

TSUBAME3.0 における計算資源分割の検証

阿部 公一¹⁾, 香月 稔¹⁾, 藤田 和宏²⁾, 鶴見 慶¹⁾, 安良岡 由規²⁾, 根本 忍²⁾,
梁井 善行¹⁾, 野村 哲弘³⁾, 三浦 信一³⁾, 渡邊 寿雄³⁾, 額田 彰³⁾, 遠藤 敏夫³⁾

1) 東京工業大学 研究推進部

2) 東京工業大学 技術部

3) 東京工業大学 学術国際情報センター

computer@o.cc.titech.ac.jp

Resource Segmentation in TSUBAME3.0 Supercomputer

Masakazu Abe¹⁾, Minoru Katsuki¹⁾, Kazuhiro Fujita²⁾, Kei Tsurumi¹⁾, Yoshinori Yasuraoka²⁾,
Shinobu Nemoto²⁾, Yoshiyuki Yanai¹⁾, Akihiro Nomura³⁾, Shin'ichi Miura³⁾,
Toshio Watanabe³⁾, Akira Nukada³⁾, Toshio Endo³⁾

1) Research Promotion Department, Tokyo Institute of Technology

2) Technical Department, Tokyo Institute of Technology

3) Global Scientific Information and Computing Center, Tokyo Institute of Technology

概要

東京工業大学学術国際情報センターにおいて2017年8月より運用を開始したスーパーコンピュータ TSUBAME3.0 では、効率的な運用のために、計算ノードを論理的に分割することにより、より多くのジョブを同時実行するよう設計されている。本稿では、過去1年間の運用実績をもとに、計算ノードの論理分割が実際にどのように使われてきたかを検証する。

1 はじめに

東京工業大学学術国際情報センターは TSUBAME3.0 を2017年8月より稼働を開始した [1]。TSUBAME3.0 は HPE SGI ICE XA を基にカスタマイズされ、540 台の計算ノードには CPU1,080 基、GPU2,160 基が搭載され、理論最大性能は倍精度で12.15ペタフロップス、半精度(以上)で47.2ペタフロップスになる。

各計算ノードには2TBのNVMe対応SSDが搭載され、合計容量1.08PBのスクラッチ領域を備える。また、全計算ノードから参照可能なストレージとして容量15.9PB、データ転送速度150GB/sの高速ストレージが用意されている。

システムの冷却方式はTSUBAME2.0/2.5の間接水冷およびTSUBAME-KFC[2]の液浸・温水冷却等で培った経験から、主要な熱源であるCPUとGPUのみを直接水冷、他のコンポーネントを間接水冷とする等の工夫を加えることにより、冷却塔による水冷効果を高め省電力性に寄与している。

本稿では、TSUBAME3.0がより多くのユーザに効率的に利用されるために行った資源分割方式の

概要と、その利用状況について報告する。



図1: TSUBAME3.0 計算ノードラック

2 TSUBAME3.0 のノード構成

TSUBAME3.0(図1)の計算ノード540台全ての構成は同一である。1台のノードには、CPUとしてIntel Xeon E5-2680 V4 (14コア, 2.4GHz) を2基搭載し、主記憶装置は256GiB、バンド幅154GB/s

を持つ。GPUとしてはNVIDIA Tesla P100を4基搭載する。ネットワークインターコネクはIntel Omni-Pathを4ポート搭載し、ノードあたり片方向400Gbps、双方向800Gbpsのインジェクションバンド幅を持つ。各ノードに搭載された2TBの大容量SSDは、NVMeサポートにより高速なアクセスが可能である。特にネットワークインターコネクとSSDはビッグデータやAI分野での処理の高速化に重要な要素である。

図2にTSUBAME3.0の計算ノードのブロック図を示す。ノード内の各プロセッサ数が偶数で、上下でほぼ対称な構造であることから後述するとおり資源分割の点で有利となっている。

