



# TSUBAME: Supercomputers for Everybody's AI/Machine Learning

東京工業大学  
遠藤 敏夫

---

2023/4/19 PCCC AI/ML ワークショップ

# 東工大TSUBAMEスパコンシリーズ



TSUBAME1.0~1.2  
(2006~2010)



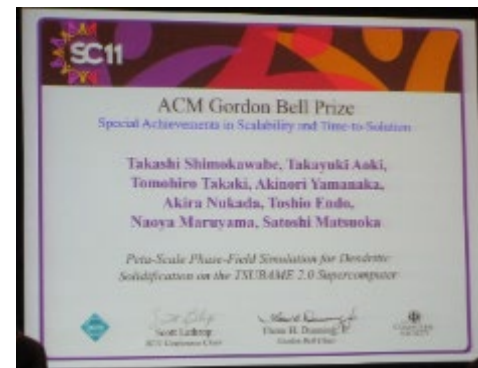
TSUBAME2.0/2.5  
(2010~2017)



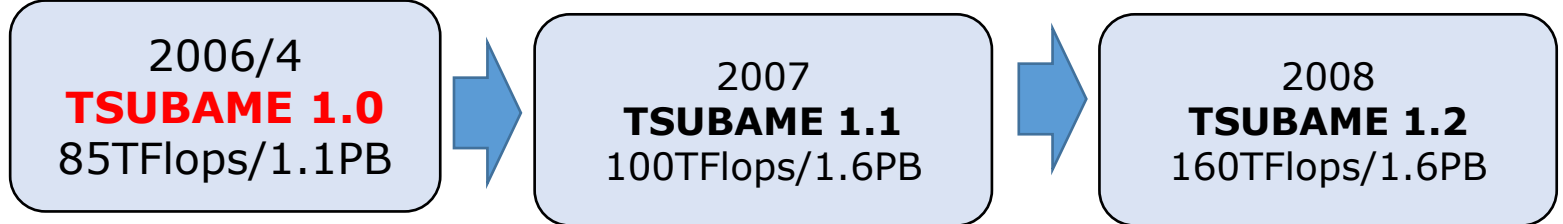
TSUBAME3.0  
(2017~)

TSUBAMEシリーズは、世界に先駆けたGPUの採用など、先進的な取り組みを行ってきた

- アジアNo.1 スパコン認定 (2006)
- 世界初大規模GPUスパコン (2008)
- Top500:演算性能世界4位 (2010)
- ACMゴードンベル賞 (2011)
  - Peta-scale Phase-Field Simulation for Dendritic Solidification on the TSUBAME 2.0 Supercomputer
- 文部科学大臣表彰科学技術賞(開発部門)(2012)
  - 運用世界一グリーンペタスパコンの開発
- Green500:省エネ性能世界一 (2017)



# 東工大TSUBAMEの歴史

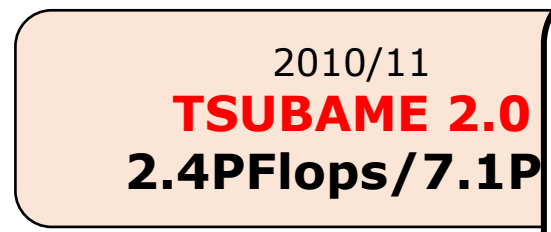
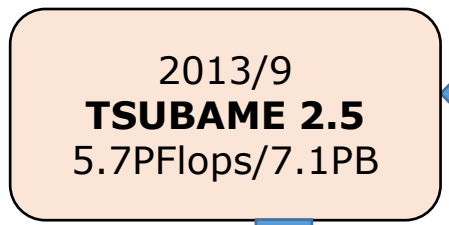


