

# エネルギー効率を考慮した電力制約下でのスループット指向ジョブスケジューリング

黄 巍<sup>1,a)</sup> 岩澤 直弘<sup>1,2</sup> カオ タン<sup>1,2</sup> 和 遠<sup>1,2</sup> 近藤 正章<sup>1,2</sup> 中村 宏<sup>1</sup>

アブストラクト：HPC データセンタにおいては、消費電力がシステムの設計や性能において最も重要な制約条件となると考えられている。そこで従来から、ある電力制約の下でシステム全体のスループットを最大化するためのジョブスケジューリングに関する研究が行われてきた。しかし、これまでは各ジョブのエネルギー効率についてはあまり考慮はされてこなかった。本稿では、システム全体の電力制約下で使用電力量を考慮しつつ、ジョブ毎の電力制約値及び使用ノード数を最適化することで、各ジョブに対してシステムスループットだけではなくエネルギー効率も同時に考慮したジョブスケジューリング手法の構築を検討する。予め用意された各ジョブの実システム上の性能・電力データをもとにしてスケジューリングを行うシミュレータを作成し実験を行い、従来のスループット重視のスケジューリングに比べ、消費エネルギーを大きく削減できることがわかった。

## 1. はじめに

現時点で最高性能を持つスーパーコンピュータシステムの Tianhe-2 は 33PFLOPS 超の性能を達成しているが、そのために 18MW 近い消費電力を必要とする [1]。2020 年あたりに実現されると予想されるエクサスケール級の HPC システムは、20~30MW と同程度の電力で現在世界トップクラスのスーパーコンピュータの 30~50 倍近い性能向上が求められる。そのため、将来の HPC システムでは、消費電力がシステムの設計や実効性能を制約する最大の要因になると考えられており、限られた電力資源を有効利用することは、今後の大規模 HPC システムにおける重要な研究課題の一つである。

これまで、HPC システム上で実行される各並列アプリケーションの電力効率を向上させるための手法や [4]、電力制約下で性能向上を図る手法についての研究が行われてきた [3], [8]。一方で、大規模なスーパーコンピュータシステムでは、同時に複数のジョブが実行されることが通常であり、アプリケーション個々の電力最適化のみならず、システム全体の電力制約を考慮したスケジューリングや電力配分も重要である。そこで、近年では電力制約下でシステムのスループットを最大化するための研究も行われている [2], [9]。

従来の電力制約下での複数ジョブ実行を考慮したスケ

ジューリングに関する研究では、スループットを向上させることのみが目的で、各ジョブのエネルギー効率に関してはあまり考慮されてこなかった。一方で、システムの電力が増大するにつれ、運用にかかる電気料金としてのコストは非常に大きく、HPC データセンタにおいてもジョブ実行に必要なエネルギーも性能同様に考慮すべき重要な指標である。

そこで、本稿では各ジョブのエネルギー効率を考慮しつつ、電力制約下でのシステムのスループットを向上させるためのジョブスケジューリング手法を検討する。具体的には、文献 [2] で提案されているジョブ毎の電力制約値と使用ノード数を最適化することで、電力制約下でのスループットを最大化させる手法をベースとし、エネルギー効率を最適化の際の評価関数に加えることで、性能とエネルギーの両者を考慮したスケジューリング手法を構築することを目指す。文献、および実ベンチマークの実行結果から得られた性能・電力モデルを基に、シミュレーションによる評価を行った結果、構築したスケジューリング手法によりジョブ実行の消費エネルギーを削減できることがわかった。また、システムが性能重視か消費エネルギー重視かに応じて、スケジューリング機能を調節できることがわかった。

## 2. 関連研究

従来から、ある 1 つの並列ジョブを対象とした電力性能最適化に関する研究は多数行われてきた。例えば、与えられた消費電力の上限を超えない範囲で、システム内の各要素間で電力を適切に配分することで、アプリケーションの実効性能を最大化する研究も行われている。文献 [3] では同期ポイントにおける各スレッドや各プロセスの待ち時

<sup>1</sup> 東京大学大学院情報理工学系研究科  
Graduate School of Information Science and Technology,  
The University of Tokyo

<sup>2</sup> 独立行政法人科学技術振興機構 CREST  
CREST, JST

a) huang@hal.ipc.i.u-tokyo.ac.jp

間を考慮して、各ノード・CPUの動作周波数を設定することで、アプリケーションの性能を低下させずに消費電力を削減する手法が提案されている [4], [7]. また、今村らはメニーコア・プロセッサを対象に、並列アプリケーション実行時の性能を最大化するための Dynamic Core-count and Frequency Scaling (DCFS) を提案している [5]. これは、プログラムの特性に応じてコア数と周波数を動的に変更し、消費電力制約下における性能を最大化するものである.

他にも、消費電力制約下での性能最適化に関する研究はいくつか行われている. 近年の Intel 社のプロセッサには、プロセッサ部と DRAM 部の消費電力を観測・制御するためのインタフェースである RAPL (Running Average Power Limit)[16], [17] が備えられている. 文献 [6] では、この機能を利用して、使用ノード数の最適化と CPU・DRAM 間での電力配分を行う手法を提案している. また、HPC 計算分野で広く用いられる GPU の消費電力は CPU と比較しても大きい、消費電力あたりの性能は高い. CPU と GPU 間に割り当てる処理量と、それぞれの動作周波数をあわせて変更することで、与えられた消費電力制約下で GPU アプリケーションの性能を最大化するための手法も提案されている [8].

個々のジョブの電力性能最適化ではなく、複数ジョブの電力性能を最適化することを目的に、Etinski らは計算機システムに供給される電力資源を最大限利用するためのスケジューリング手法を提案している [9]. Dynamic Voltage Frequency Scaling (DVFS) を利用して、HPC 計算機システム内で同時に実行されている各ジョブのプロセッサの周波数・電源電圧を変更することで、各ジョブを適切な消費電力で動作させつつ、最大限に電力を使用するものである. ただし、一般的にはメモリアクセス頻度などに依存して、同じ周波数を用いても各ノードの消費電力はジョブ毎に異なる可能性がある. また、当該論文の目的は電力資源の最大限の利用であるため、ジョブの実行速度や消費エネルギーについては十分に考慮されていない.