表1はTSUBAME2.5とTSUBAME3.0の計算ノードの比較である。

表1: TSUBAME2.5とTSUBAME3.0の計算ノードの比較 (数値はノード内合計)

| 指標 | T2.5 | T3.0 | 倍率 |
|---------------|-------|--------|-------|
| CPU | | | |
| コア数(HTを除く) | 12 | 28 | 2.33 |
| コア数×周波数(GHz) | 35.16 | 72.8 | 2.07 |
| メモリ容量(GiB) | 54 | 256 | 4.74 |
| メモリバンド幅(GB/s) | 64 | 153.6 | 2.40 |
| GPU | | | |
| CUDAコア数 | 8064 | 14,336 | 1.78 |
| FP64(TFlops) | 3.93 | 21.2 | 5.39 |
| FP32(TFlops) | 11.85 | 42.4 | 3.58 |
| FP16(TFlops) | 11.85 | 84.8 | 7.16 |
| メモリ容量(GiB) | 18 | 64 | 3.56 |
| SSD | | | |
| 容量(GB) | 120 | 2,000 | 16.67 |
| READ(MB/s) | 550 | 2,700 | 4.91 |
| WRITE(MB/s) | 500 | 1,800 | 3.60 |
| ネットワーク | | | |
| 転送速度(Gbps) | 80 | 400 | 5.00 |
| ノード数 | 1408 | 540 | 0.38 |

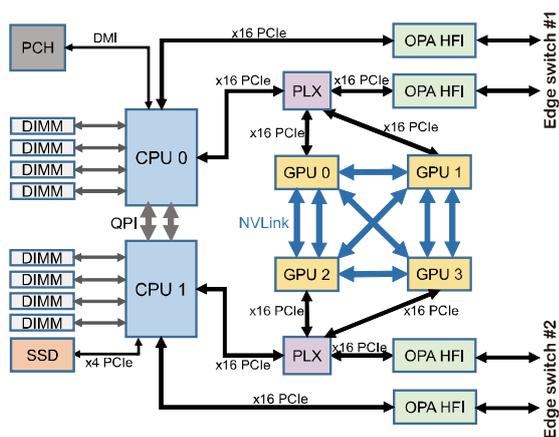


図2: 計算ノードのブロックダイアグラム

このように、TSUBAME3.0の計算ノードはあらゆる面でTSUBAME2.5の計算ノードより強化されているが、ノードの総数はTSUBAME2.5の約38%しかなく、学内外の多数のユーザの計算需要を満たすためにより効率的な運用が求められる。

このため、より先進的な資源分割を実施し、より多くのユーザが効率的に計算ノードを利用できるようにすることが求められている。

3 TSUBAMEにおけるノード動的分割

前節に示したとおり、TSUBAMEシリーズの計算ノードは、多数のコアを持つCPUや複数のGPUなど、ノード単体で見ても強力な計算資源である一方で、TSUBAMEの利用者はスーパーコンピューティングの初心者である学生や、ISVアプリケーションユーザ、GPUコンピューティングの研究者など多岐にわたり、すべてのユーザがTSUBAMEの計算ノードを余すことなく利用できるとは考えにくい。たとえば、GPUを利用しないプログラムの実行中はそのノードのGPUは利用されない遊休資源となる。逆にGPUで計算を行うプログラムは一般的には多数のCPUコアを必要としない。また、1つのプログラムからノード内の複数のGPUを同時に使用するには複数GPUに対応したプログラムを記述する必要があり、プログラミングの難度が上がるため、複数のGPUの利用を望まないユーザも存在する。このような要求資源の性質が異なるユーザ・プログラムを同一の計算ノードに割り当て、計算ノードを共有させることにより限られたTSUBAMEの計算資源をより多くのユーザへ提供可能とする。

TSUBAME2.0/2.5では全体の約1/3のノードにおいてVM技術を用い計算ノードをGPU部分とCPU部分に分割し、それぞれを仮想的に別の計算ノードとして資源提供することで、GPUを中心に利用しCPUはさほど利用しないユーザや、GPUを利用せずCPUのみを利用するユーザの需要に依ってきた。しかしながら、VMおよびスケジューラの資源分割が固定であるため、利用者の需要と合致しない資源が一時的に遊休化するなど、資源利用率の面の課題があった。また、当時のVM技術の制約上、性能を犠牲にすることなくGPUを仮想化することが出来なかったため、ノードのGPU部分は仮想化の対象とせず、GPUジョブはベアメタルノードで実行しつつ、そのノードの中に

CPU ジョブのための VM を作成する形をとる必要があり、ノード内の 3 基の GPU を分割することはできなかった。

TSUBAME3.0 では、TSUBAME2.5 で利用してきた VM 技術に代わり、Linux cgroup を用いた資源分割を導入した。分割された各資源からは GPU および Omni-Path HFI に直接アクセスすることが可能であり、それぞれから利用できるデバイスを制限することもできるため、TSUBAME2.5 より柔軟に資源の分割が可能となる。例えば図 3 に示すように、ノード内の一部の GPU のみを切り出してノードの分割を行うことができる。