アジアNo.1  
「みんなのスパコン」  
AMD+ClearSpeed

ストレージ・アクセラ  
レータ増強

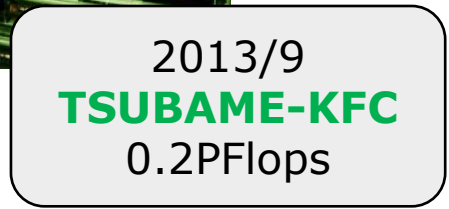
GPUアクセラレータ増強  
→ Top500で初GPUスパコン

GPUをアップグレード

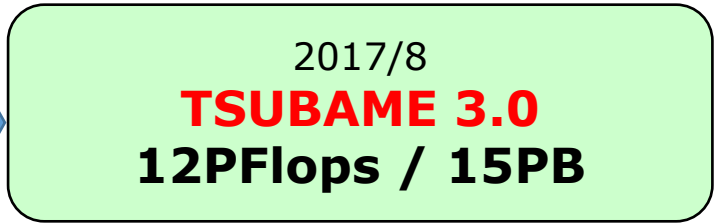


日本初のペタフロップス超

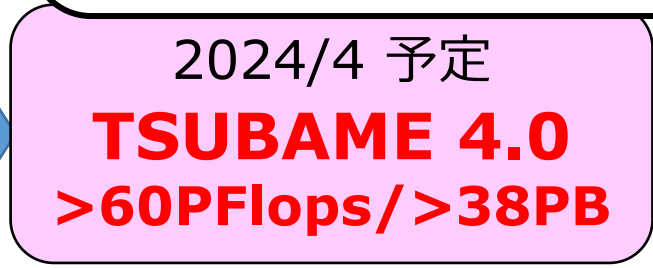
**おことわり：**  
本日時点ではまだ  
**TSUBAME4.0に関する情報**  
は仕様書ベースのものです



