

How to Utilize multi-GPU system better for AI model training

AIモデルのトレーニング、マルチGPUシステムの活用方法

Jerry Huang, シニアディレクター

2018/10/17



Faceswap - マルチGPU構成の例





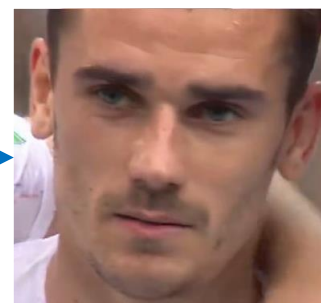
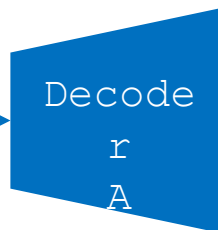
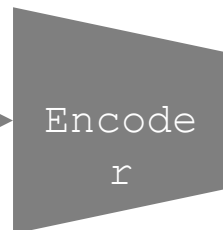
Faceswap – マルチGPU構成の例





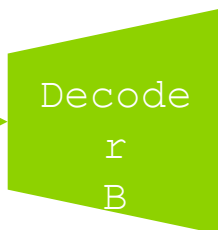
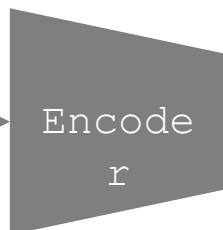
Faceswap – マルチGPU構成の例

元の顔
Antoine



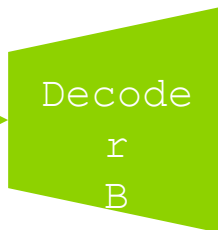
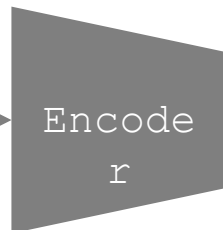
再構成された顔
Antoine