複数ジョブが実行されるシステム上で、電力配分の最適化とともに、ジョブの配置に関する研究も行われている. 特に、各ノードの冷却コストを考慮したジョブ配置の研究として、冷却効率向上のために休止状態ノードを配置しつつ、各ジョブを実行するノード群をネットワーク的に隣接して配置するジョブ配置ポリシーも提案されている [10], [11].

これに対して、Sarood らは消費電力制約下で、システム全体のスループットを最大化するためのスケジューリング手法を提案した [2]. このスケジューリング手法により、ジョブごとに適切なノード数・電力キャップ値を設定することができ、システム全体のスループットを最大化することが可能となる.

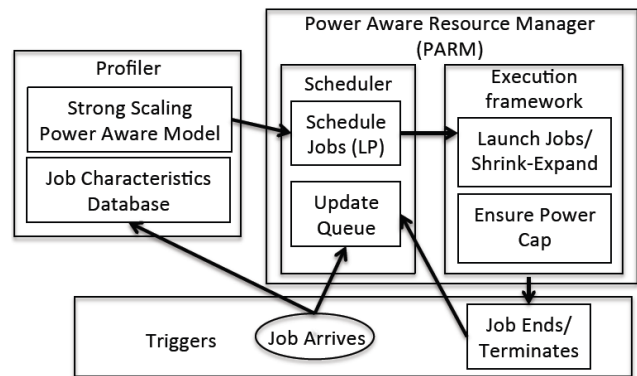


図 1 Power Aware Resource Manager の概要 (出典 [2])

本稿は、システム全体のスループットの最大化と、エネルギー効率の向上の両者を考慮したジョブスケジューリング、および電力配分手法を提案するものである. この点で、従来の研究とは異なるものである.

### 3. 電力制約下でのジョブスケジューリング

本章では、従来の電力制約下でのジョブスケジューリング手法として文献 [2] で提案されている手法について述べる.

#### 3.1 ジョブスケジューリング手法の概要

文献 [2] では、ある計算機センターのシステムが利用可能な総消費電力の制約が与えられた際に、システム全体のスループットを最大化するためのスケジューリング手法である *Power Aware Resource Manager (PARM)* を提案している. PARM の概要を図 1 に示す. PARM はイベント発生時にどのジョブをどのような設定で実行するかを決定するスケジューラ (*scheduler*) と、実際に与えられた設定で指定されたジョブを実行させる実行フレームワーク (*Execution framework*) で構成される. スケジューラは新たなジョブがキューにサブミットされた際といずれかのジョブの実行が完了、またはエラーなどの原因で終了した際をトリガーとして起動される. スケジューラは既に実行中の全ジョブ、およびキューで待ち状態のジョブの中から新たに実行するジョブを選択し、各ジョブの実行設定として利用ノード数と使用電力の上限 (電力キャップ) を決定する. なお、あるジョブが一旦実行を開始すると、そのジョブは完了まで止めないと仮定する. キューの中からどのジョブを選択するか、および各ジョブにどのような実行設定を与えるかには多数の候補があり得るが、組み合わせ最適化問題を解くことで計算機センターのスループットを最大化する最適解を求めることが PARM の役割である.

一方で、実行フレームワークはスケジューラから与えられた設定を用い、各ジョブを指定ノード数と指定電力キャップ値で実行させる. 電力キャッピングの実現には、Intel 社のプロセッサに備えられている RAPL を利用する

ことを想定する。通常、実行ノード数を動的に変更することは、大規模並列処理を行うプログラムとして一般的に利用されている並列プログラミングインタフェースの MPI を用いた場合には簡単ではない。そこで、文献 [2] では、キュー内のジョブがシステム上で実行される際に、ある範囲からノード数を選択して実行することが可能なケース (Moldable ジョブ)、および Charm++ [18], [19] などの runtime を用いることで実行時にも使用ノード数を変更可能なケース (Malleable ジョブ) を主に仮定している。

### 3.2 最適化問題としての定式化

前述のように、スケジューラはトリガーイベントの発生の度に組み合わせ最適化問題を解く必要がある。文献 [2] で用いられている整数線形計画問題としての定式化を式 (3.1)-(3.5) に示す。なお、当該問題への入力は、実行中、あるいはキュー内で実行待ちのジョブの集合  $\mathcal{J}$ 、あるノード数  $n$  とある電力制約  $p$  で実行した際の各ジョブ  $j \in \mathcal{J}$  の実行時間であり、出力は実行すべきジョブとその実行設定  $(n, p)$  である。また、表 1 に式中で利用する記号の意味を示す。

目的関数

$$\sum_{j \in \mathcal{J}} \sum_{p \in \mathcal{P}_j} \sum_{n \in \mathcal{N}_j} s_{j,n,p} * w_{j,n,p} * x_{j,n,p} \quad (3.1)$$

2 値変数  $x$  の満たす条件

$$\sum_{p \in \mathcal{P}_j} \sum_{n \in \mathcal{N}_j} x_{j,n,p} \leq 1 \quad j \in I \quad (3.2)$$

$$\sum_{p \in \mathcal{P}_j} \sum_{n \in \mathcal{N}_j} x_{j,n,p} = 1 \quad j \in \mathcal{I} \quad (3.3)$$

制約条件

$$\sum_{j \in \mathcal{J}} \sum_{p \in \mathcal{P}_j} \sum_{n \in \mathcal{N}_j} n * x_{j,n,p} \leq N \quad (3.4)$$

$$\sum_{j \in \mathcal{J}} \sum_{p \in \mathcal{P}_j} \sum_{n \in \mathcal{N}_j} (n(p + W_{base}))x_{j,n,p} \leq W_{max} \quad (3.5)$$