液浸冷却による  
小型実験スパコン  
Green世界一



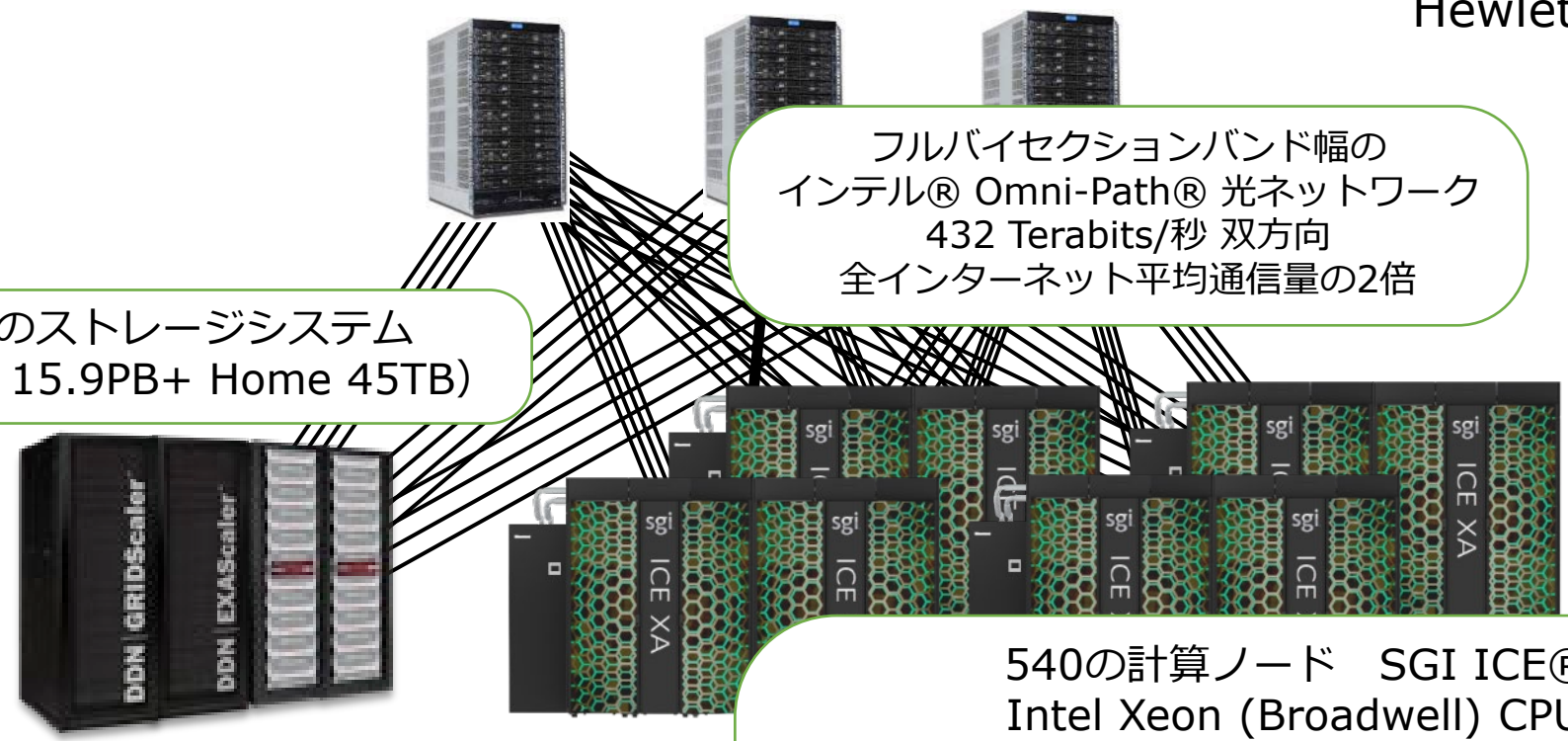
省エネ型「ビッグデータ」スパコン  
Green世界一



# 現行TSUBAME 3.0 のシステム概要

2017/8~

Integrated by  
Hewlett Packard (HPE)



DDNのストレージシステム  
(並列FS 15.9PB+ Home 45TB)

フルバイセクションバンド幅の  
インテル® Omni-Path® 光ネットワーク  
432 Terabits/秒 双方向  
全インターネット平均通信量の2倍

540の計算ノード SGI ICE® XA  
Intel Xeon (Broadwell) CPU×2  
+ NVIDIA P100 GPU×4  
256GBメモリ、2TBのNVMe対応インテル®SSD  
**12PFlops (倍精度), 47PFlops (FP16)**

- ユーザには学内・学外研究者・産業利用を含み、アクティブユーザ数1,400
- 東工大では200近くの研究室が利用
- 深層学習ユーザが大幅増加

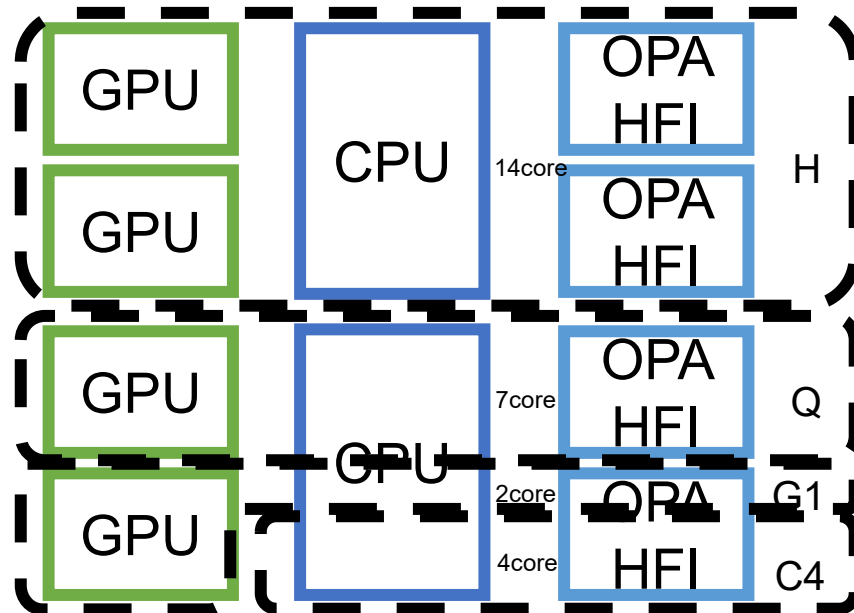
AI/機械学習を強かにサポートするためのTSUBAME3.0の方策：計算科学・シミュレーションとの共存

- ハードウェア面 → **GPU!!**
- ソフトウェア・運用面
  - 機械学習/深層学習フレームワークの提供 (TensorFlow, PyTorch, Chainer...)
  - ↑バージョンアップが早い。ユーザが必要なバージョンを入れられるよう、Singularityコンテナの用意
  - ばらばらな**ジョブの粒度**への対応：>100GPUのMPIジョブもあれば、1コア+1GPUで十分な場合、Pythonスクリプトが動けばよい場合・・・
  - **インタラクティブ利用**需要への対応：バッチスケジューラだけでなく、Webベース利用(2020～)