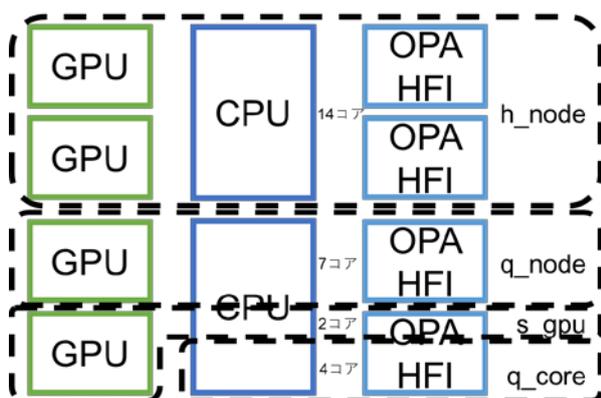


図 3: TSUBAME3.0 計算ノードの資源分割イメージ

ユーザは実行するプログラムの性質に基づき、下記の資源タイプから利用したい資源タイプを選択する。ユーザへの課金(トークンの消費)も、利用する資源タイプにもとづいて行われる。なお、下記 CPU コア数は物理コア数であり、HyperThreading により 2 倍の論理コア数を利用可能である。

- f_node: ノード全体を使用、4GPU, 28 コア
- h_node: 1/2 ノードを使用、2GPU, 14 コア
- q_node: 1/4 ノードを使用、1GPU, 7 コア
- s_gpu: GPU をメインに使用、1GPU, 2 コア
- q_core: CPU コアをメインに使用、4 コア
- s_core: CPU1 コアのみを使用、1 コア

以下、本稿では分割する前の計算ノードのことを物理ノード、分割後の各資源タイプの仮想的なノードのことを論理ノードと呼ぶこととする。

スケジューラである Univa Grid Engine[3]と連携して Linux cgroup 機構を用いた資源分割を動的に行うことにより、TSUBAME2.5 で行われていた静的なパーティショニングは不要であり、あるときは h_node として提供していた部分を次のジョブ

では q_node と s_gpu と q_core に分割して提供するなど、需要に応じた物理ノードの切り出しが行えるようになったために、システム全体の計算資源が有効に活用可能となった。

4 ノード分割における制約事項

TSUBAME3.0 でノードを分割するにあたり、以下のような制約が発生する。これらの制約により実行できなくなるプログラムのユーザは、資源分割を行わない f_node 資源を使用する必要がある。

- 計算ノードへの SSH の利用
デバッグ目的、もしくは X アプリケーションの利用のために、TSUBAME ではジョブ実行中の計算ノードへの SSH を許可しているが、cgroup による資源分割で複数のジョブが実行されている物理ノードへの SSH は、どの論理ノードへの接続かを区別できないため、利用することができない。
対応策として、X アプリケーションの実行については、インタラクティブジョブに限定して、SSH を利用せずに実行する方法を整備し、2018 年 1 月に公開した[4]。これにより X アプリケーションのユーザも資源分割の利点を楽しむことができるようになった。
- ISV アプリケーションの cgroup 対応
ISV アプリケーションによっては、cgroup によって利用できる CPU・GPU 等に制約を受けている状態では動作しないものがある。これらを利用する際には、cgroup による制約がかからない f_node を使用する必要がある。

5 ノードの利用状況

TSUBAME3.0 における 2017 年 9 月から 2018 年 7 月における資源タイプごとのジョブの論理ノード時間積の積み上げを図 4 に示す。各月の積み上げグラフの左側には実際に使用された物理ノード時間積を、右側には当月に提供していた物理ノード時間積が示されている。なお、others とされている計算時間は、TSUBAME グランドチャレンジ制度[5]によるスケジューラの利用を伴わないノード占有利用を示している。

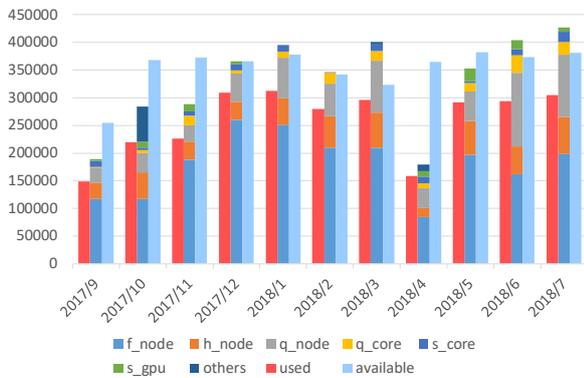


図 4: TSUBAME3.0 の資源タイプ別使用論理ノード時間および使用物理ノード時間

TSUBAME3.0 が提供可能な物理ノード時間は 30 日あたり 388,800 ノード時であるため、たとえ障害などによるノード提供不能時間がなかったとしても、図中に示す全てのジョブを資源分割することなく TSUBAME3.0 上で実行することは不可能である。計算資源を分割することにより、このようなワークロードを論理ノードで実行することが可能となった。