本定式化においては、ジョブ  $j$  をノード数  $n$ 、電力  $p$  で実行した場合の最小構成に対するスピードアップを見積もる必要がある。この組み合わせは非常に膨大な数になるため、いかに効率よく見積れるかが重要である。そこで、文献 [2] では *Strong Scaling Performance Model* という性能モデルを導入している。以下、その性能モデルについて説明する。

まず、ジョブ  $j$  のスピードアップを式 (3.6) のように定義する。

$$s_{j,n,p} = \frac{t_{j,\min(N_j, \min(N_j), \min(P_j))}}{t_{j,n,p}} \quad (3.6)$$

ここで、 $s_{j,n,p}$  はジョブ  $j$  をノード数  $n$ 、各ノードの CPU 電力  $p$  で実行した場合の実行時間と、ジョブ  $j$  が実行可能

表 1 定式化で用いる記号

記号	説明
$N$	ノードの総数
$J$	すべてのジョブの集合
$\mathcal{I}$	現在実行中のジョブの集合
$I$	キュー内にある実行待ちジョブの集合
$\mathcal{J}$	サブミットされたが、実行し終わっていないジョブの集合
$\mathcal{N}_j$	ジョブ $j$ が実行できるノード数の集合
$\mathcal{P}_j$	ジョブ $j$ が実行できる電力設定の集合
$w_j$	ジョブの優先度付けのための重み
$\alpha$	$w_j$ の中で、優先度と性能を調節するための変数
$x_{j,n,p}$	ジョブ $j$ が $n$ ノード上で電力 $p$ で実行する場合は 1, その他の場合は 0
$t_{now}$	現在の時刻
$t_j^a$	ジョブ $j$ がサブミットされた時刻
$W_{base}$	マシンのベース電力
$t_{j,n,p}$	ジョブ $j$ をノード数 $n$ 、電力 $p$ で実行した場合の実行時間
$s_{j,n,p}$	ジョブ $j$ をノード数 $n$ 、電力 $p$ で実行した場合の最小構成に対するスピードアップ

な最小ノード数  $\min(N_j)$  と最小 CPU 電力  $\min(P_j)$  で実行した際の比である。式 (3.1) で示した目的関数は、全実行ジョブのスピードアップ  $s_{j,n,p}$  を足し合わせたものである。そのため、目的関数を最大化することは、実行ジョブ全体を高速に実行するのに適するノード数と電力キャップ値を各ジョブに対して選ぶことに相当する。 $t_{j,n,p}$  の性能モデルの詳細に関しては、文献 [2] を参照されたい。なお、4.2 節では、文献 [2] の性能モデルを本稿のために改変したモデルについて述べる。

目的関数として  $s_{j,n,p}$  のみを使用すると、大規模なジョブを実行するよりも複数の小さなジョブを実行する場合の方が目的関数の値が大きくなり、大規模ジョブが実行され難くなる傾向がある。その結果、大規模ジョブの待ち時間が長くなり、ジョブ間の公平性が悪化する。この対処のために、キューで長く待たされているジョブの優先度を高くするべく、以下のように重み  $w_j$  を定義する。

$$w_j = (t_{j,\min(N_j), \min(P_j)}^{rem} + (t_{now} - t_j^a)^\alpha)^\alpha \quad (3.7)$$

ここで、 $t_{j,\min(N_j), \min(P_j)}^{rem}$  は、ジョブ  $j$  を最小ノード数  $\min(N_j)$ 、各ノードの CPU を最小電力  $\min(P_j)$  で実行させた際に、現時刻からジョブ完了までかかる残り時間の見積もり値である。これとジョブ  $j$  がサブミットされた時間を足し合わせたものをジョブの重要度とすることで、より実行時間の長いジョブ、あるいは待ち時間の長いジョブの優先度が高くなる。なお、 $\alpha$  は 0 から 1 までの値域を持つ変数であり、0 に設定した場合は性能重視で公平性を無視したスケジューリング、1 に設定した場合は公平性を重視し性能低下を許容するスケジューリングとなる。

### 3.3 従来のジョブスケジューリング手法の課題

前節で述べた最適化の目的関数として、スピードアップ

表 2 エネルギーの定式化に用いる記号

記号	説明
$Er_{j,n,p}$	消費エネルギーの削減度合い
$E_{j,n,p}$	ジョブ $j$ をノード数 $n$ , 電力 $p$ で実行した場合の消費エネルギー
$\beta$	性能と消費エネルギーの重要度の重み

を表す  $s_{j,n,p}$  のみが存在することからもわかるよう、従来の電力制約下でのジョブスケジューリング手法は、スループットを向上させることのみが目的であり、各ジョブのエネルギー効率に関しては考慮されていない。一方で、計算機センターで消費される計算機システムの電力が増大するにつれ、運用にかかる電気料金コストは非常に大きく、ジョブ実行に必要な消費エネルギーも性能と同様に考慮すべき重要な指標である。例えば、スループット最大化のために、高いクロック周波数を用いてプロセッサを動作させるよりも、低い周波数と電源電圧を用い、消費電力を抑えて実行する方が一般的にエネルギー効率は高くなる。そのため、エネルギー効率も考慮した場合には、従来のスケジューリングとは異なるジョブを異なる設定で実行させる場合が最適となる可能性がある。

#### 4. エネルギーを考慮した電力制約下でのジョブスケジューリング

本稿では、従来のスケジューリング手法の課題をふまえて、PARM スケジューリング手法の拡張を検討する。一般的に、性能とエネルギー効率にはトレードオフの関係がある。計算機センターの運用ポリシーやエネルギーコストの違いにより、性能とエネルギー効率のどちらを重視するかは運用ポリシーに依存すると考えられる。そこで、性能と消費エネルギー効率のどちらをどれだけ重要視するかを重みとして与えられるようにしつつ、両指標を考慮したスケジューリング手法の構築を目指す。

##### 4.1 目的関数の拡張

3.2 節で述べた PARM 定式化の拡張として、消費エネルギーに関する項目を目的関数に含めることを考える。拡張した目的関数を式 (4.1) に示す。また、ここで新たに使用する記号の説明を表 2 に示す。

$$\sum_{j \in \mathcal{J}} \sum_{p \in \mathcal{P}_j} \sum_{n \in \mathcal{N}_j} (\beta * s_{j,n,p} + (1 - \beta) * Er_{j,n,p}) * w_{j,n,p} * x_{j,n,p} \quad (4.1)$$