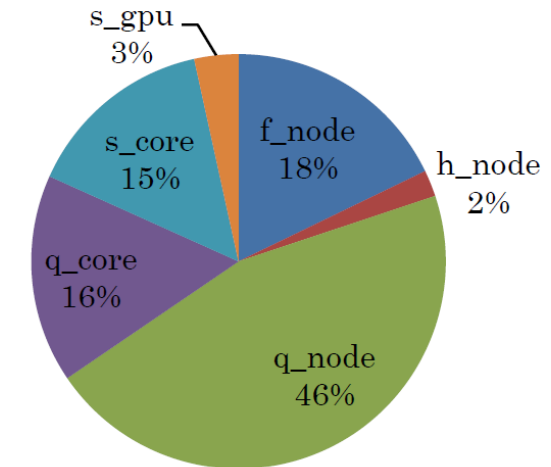
# CPU・GPU資源を有効活用しつつジョブ粒度対応

TSUBAME3.0は「ファットノード」：28 core + 4 GPU + 256GB memory

- 「小さい」ジョブにノード丸ごととはもったいない → **ノード分割**の導入
- TSUBAME3上で、複数の「インスタンスタイプ」を定義した
  - GSIC野村准教授(現在)が設計



- f\_node: 分割しないノード丸ごと
- h\_node: 14 core + 2GPU
- q\_node: 7 core + 1GPU
- s\_core: 1 Core
- q\_core: 4 Core
- s\_gpu: 2 Core + 1 GPU



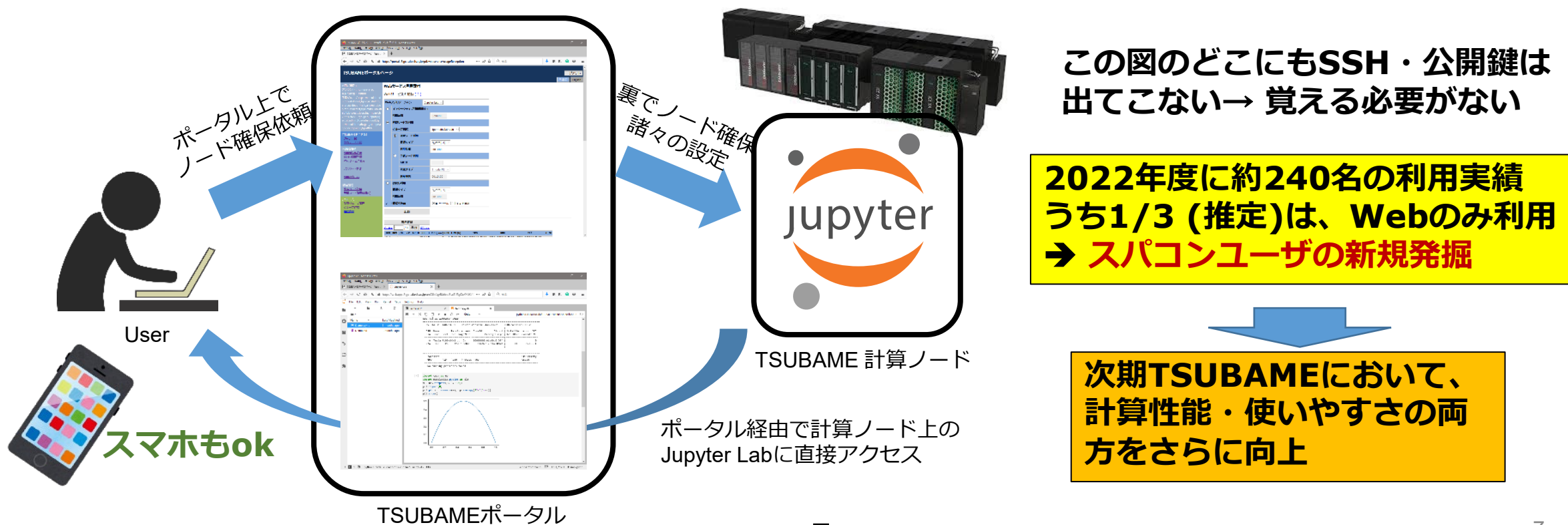
ユーザは、ジョブ投入時に「タイプ」×「インスタンス数」を指定  
→ ユーザは日常的にタイプを選んでいる