元の顔
Harry



再構成された顔
Harry

元の顔
Antoine



再構成
Antoineの顔から
Harryの顔



AIトレーニングを行う方法

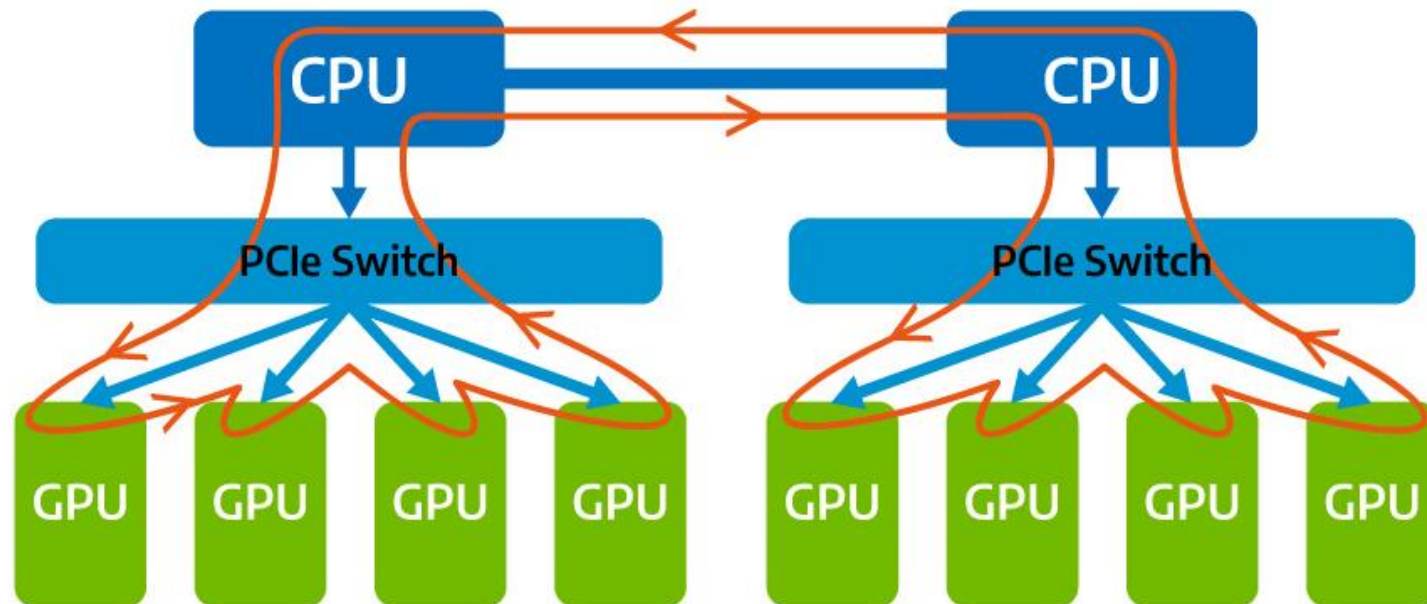


- 1人につき10万枚以上の写真
- 1人につき12～15時間
- トレーニング時間を短縮する方法

複数の**GPU**アクセラレータを展開する

機械学習用のマルチGPU

- データの並列性
 - 同じモデルを各GPUカードに入れますが、データの異なる部分でフィードします。





Wiwynn GPU ソリューション

必要ですか？

・21インチ製品



・19インチ製品



・更なるGPU



SV7400 シリーズ

4U8G GPUサーバー

SV500 シリーズ

4U8G GPUサーバー

XC200 シリーズ

4U16X GPUアクセラレータ



アプリケーションのワークロード

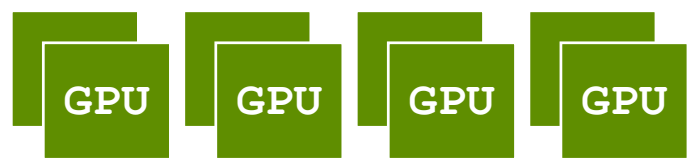
- ワークロードは計算能力 **GPU数** に依存する

ワークロード **1:4** ワークロード **1:8** ワークロード **1:16**



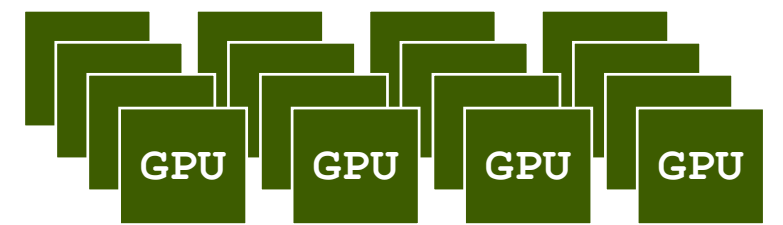
GPUサーバー

サーバー: GPU = 1 : 4



GPUサーバー or
アクセラレータ

サーバー: GPU = 1 : 8



GPUアクセラレータ

サーバー: GPU = 1 : 16



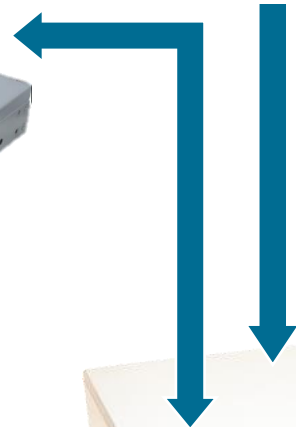


システム構成



SV300G3 (1U Multi-purpose)	
CPU	Intel Xeon Gold 6130 *2
メモリ	DDR4 2400MHz 32GB *16
SSD	Intel® SSD DC S3500シリーズ480GB *1

Mini SAS HDケーブル (PCIe x16)