ここで  $Er_{j,n,p}$  はジョブ  $j$  を設定  $(n, p)$  で、すなわちノード数  $n$ , 電力  $p$  で実行した場合に対する、最小ノード数と最小電力で実行した場合のエネルギー削減率であり、以下のように定義できる。

$$Er_{j,n,p} = \frac{E_{j,\min(N_j),\min(P_j)}}{E_{j,n,p}}$$

$$\begin{aligned} &= \frac{t_{j,\min(N_j),\min(P_j)} * \min(N_j) * (\min(P_j) + W_{base})}{t_{j,n,p} * n * (p + W_{base})} \\ &= s_{j,n,p} * \frac{\min(N_j) * (\min(P_j) + W_{base})}{n * (p + W_{base})} \quad (4.2) \end{aligned}$$

$(n, p)$  で実行した場合の消費エネルギーが少ないほど  $Er_{j,n,p}$  は大きくなることから、スピードアップ  $s_{j,n,p}$  が大きいほど、また  $Er_{j,n,p}$  が大きいほど目的関数の値が大きくなり、それを最大化するように最適化することで、エネルギーを考慮したスケジューリングを行う事ができる。

なお、式 (4.1) にある係数の  $\beta$  は性能と消費エネルギー効率を調整する重みである。 $\beta$  を設定することで、性能重視か消費エネルギー効率重視かを変更することができる。 $\beta$  が 0 の場合では、従来研究と同じ目的関数となり、完全に性能を重視したスケジューリングとなる。一方、 $\beta$  が 1 の場合は、完全に消費エネルギー効率を重視したスケジューリングとなる。本稿では、 $\beta$  がジョブに対して個別に決定されるのではなく、一律に同じ値を用いることにする。

##### 4.2 スピードアップモデル

本節では、目的関数にあるスピードアップを見積もるための性能モデルについて説明する。文献 [2] でも性能モデルが提案されているが、より簡易にモデルを構築・評価できるように、本稿では一部性能モデルを変更する。

ここでのスピードアップモデルは、以下のようにノード数・実行性能モデルと電力・実行性能モデルの二つに分けて構成される。

$$s_{j,n,p} = s_j(n) * s_j(p) \quad (4.3)$$

以下、それぞれについて詳述する。

###### 4.2.1 ノード数・実行性能モデル

あるジョブは、平均並列度  $A$  と並列度の分散  $\sigma$  を用いて、使用ノード数を変えさせた場合の実行時間  $t(n)$  を次のように表現できる [12]。

$$t(n) = \begin{cases} \frac{T_1 - \frac{T_1\sigma}{2A}}{n} + \frac{T_1\sigma}{2A} & n < A \quad (4.4a) \\ \frac{\sigma(T_1 - \frac{T_1}{2A})}{n} + \frac{T_1}{A} - \frac{T_1\sigma}{2A} & A \leq n \leq 2A - 1 \quad (4.4b) \\ \frac{T_1}{A} & n \geq 2A - 1 \quad (4.4c) \end{cases}$$

ここで、 $n$  はノード数であり、 $T_1$  はジョブが単一ノードで実行される場合の実行時間である。式 (4.4a) では、ノード数が  $A$  より小さい場合には、ノード数に対するスケラビリティが高いことを示しており、式 (4.4b) では、ノード数が  $A$  より大きい場合にはスケラビリティが大幅に減少することを意味している。さらに、式 (4.4c) にあるように、 $n$  が  $2A - 1$  よりも大きい場合は実行時間が  $T_1/A$  のままとなり、ノード数を追加しても実行時間は削減されな

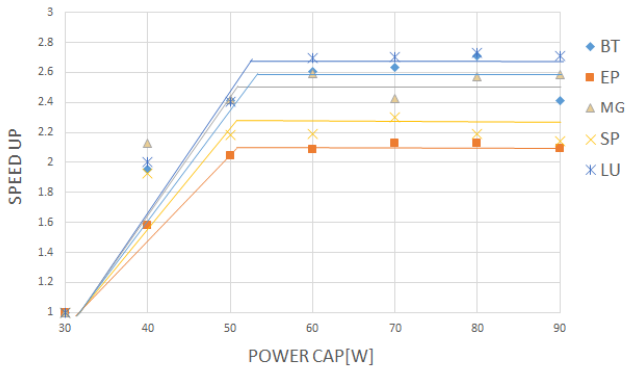


図 2 16 コアでの各ジョブの電力・実行速度モデル

いことを示している。このモデルにより、ジョブの特徴として  $A, \sigma, T_1$  を求めることができれば、任意のジョブに対してノード数を変化させた際の性能を見積もることが可能となる。ここで、スピードアップ  $s_j(n)$  は  $t(n)/T(1)$  として表すことができる。

なお、次章の評価では実行するジョブの  $A, \sigma$  の値は文献 [12] に記載されている値を用いる。

#### 4.2.2 電力・実行性能モデル

本研究では RAPL を用いて、CPU パッケージの電力ドメインに対して電力キャップ値を設定し、実際にジョブを実行して実行時間を測定することで、電力・実行性能モデルを構築する。

図 2 に Nas Parallel Benchmarks (NPB) 中の 5 種類のプログラム (BT, MG, SP, LU, EP) を選択し、1 ノード 16 コアのシステムにおいて電力キャップ値を 30W から 90W に変化させた場合のスピードアップを示す。なお、30W 制約時の実行時間を規準とした場合のスピードアップである。図より、電力キャップ値を大きくしていくと、5 種類のプログラムとも 50W 電力キャップまでは性能向上が見られるが、50W より高いキャップ値では、速度向上が見られない。また、30W から 50W の間では、電力キャップ値に対してほぼ線形にスピードアップすることがわかる。本研究では、次章の評価にこの 5 種類のジョブを用い、電力・実行性能モデル  $s(j, p)$  を次のように定義する。

$$s_j(p) = \begin{cases} (p - 30) * m_j + 1 & p < 50 \quad (4.5a) \\ s_j(50) & 50 \leq p \quad (4.5b) \end{cases}$$