2018 年 7 月時点でも、論理ノード時間積で 47% のジョブは物理ノードを占有して実行しているが、その割合は徐々に減少してきている。これは、課金における優位性だけではなく、X アプリケーションなどの資源分割に伴う制約事項を解消していったことによるものと考えられる。また、ISV アプリケーションのなかに、運用期間中に GPU 対応版がリリースされたもの(例: Gaussian)があり、それらによる需要の変動も考えられる。

なお、2017 年 1 月から 3 月の間、スケジューラの資源分割に関する不具合があり、s_gpu の資源が適切に分離できなかつたため、この期間は s_gpu の利用を禁止しており、利用量が 0 となっている。

5 今後に向けて

TSUBAME3.0 では現在、cgroups によるアクセス制限によりノードを論理的に分割しているが、これをより推し進めて、コンテナ技術によるノードの仮想化を行うことを計画している。コンテナによる仮想化を行うことで、ユーザは自らのプログラムに最適な環境を持ち込むことが可能になることに加え、コンテナで実現されるアプリケーション用ユーザランドと基盤となる物理ノードのユーザランドを分離することにより、それぞれを独立に更新することができるようになり、セキュリティ問題の対処が行いやすくなると期待している。

TSUBAME3.0 導入当初は、スケジューラと連携させて Docker を利用できるように計画していたが、こちらはテスト段階においてスケジューリングに関する不具合が発生したため、ユーザへの提供に至っていない。また、Docker のセキュリティモデルと複数のユーザが計算ノードや並列ファイルシステムを共有するスパコン環境の親和性は高くなく、データセキュリティのため、計算ノード上で実行できる Docker コンテナを何らかの形で制限することが必要となり、コンテナの利点を最大限に活用できる形でユーザに提供できる道筋が立っていない。

一方、コンテナの実行に root 権限を常用しない Singularity がこのような環境におけるコンテナ実現手段として有力となっており、本稿執筆時点でユーザへの提供を開始したところである。コンテナ内から OmniPath の高速な通信を利用する方法など、多少の困難はあるものの、ユーザが任意の環境を持ち込んでプログラムを TSUBAME で実行できるというコンテナ技術の利点を十分に果たすものとして期待している。

6 おわりに

本稿では、東京工業大学のスーパーコンピュータ TSUBAME3.0 において、ノードの論理分割を用いた限りある計算資源の有効活用の取り組みと、その利用状況について報告した。計算資源の仮想化を行うことにより、ユーザはより扱いやすい形で計算ノードにアクセスすることができるとともに、ジョブスケジューリングの観点からも物理的なノード数を超えるジョブを実行することが可能となり、ジョブの実行開始時間の減少も期待される。また、ノードの仮想化を推進することによりユーザがプログラムを実行環境ごとを持ち込めるようにするコンテナ化の展望もあり、実装を推し進めているところである。

謝辞

TSUBAME3.0 の設計および運用には、東京工業大学学術国際情報センターが推進してきた文部科学省「スパコン・クラウド情報基盤におけるウルトラグリーン化技術」および「スマートコミュニティ実現のためのスパコン・クラウド情報基盤のエネルギー最適化の研究推進」、JST CREST(JPMJCR1303, JPMJCR1501) などのプロジ

エクトの研究成果が活用されている。

参考文献

- [1] 藤田 和宏 ほか 「新スーパーコンピュータ TSUBAME3.0 の概要」大学 ICT 推進協議会 2017 年度年次大会 (2017)
- [2] T. Endo, A. Nukada, S. Matsuoka. TSUBAME-KFC: a Modern Liquid Submersion Cooling Prototype towards Exascale Becoming the Greenest Supercomputer in the World. IEEE ICPADS 2014, pp.360-367 (2014)
- [3] Univa Grid Engine,
<http://www.univa.com/products/>
- [4] f_node 以外でも X 転送(GUI)が利用できるよ
うになりました(TSUBAME 計算サービス お
知らせ), <http://www.t3.gsic.titech.ac.jp/node/149>
- [5] 東京工業大学学術国際情報センター：
TSUBAME グランドチャレンジ大規模計算
制度、
<http://www.gsic.titech.ac.jp/GrandChallenge>