ジョブ件数の統計(2023/3):  
>80%が分割タイプ指定

# TSUBAME + シングルサインオン + Jupyter = みんなのビッグデータ・AI/ML基盤 [野村 SWoPP2020]

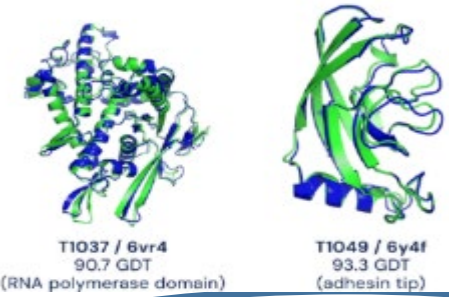
JupyterNoteBook(当時)をはじめ、Web上での計算資源利用が急激に普及

- TSUBAME上にWebインターフェースを実装: JupyterLab / CodeServer / noVNC
  - ブラウザだけでGPUを含めTSUBAMEを直接**インタラクティブ利用**できる
  - 注: OpenOnDemandがメジャーになる少し前だった



# TSUBAME4.0:

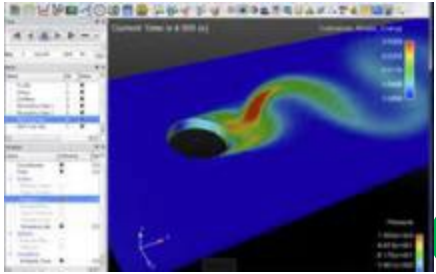
データ・計算科学・AI融合のための「もっと」みんなのスパコン  
により、コンバージェンス・サイエンスの中核インフラへ



深層学習との融合による  
シミュレーション革新



SNSのフォロー関係解析



対話的データ解析

ビッグデータ解析

計算科学・  
シミュレーション

AI・深層学習



深層学習による  
画像等認識

シミュレーションと  
リアルタイム可視化

**TSUBAME4.0**  
**2024/4運用開始予定**  
(写真はイメージ図)



- **現行TSUBAME3と比べ、5倍以上**

**の演算速度** ※倍精度・行列演算にて

- AI・シミュレーションにおいてさらに増大する計算量への対応
- 混雑の緩和へ

- **対話的利用・コンテナ技術の拡充**

- ビッグデータ解析や可視化を容易化、研究のPDCAを加速
- 各ユーザの欲しいソフトウェア環境を迅速に準備
- 待ち時間を短縮するスケジューリングにより、ライトユーザへも恩恵

- **GPUの大幅利用による加速**

- TSUBAMEシリーズではGPUの利用により、演算速度効率が数倍に
- 投資あたりの研究成果の増大

「GPU利用には工夫が必要」という課題に対しては、以下の対応

- 東工大の長年のGPUに関する教育・研究コミュニティの実績
- 深層学習分野ではGPUがデファクトスタンダードになっており、急速に整備

**データ・計算科学・AIを中心とした  
研究成果創出の支援を大幅強化**

下記は仕様書時点の情報

- 総演算性能

- 倍精度:  $\geq 60$ PFlops
  - 半精度:  $\geq 240$ PFlops
- } *TSUBAME3.0の5倍以上*

- 計算ノード(少なくとも一部)は、x86互換CPU + CUDA対応GPU

- 現行機との連続性
- 単体GPUは、少なくとも以下の性能
  - 倍精度  $\geq 36$ TFlops ※A100は要件を満たさない
  - メモリ容量  $\geq 64$ GB