ここで、 $m_j$  はジョブ  $j$  の電力キャップの値が 30W から 50W の間の近似線の傾きである。

## 5. 評価実験

本稿で検討したエネルギー効率を考慮したスケジューリング手法の効果を調べるために、前節で示したモデルとスピードアップに関するパラメータを利用して、シミュレーションにより評価する。評価では、前節で述べた NPB の 5

種類のプログラムのの中からジョブをランダムに選択して生成するジョブセットと、より大規模の HPC システムを評価するために、Intrepid[13] のジョブ実行ログを用いて生成したジョブセットの二種類を用いて実験を行う。また、目的関数のパラメータである  $\beta$  を変化させながらシミュレーション評価を行い、その結果を単純な FIFO ポリシー、および文献 [2] の従来手法でスケジューリングした場合と比較する。

### 5.1 実験環境

本実験用に、ジョブ実行のシミュレーション環境を独自に構築した。組み合わせ最適化問題の解法として、整数線形計画法 (従来手法)、および非線形整数計画法 (本稿での拡張手法) を解くと、探索に時間がかかり解が得られないことがあるため、今回は遺伝的アルゴリズムを利用して近似最適解を求めることとした。遺伝的アルゴリズムのパラメータは、1 世代の個体数を 300、親遺伝子数を 20、世代数を 100、突然変異率を 90%、組み換え方式は多点交叉とした。遺伝的アルゴリズムによる最適化にかかる時間は 1~5 秒程度であり、実用的には十分短い時間であると考えられる。

ジョブセット 1: NPB の BT, MG, SP, LU, EP からランダムに 50 個を選択してジョブセットを生成し評価する。表 3 は実験に用いた各プログラムのパラメータである。これらは、Intel Xeon E5-2670-v2 2.50GH を 2 ソケット搭載する計算機で実行することで得られたデータ、および文献 [12] に記載されたデータを基に作成した。また、ジョブがサブミットされるタイミングは平均到着率  $\lambda=10$  のポアソン分布に従うと仮定した。NPB は実際には実行時に使用ノード数を変更することはできないが、本実験ではスケジューリングの効果を検証するために、各ジョブが実行時に使用ノード数を変更可能な Malleable ジョブを仮定して評価を行った。

表 3 プログラムのパラメータ ( $\min(P)=30W, \min(N)=1$ )

ジョブ名	$t_{\min(P), \min(N)}$	$A$	$\sigma$	$m$
BT	291.747[s]	7	0.02	0.055
MG	22.473[s]	18	0.04	0.070
SP	297.194[s]	20	0.05	0.080
LU	228.374[s]	18	0.01	0.075
EP	43.867[s]	Inf	0	0.065

ジョブセット 1 の評価対象 HPC 計算機システムの仮定として、16 コアの CPU を持つ計算機を 1 ノードとし、最大ノード数は 64 ノードとした。なお、CPU ソケット以外に常に消費されるベース電力を 10W と仮定し、評価では CPU の電力のみを対象に評価を行う。この際に、最大許容電力は 3850W に設定して実験を行った。

ジョブセット 2: 大規模なシステムを対象として評価す

るため、Intrepid のジョブ実行ログ [14] に基づき生成した。Intrepid はアルゴンヌ国立研究所にある 40,960 ノード、163,840 コアからなる HPC システムであり、ジョブセット 2 ではこれと同じ環境の HPC 計算機システムを仮定し、実験を行った。なお、この際に、最大許容電力は 2.048MW と設定した。

本ログに含まれるジョブトレースには、8ヶ月間で実行された 68,936 個のジョブの実行ログがあるが、本実験ではこの中の 1000 個分を利用して実験を行う。また、本ログには、電力や並列度の情報は含まれないため、前述の NPB の特徴をランダムに選択し、それらをスケールさせて各ジョブ当てはめることで、評価を行った。

## 5.2 評価結果

比較評価のための FIFO スケジューリングには各ジョブの実行設定  $(n, p)$  を最適化する枠組みはないため、ノード数  $n$  と電力キャップ  $p$  を固定して評価を行う。この際、各ジョブを最小ノード数  $\min(N_j)$  と各ノードの CPU を最小電力  $\min(P_j)$  に設定して実行する場合  $FIFO_{min}$ 、各ノードの CPU を最大電力  $\max(P_j)$  で実行できるようにノード数を設定した場合  $FIFO_{max}$ 、最大 CPU 電力と最小 CPU 電力の中間の電力で実行できるようにノード数を設定した場合  $FIFO_{mid}$  のそれぞれを評価した。

まず、ジョブセット 1 の結果について議論する。図 3 に、作成したジョブセットに対して  $FIFO_{max}$ 、 $FIFO_{mid}$ 、 $FIFO_{min}$ 、およびスループットを重視した従来のスケジューリング (Performance Oriented)、本稿で拡張したエネルギーを重視したスケジューリング (Energy Oriented) のそれぞれにおける総実行時間と総消費エネルギーを示す。ここで、総実行時間とは、最初のジョブがキューにサブミットされてから最後のジョブの実行が完了するまでの時間である。

