 2.2 バッチ処理環境と運用管理ソフトウェア

ES2 の運用では、高並列プログラムの効率的な実行と、システム全体の利用率の向上を、同時に実現することを目指してきた。初代地球シミュレータでは、ジョブスケジューラを含む運用管理ソフトウェアを独自開発したが、その基本的な考え方は、現在運用中の ES2 にも受け継がれている。

ES2 におけるプログラムの実行は、バッチ処理システムを介して行われる。バッチ処理システムとして NEC 製「NQSII」が、ジョブスケジューラとして同じく NEC 製「JobManipulator」が使用されている。図 3 に、ES2 におけるバッチ処理概念図を示す。

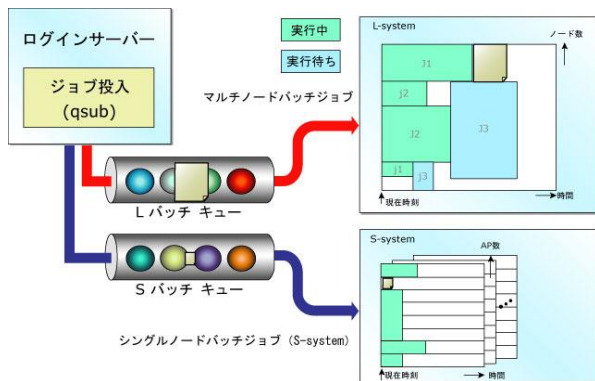


図 3 バッチ処理概念図

## 2.3 S系バッチジョブ

S 系バッチノード及び L 系バッチノードは、それぞれ単一のバッチキューにより管理される。利用者は、S 系バッチキューまたは L 系バッチキューのどちらかを選択する。

S 系バッチジョブは、1 台の計算ノード内で実行される。S 系バッチジョブのジョブスケジューリングには、使用する CPU 数と CPU 時間、メモリ使用量が用いられる。S 系バッチジョブの実行時、バッチスクリプト中では以下の指定を行う。

- 使用 CPU 数
- CPU 時間
- 使用メモリ量

S 系バッチジョブでは、8CPU まで使用した並列処理が可能であり、最大使用可能なメモリ量は約 127GB である。計算ノードは複数のバッチジョ

ブで共有される場合があるが、搭載する CPU 数またはメモリ量を超えてバッチジョブがアサインされることはない。また、S 系バッチノードではユーザディスクを NEC GFS によりマウントしており、ジョブ実行に際してステージングは行われない。

## 2.4 L系バッチジョブ

L 系バッチジョブは、複数台の計算ノードを用いた実行が可能である。1 つの計算ノードに複数のバッチジョブがアサインされることはなく、計算ノードを専有利用できる。運用上、L 系バッチジョブで最大使用可能なノード数は 128 ノードとしている。利用者が L 系バッチジョブの実行をシステムにリクエストする際、バッチスクリプト中に以下の指定を行う。

- 使用ノード数
- 経過時間
- ワークディスク容量
- ステージンファイル (ロードモジュールや入力ファイル等)
- ステージアウトファイル (出力ファイル等)

リクエストされたバッチジョブは、バッチ処理システムと密接に連携する運用管理ソフトウェアによって、研究課題で保有するノード時間との比較やステージインファイル指定の検査が行われた後に、ノード数と経過時間、及びワークディスク容量を元に、ジョブスケジューラが計算ノードへのアサインを行う。図 4 に、L 系バッチジョブの状態遷移概念図を示す。