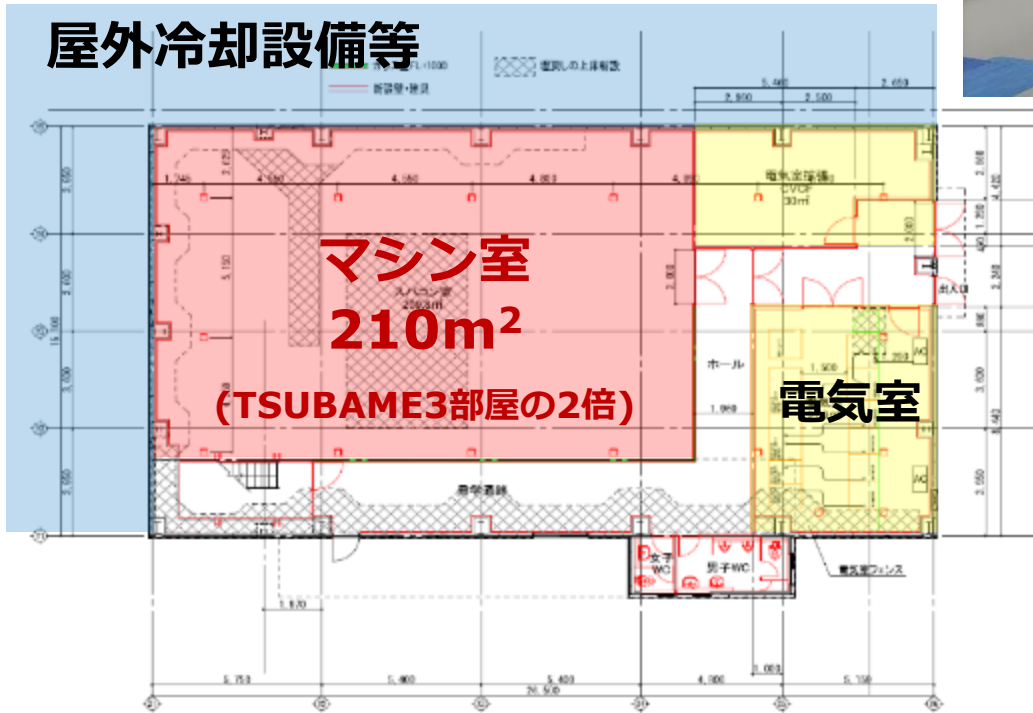
- ストレージ

- 共有ストレージ  $\geq 38$ PB
- 一部、SSDベースの高速共有ストレージ

# 2024/4 TSUBAME4はすずかけ台キャンパスで稼働

TSUBAME1~3は大岡山キャンパスに設置

TSUBAME4: **すずかけ台キャンパス**G4A棟(旧MHD棟)  
を**新データセンター**として設置予定



# TSUBAME4.0のもう少し先へ向けて

HPCシステムへの期待の高まりから、計算資源利用法も広がる

- 機械学習などの対話利用 (LLMも含む)
- リアルタイム高速可視化
- IoTデバイス・エッジからのデータ解析

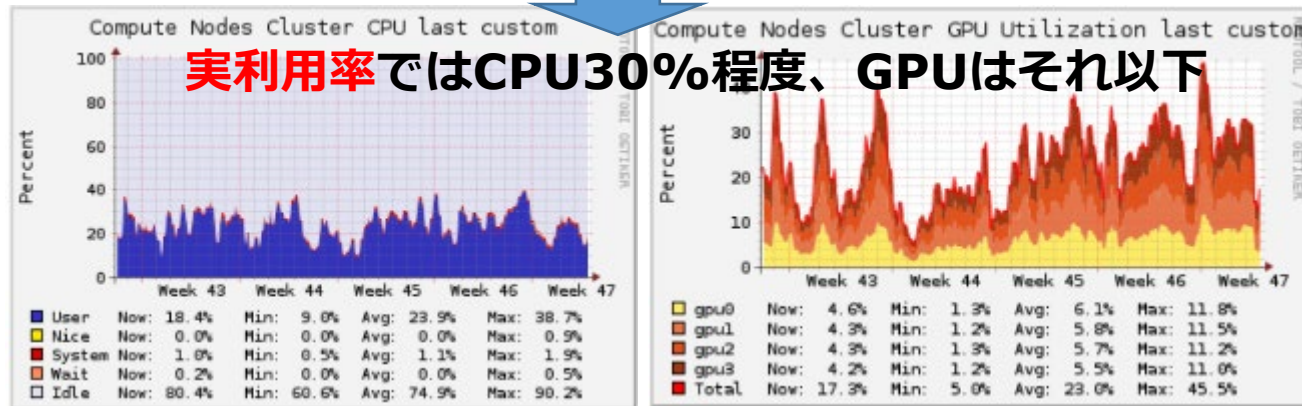
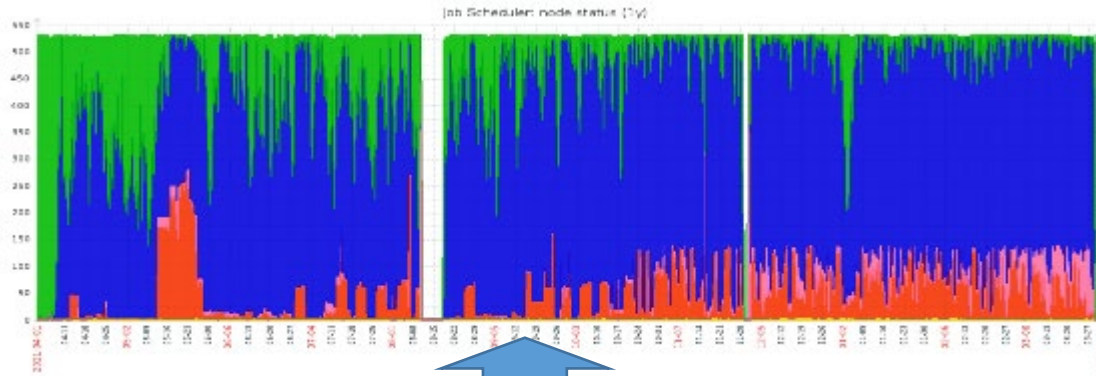