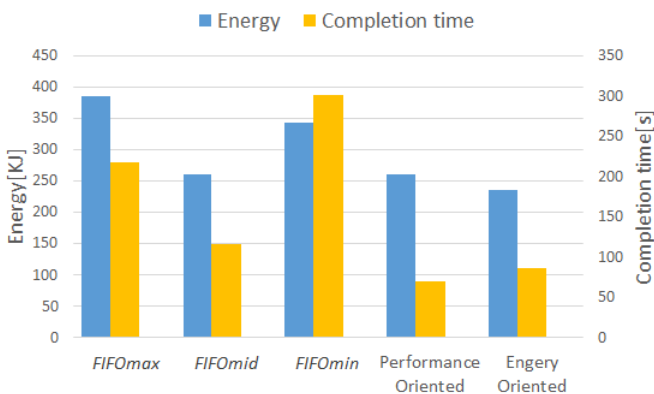
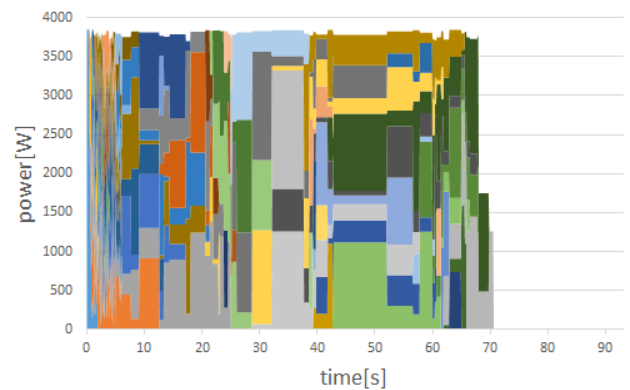


図 3 総実行時間と総消費エネルギーの比較 (ジョブセット 1)

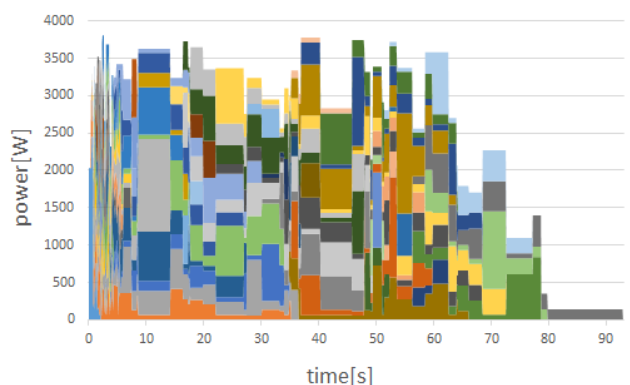
図 3 より、FIFO では単純にジョブのサブミット順に実行され、またノード数と電力キャップ値も固定の値が使われ

ることから、実行時間が Performance Oriented や Energy Oriented に比べて長く、また消費エネルギーも大きいことがわかる。 $FIFO_{max}$  よりも  $FIFO_{mid}$  の実行時間が短いのは、 $FIFO_{max}$  は平均並列度  $A$  を超えても多くのノードを利用して実行するために、システム全体のノードの並列化効率が悪くなったためと考えられる。常に消費されるベース電力があるため、実行時間が長くなると消費エネルギーも増加する傾向にある。

性能重視のスケジューリング (式 (3.1) の目的関数を使用) では、FIFO の中では最も性能が良い  $FIFO_{mid}$  に対しても、総実行時間が 39% も減少していることがわかる。しかし、総消費エネルギーは  $FIFO_{mid}$  とほぼ同等である。一方、本稿で検討したエネルギーを考慮したスケジューリングで省エネルギー重視の戦略をとった場合 ( $\beta=0.0$ ) には、性能重視のスケジューリングに比べて総実行時間は長くなっているものの、約 10% もの消費エネルギーを削減することができている。また、総実行時間も  $FIFO_{mid}$  より 25% 減少していることがわかる。この結果より、本稿で拡張したスケジューリングは、計算機センターにおいて大きなコスト要因となる消費エネルギーの削減に効果的であると結論付けることができる。



(a) 性能重視 ( $\beta = 1.0$ ) の場合



(b) 消費エネルギー重視 ( $\beta = 0.0$ ) の場合

図 4 各ジョブに分配された電力

表 4 待ち時間と実行時間の比較

	$FIFO_{max}$	$FIFO_{mid}$	$FIFO_{min}$	Performance oriented( $\beta=1.0$ )	Energy oriented( $\beta=0.0$ )
Avg. Wait Time[s]	54.57	25.11	0.00	17.21	20.87
Max Wait Time[s]	213.29	106.32	0.00	61.66	55.08
Var. of Wait Time	4012.98	909.07	0.00	406.06	370.44
Avg. Relative Run Time	1.00	1.46	31.39	2.13	3.40

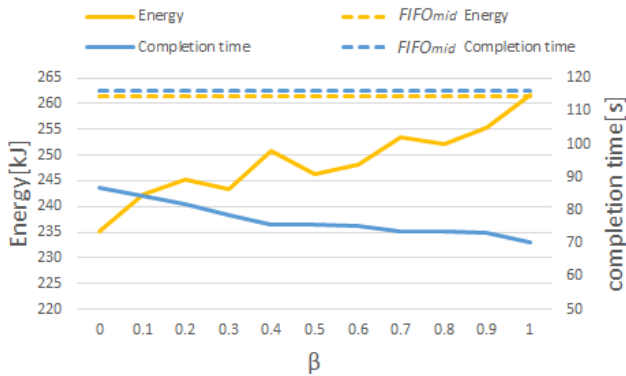


図 5  $\beta$  を変化させた場合の結果

図 4 はジョブセット 1 における各ジョブに分配された電力を示す。図より、 $\beta=1.0$ 、すなわち性能重視(従来手法と同様)の場合は最大許容電力に到達することが多いが、 $\beta=0.0$  の消費エネルギー重視の場合では最大許容電力に到達していない時間の割合が多いことがわかる。この原因は次のように考えられる。従来手法ではスループット最大化が目的であるため、ジョブを実行しないよりは、エネルギー効率が悪くてもジョブを投入して実行しようとするため、電力が常に最大限使用される。しかし、少ない電力で無理にジョブを実行させようすると実行時間が長くなり、ベース電力のために合計の消費エネルギーが増大する可能性がある。したがって、消費エネルギー重視の場合、余剰電力があってもジョブを実行しない場合が有利と判断したためである。