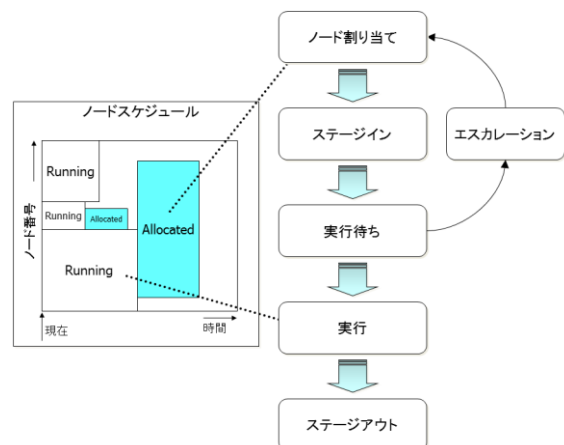


図 4 L系バッチジョブ状態遷移概念図

計算ノードの割当てには、FIFO + バックフィルのスケジューリングポリシーを用いる。計算ノードの使用効率を高めるためと、計算ノードに外乱を与えないために、ステージング処理は、IOCSと呼ぶ専用の装置によってジョブの実行に先んじて行われる。計算ノードへのバッチジョブ割当て状況は、スケジューリングマップと呼ぶ図によって表す。図5に、48時間後までのスケジューリングマップを示す。

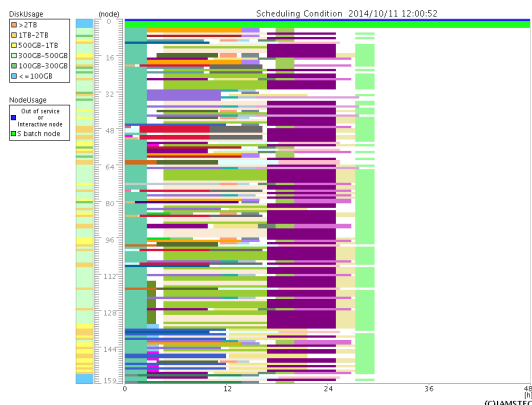


図5 スケジューリングマップ

## 2.5 ノード使用状況

ES2では、計算ノードの状態を以下のように定義する。

- RUN ジョブ実行中
- SAA ジョブ実行準備・後処理中
- ECO 省電力休止中
- SBY 待機中
- STP 停止中
- MTN 計画保守

運用管理ソフトウェアは、バッチ処理システムと連携し、定期的に計算ノードの状態を記録する。計算ノードへの外乱を排除するため、この処理に際して計算ノードへの問合せは行わない。図6に、計算ノードの使用状況を示す。これは2014年4月1日から10月26日までのノード状態を表している。計算ノードは、90%の割合でジョブ実行に利用されている。

スケジューリングマップや計算ノード使用状況は、計算資源の計画的な利用の一助となるよう、地球シミュレータ利用者サポートウェブページから閲覧可能としている。

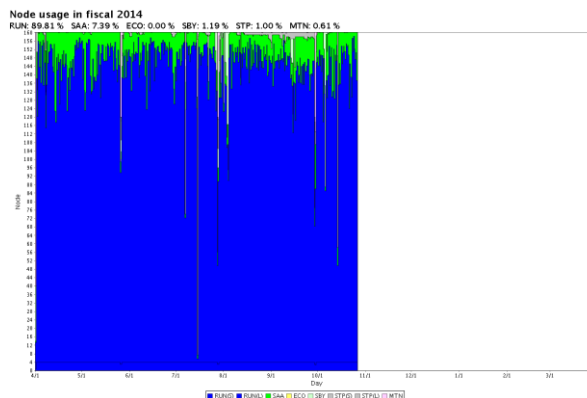


図6 計算ノード使用状況

## 2.6 大規模共有メモリシステム

ES2の周辺システムとして、2014年4月から大規模共有メモリシステムの利用が可能となっている。大規模共有メモリシステムは、SGI製UV2000を中心としたシステムである。商用アプリケーションを利用した大規模計算など、主に産業界からの利用を目的として導入された。CPUにはIntel E5-4650v2プロセッサを採用し合計2560CPUコアを搭載するとともに、32TBの巨大な共有メモリ空間を利用可能である。表1に大規模共有メモリシステムの概要を示す[3]。

表1 大規模共有メモリシステムの概要

理論性能	49.152TFLOPS
総プロセッサ(コア)数	256CPU socket (2,560core) 【Intel Xeon E5-4650v2(2.4GHz, 10core)】
総メモリ容量	32TB
OS	SUSE Linux Enterprise Server 11 SP3
ノード内接続	SGI NUMalink6