- 現状の多くのスパコンでは、バッチスケジューラによる固定的な運用がほとんど
  - ユーザはジョブスクリプトを前もって作成し、システムに実行依頼（投入）
  - システムがすでに満杯だったら待たされる
  - ➔ **対話利用・リアルタイム可視化×**  
**(2020年から一部TSUBAME3で対応)**
  - 各ジョブは有限時間で終わり、計算資源を明け渡す（時間方向に固定）
  - ➔ **外部データの継続的解析 ×**
- 各ジョブ投入前に、資源量(コア数・GPU数...)指定（空間方向に固定）
- ➔ **計算資源の有効活用 ×**

# スパコンの利用率とは何か

- 各ジョブが利用を宣言した資源と、時間ごとの実利用率にはギャップ

**TSUBAME3の平均占有率85%**  
→ 繁忙期は待ち時間>24時間も



ムーア則が止まる時代に富豪的  
利用はふさわしいか・・・？

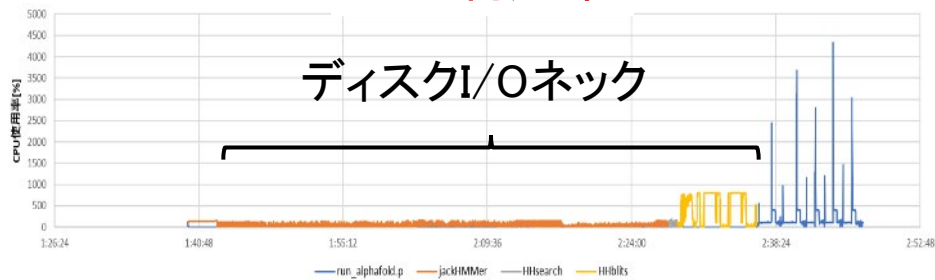
# 利用率のギャップの原因

- AI/機械学習ジョブにおける実利用率の変動

- ビッグデータのディスク読み込み
- 複数の計算フェーズ
  - GPUメインの部分、CPUメインの部分・・・