表 4 は、各手法でスケジューリングした場合の最小待ち時間、最大待ち時間、待ち時間の分散および  $FIFO_{max}$  に対するジョブの相対実行時間を示している。表 4 により、 $FIFO_{max}$  ではジョブの平均実行時間は短いものの、待ち時間に関しては最大となつている。一方で、 $FIFO_{min}$  は今回のシステム設定のもとでは、全ジョブを実行したとしても余剰電力があるため、待ち時間はすべて 0 となっている。しかし、ジョブの実行時間は非常に長くなってしまふ。 $FIFO_{min}$  を除くと、 $\beta = 1.0$  の従来手法では平均待ち時間が最小となり、また  $\beta = 0.0$  の提案手法の場合では、最大待ち時間、および待ち時間の分散は最も小さくなっている。これにより、提案スケジューリング手法では各ジョブを公平にスケジューリングできていることが分かる。

図 5 は、 $\beta$  の値を 0.0~1.0 の間で 0.1 刻みで変化させな

がらシミュレーションを行った結果である。同じく、総実行時間と総消費エネルギーを示している。参考に、FIFO で最も性能の良かった  $FIFO_{mid}$  の結果も点線で示している。

図 5 より、 $\beta$  の値が小さくなるにつれ消費エネルギー重視のスケジューリングになることから、消費エネルギーが減少している。一方で、総実行時間は長くなる傾向にある。図の中にこれらの傾向と違って変化する部分があるが、これは遺伝的アルゴリズムにより生じる誤差である。このように、 $\beta$  のパラメータ値を変更することで、各計算機システムの運用に合わせて性能と消費エネルギーを調整することができる。これも本稿で拡張したスケジューリング手法の利点になると考えられる。

図 6 はジョブセット 1 を利用し、到着率  $\lambda$  の値を 10 以外に 1 と 0.2 に変化させて実験を行った結果である。図 6 により、 $\lambda$  の値が小さくなると性能重視の場合の実行時間の削減効果は低くなるが、エネルギー重視の場合の消費エネルギー削減度は  $\lambda$  によってほとんど変わらないことがわかる。これはジョブの到着率の減少により、次のジョブが到着するまでの時間が長くなることから、個々のジョブの実行時間が伸びても総実行時間への影響が少なくなるためである。

次に、ジョブセット 2 の結果について議論する。図 7 に、 $FIFO_{max}$ 、 $FIFO_{mid}$ 、およびスループットを重視した従来のスケジューリング (Performance Oriented)、本稿で拡張したエネルギーを重視したスケジューリング (Energy Oriented) のそれぞれにおける実行時間と消費エネルギーを示す。なお、ジョブセット 2 では各ジョブサイズが大きく、最後に実行されるジョブの実行時間により総実行時間が大きく変化してしまう。HPC システムでは常にジョブがサブミットされ、最後のジョブの実行が完了するまでの時間はあまり重要ではなく、ジョブ実行のスループットが重要視される。そこで、1000 ジョブのうち最初の 950 ジョブが終了した際の時間 (スループットを評価することとほぼ等価となる) を評価指標とする。また、最小ノード数で実行する  $FIFO_{min}$  ではノード数が 40,000 ノード以上あるシステムでは、ノード使用率が 0.03% と非現実的な値になってしまうため、ここでは結果を省略する。

図 7 より、 $FIFO_{min}$  と比較すると、スループットを重視したスケジューリング手法では実行時間が 71% が減少

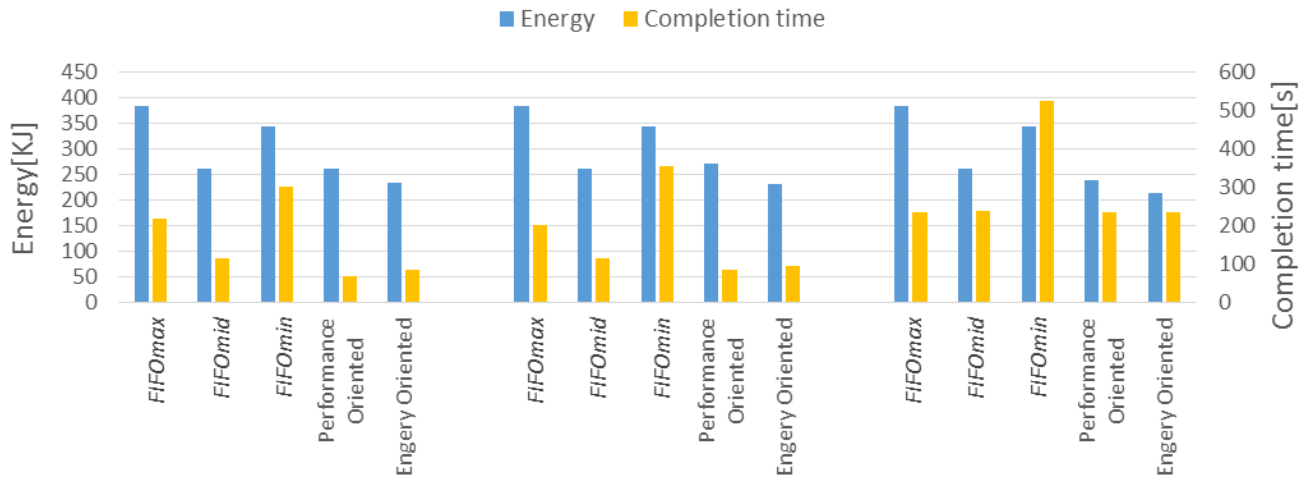


図 6 λ を変化させた場合の結果

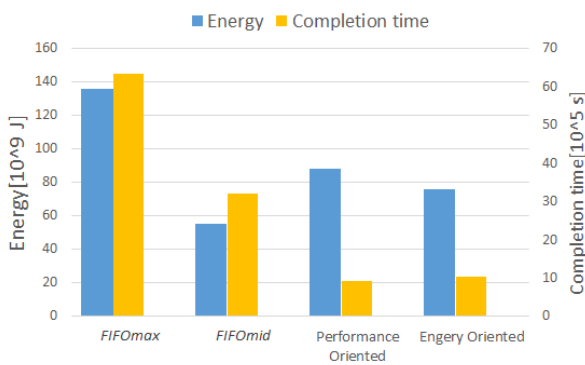


図 7 総実行時間と総消費エネルギーの比較 (ジョブセット 2)