大規模共有メモリシステムの利用者情報や計算資源は、ES2の運用管理ソフトウェアによって一元管理される。バッチ処理での利用を基本とし、ES2のバッチ処理システムによって、統一的なジョブ管理と操作を実現している。

## 3 ES2 後継システム

現行システムであるES2の運用は、2009年3月から2015年2月末までの6年間の予定で行われている。2015年3月からは、NEC製SX-ACE

を中心とした後継システムの運用を開始する予定である。表2に、後継システムの概要を示す。また、図5にシステム概念図を示す。後継システムの運用では、大規模共有メモリシステムを含む周辺システムとの連携機能についても強化を図る予定である。

表2 後継システムの概要

システム		SX-ACE
CPU数(計算ノード数)		5120CPU・ノード
総ラック数(計算ノード部)		80
ノード構成	CPU名×数	専用CPU×1
	CPU周波数	1GHz
	CPUコア数	4
	コア最大性能	64GFLOPS
	CPU最大性能	256GFLOPS
	メモリ容量	64GB
システム総計	メモリバンド幅	256GB/s
	総合演算性能	1,3TFLOPS
	総メモリ容量	320TB
	総メモリバンド幅	1.3 PB/s

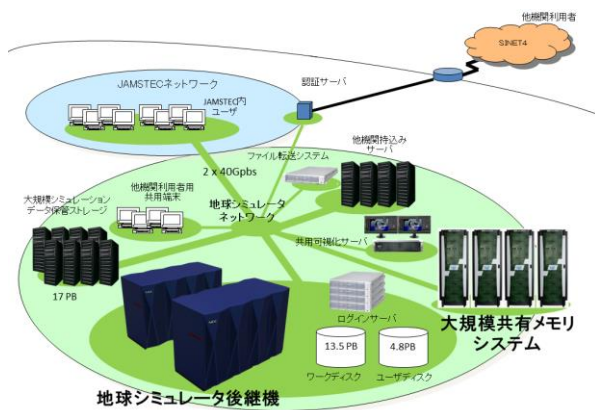


図5 地球シミュレータ後継システム概念図

後継システムの選定は、総合評価方式による一般競争入札によって行われた。総合評価方式とは、システムの性能・機能と入札価格を点数化し、最も高得点を獲得した者が落札者となる方式である。後継システムでは、アプリケーションプログラムの実行性能を重視し、性能評価試験の結果に最も多くの配点を与えている。ベンチマークプログラムには、ES2で実行されている代表的な以下の7本のアプリケーションプログラムを用いた。プログラムの実行に当たっては、提案システムへの最適化を許可している。

- ・ 全球大気海洋結合気候モデル (低解像度)
- ・ 全球大気海洋結合気候モデル (高解像度)
- ・ 全球雲解像大気大循環モデル

- ・ 3次元地震伝搬解析コード
- ・ 地震発生サイクルコード
- ・ 大気海洋結合モデル
- ・ データ同化コード

この性能評価試験の結果から、後継システムでは、単体プログラムの実行性能において現行システムの2倍から2.5倍、スループットにおいて4倍、システム全体の計算能力として8倍から10倍程度の向上が期待されている。

#### 4 おわりに

地球シミュレータ(ES2)では、ES2向けに最適化された運用管理ソフトウェアとジョブスケジューラにより、高い計算ノード利用率を実現してきた。後継システムにおいても、システム規模の拡大と性能の向上にあわせ、貴重な計算資源を有効活用し、かつ利用者にはより使いやすいシステムを構築するべく、準備が進められているところである。

#### 参考文献

- [1] 独立行政法人海洋研究開発機構地球シミュレータ開発史編集チーム、「地球シミュレータ開発史」、2010年12月
- [2] 地球シミュレータ、  
<http://www.jamstec.go.jp/es/jp/index.html>
- [3] 大規模共有メモリシステム、  
<http://www.jamstec.go.jp/es/jp/uv/2000index.html>