## AlphaFold2実行中の利用率

### CPU利用率

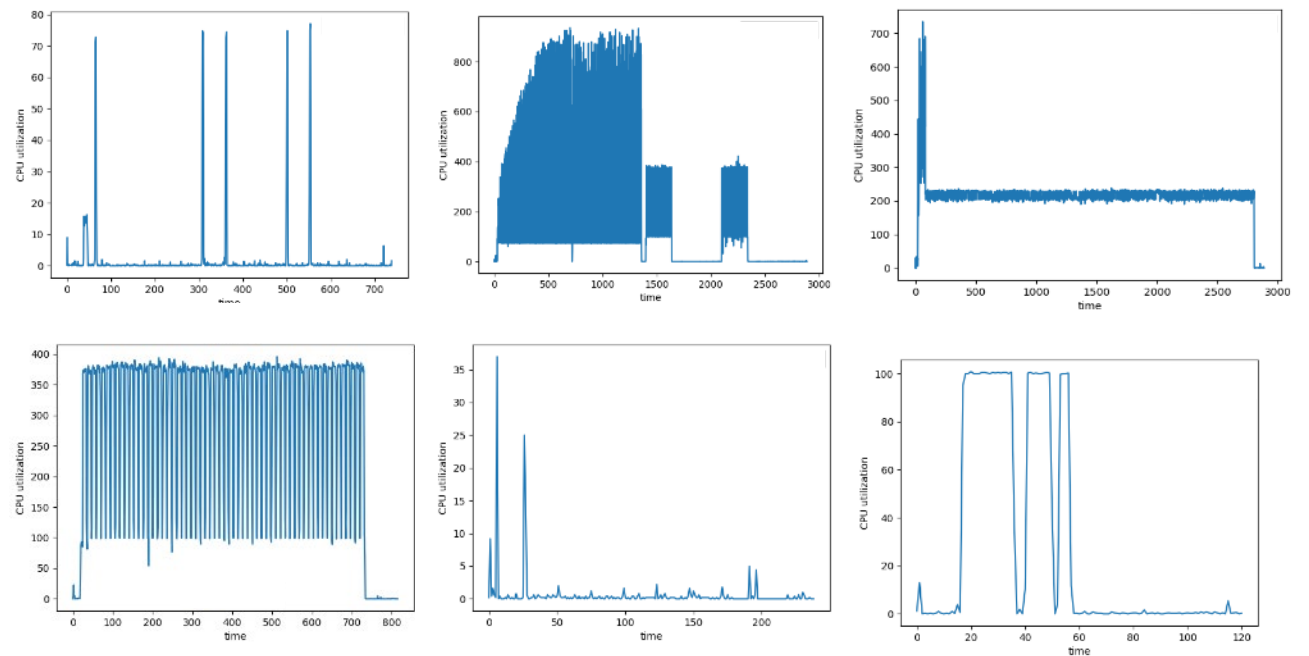


### GPU使用 GPU利用率 [数:107]



- インタラクティブ利用ではさらに、「ユーザが考える」時間が入る

## 複数のインタラクティブ利用中のCPU利用率



# 進行中の研究(一部)

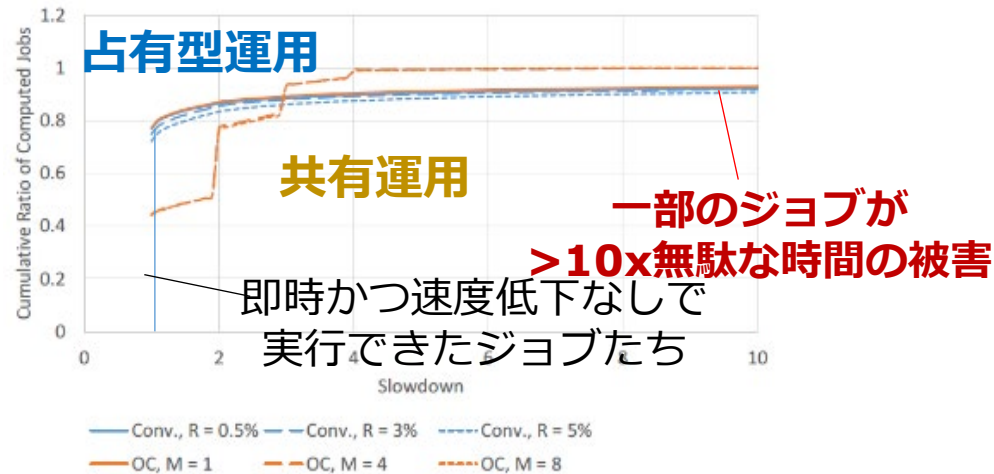
- ノード分割よりもさらにアグレッシブに資源共有を行う試み

## ➔ 「Oversubscribing」スケジューリングの導入？

- コアレベルでも、複数ジョブが計算資源を共有

シミュレーションによる影響評価

Minami et al. HPCAsia23



課題は山積み

- (当然)各ジョブの速度は低下 ➔ 特に並列ジョブへの影響のアセスメント
- スケジューリングアルゴリズム提案
- 公平性・即時性とのトレードオフをどう考えるか
- GPUへの対応！
- 外部データ解析ジョブや、Microserviceとの親和性調査
- メモリ容量の不足をどうカバーするか (Pegasusすごい)

これらを含み、さらなる次世代高性能システムへ向け研究中

TSUBAME4予定地  
このあたり



ありがとうございました

---



Tokyo Tech