し、エネルギー重視したスケジューリング手法は実行時間が68%減少することがわかる。また従来手法に比べ、本稿で検討しエネルギーを重視した場合は14%のエネルギー削減ができた。これにより、大規模なHPCシステムにおいても本稿で検討したエネルギー重視手法は効果があることがわかった。

## 6. おわりに

本稿では、性能とエネルギー効率の両方を考慮した、HPC 計算機システム向けの電力制約下でのスケジューリング手法について検討し、性能のみを考慮した従来手法と比較しつつ評価を行った。従来研究で用いられていた最適化問題の目的関数にエネルギー効率の項を追加して再定義することで、消費エネルギーも考慮したスケジューリングを行うことが可能になる。シミュレーションにより評価実験を行い、従来手法に比べて消費エネルギーを削減できることがわかった。また、目的関数内の重み変数を変化させることで、スケジューリングにおいて性能またはエネルギー効率の重視度を変化させることが可能であると分かった。今後の課題として、組み合わせ最適化問題の解法を工夫し、より多くの探索候補の中から最適な実行設定を選択すること、また性能電力モデルの妥当性を検証することなどが挙げられる。

**謝辞** 本研究の一部は、科学技術振興機構・戦略的創造研究推進事業 (CREST) の研究プロジェクト「ポストベタスケールシステムのための電力マネジメントフレームワークの開発」の助成により行われたものである。

## 参考文献

- [1] Top 500 Supercomputer Sites: <http://www.top500.org/>.
- [2] O. Sarood, A. Langer, A. Gupta, and L.V. Kale: Maximizing Throughput of Overprovisioned HPC Data Centers Under a Strict Power Budget, *Proc. of the International Conference on High Performance Computing, Networking, Storage and Analysis*, pages 807-818, SC'14, (2014).
- [3] 會田 翔, 三輪 忍, 中村 宏: ロードバランスを考慮した電力制約下における CPU の DVFS 制御, 情報処理学会研究報告ハイパフォーマンスコンピューティング (HPC), Vol.2014-HPC-143, No.23, pp. 1-8, (2014).
- [4] O. Sarood, A. Langer, L. Kale, B. Rountree, and B. de Supinski: Optimizing Power Allocation to CPU and Memory Subsystems in Overprovisioned HPC Systems, *Proc. of 2013 IEEE International Conference on Cluster Computing*, pp. 1-8, (2013).
- [5] S.Imamura, H. Sasaki, N. Fukumoto, K. Inoue, and K. Murakami. Optimizing Power-Performance Trade-off for Parallel Applications through Dynamic Core and Frequency Scaling. In *The 2nd Workshop on Runtime Environments/Systems, Layering, and Virtualized Environments*. RESoLVE'12, March (2012).
- [6] 稲富雄一, 吉田匡兵, 深沢圭一郎, 上田将嗣, 青柳 睦, 井上弘士: 電力指向型次世代スーパーコンピュータを想定した HPC アプリケーションの性能最適化~量子化学計算の場合~, 情報処理学会研究報告ハイパフォーマンスコンピューティング (HPC), Vol.2013-HPC-142, No.30, pp. 1-6, (2013).
- [7] T. Patki, D. K. Lowenthal, B. Rountree, M. Schulz, and B. R. de Supinski. Exploring hardware overprovisioning in powerconstrained, high performance computing. In *International Conference on Supercomputing (ICS)*, (2013).
- [8] T. Komoda, S. Hayashi, T. Nakada, S. Miwa, and H. Nakamura: Power Capping of CPU-GPU Heterogeneous Systems Through Coordinating DVFS and Task Mapping, *Proc. of 2013 IEEE 31st International Conference on Computer Design*, pp. 349-356, (2013).



- [9] M. Etinski, J. Corbalan, J. Labarta, M. Valero: Parallel job scheduling for power constrained HPC systems, *Parallel Computing*, Volume 38, Issue 12, pp. 615-630, (2012).
- [10] F. Kaplan, J. Meng, and Ayse K.Coskun. Optimizing communication and cooling costs in HPC data centers via intelligent job allocation. In *Proceedings of the International Green Computing Conference (IGCC)*, pp. 1-10, (2013).
- [11] J. Meng, S. Mc Caule, F. Kaplan, Vitus J.Leung, and Ayse K.Coskun. Simulation and optimization of HPC job allocation for jointly reducing communication and cooling costs. In *Sustainable Computing: Informatics and Systems*, pp. 1-10, (2014).
- [12] A. B. Downey: A Model for Speedup of Parallel Programs, U.C. Berkeley Technical Report CSD-97-933 , (1997).
- [13] "Intrepid - Blue Gene/P Solutione", <http://www.top500.org/system/176322>.
- [14] "Parallel Workloads Archive",<http://www.cs.huji.ac.il/labs/parallel/workload/>, Tech. Rep.
- [15] Nas parallel benchmarks. <http://www.nas.nasa.gov/publications/npb.html>
- [16] Rotem, E., Naveh, A., Rajwan, D., Ananthkrishnan, A. and Weissmann, E.: Power-Management Architecture of the Intel Microarchitecture Code-Named Sandy Bridge, *IEEE Micro*, Vol.32, No.2, pp. 20-27, (2012).
- [17] David, H., Gorbato, E., Hanebutte, U. R., Khanna, R. and Le, C.: RAPL: Memory Power Estimation and Capping, *Proceedings of the 16th ACM/IEEE International Symposium on Low-Power Electronics and Design (ISLPED)*, pp. 189-194, (2010).
- [18] L. Kale: The Chare Kernel Parallel Programming Language and System, *Proc. of the International Conference on Parallel Processing*, vol. II, pp. 17-25, (1990)
- [19] L. Kale, A. Arya, N. Jain, A. Langer, J. Lifflander, H. Menon, X. Ni, Y. Sun, E. Toton, R. Venkataraman, and L. Wesolowski: Migratable Objects + Active Messages + Adaptive Runtime = Productivity + performance a Submission to 2012 HPC Class II Challenge," *Parallel Programming Laboratory*, Tech. Rep. 12-47, (2012).