XC200 (4U16X GPU)	
PCIe スイッチ	Broadcom PEX9797
GPU	NVIDIA V100 *16





XC200シリーズ

19インチ **4U16X GPU** アクセラレータ

Disaggregated PCIe アクセラレータ

16 PCIe 3.0 x 16 アドインカード

最大4台のサーバー接続用の柔軟な構成

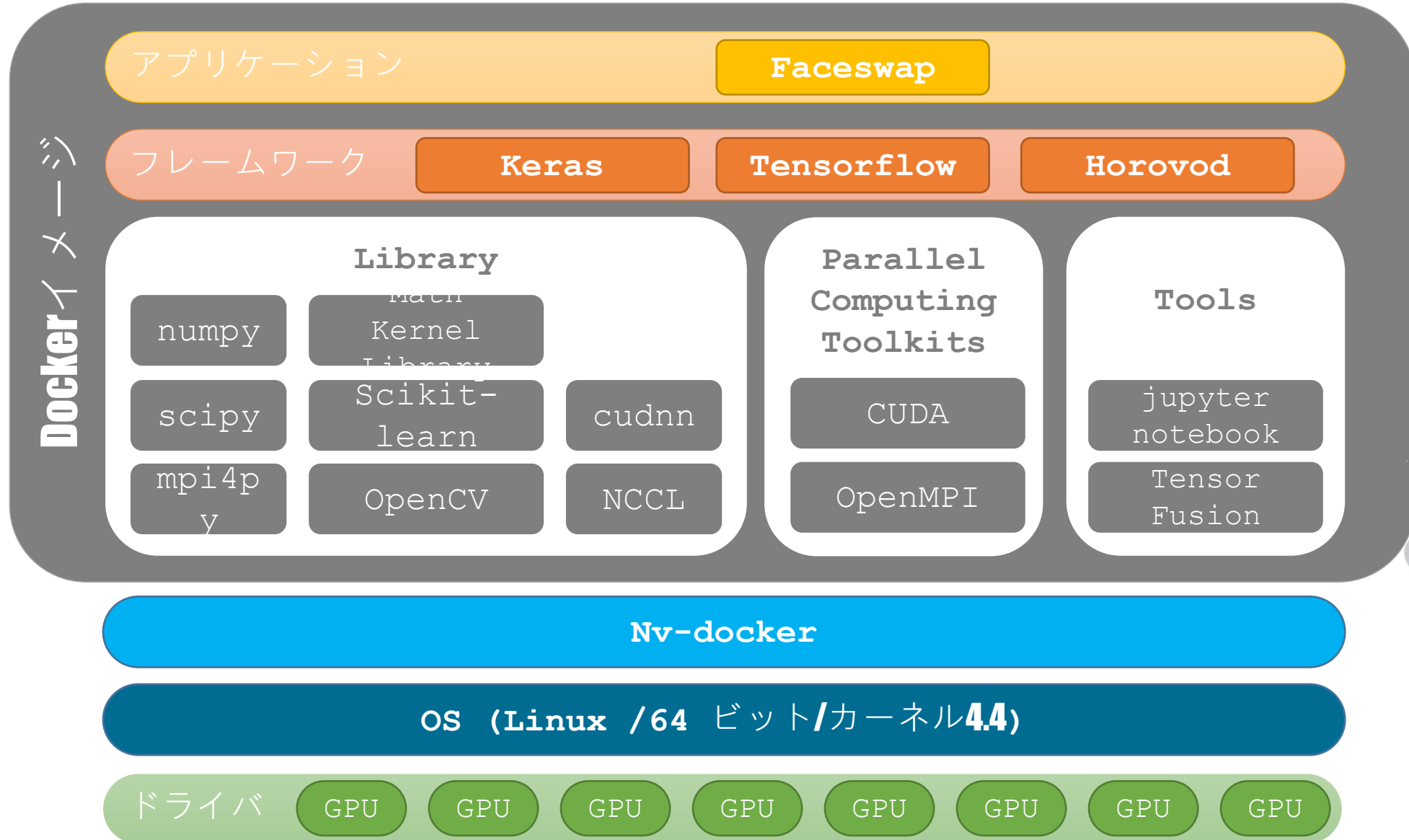
工具不要でメンテナンスが簡単な4つのドロア

ワークロード1:4、1:8、1:16をサポート

4U; 176 (H) x 448 (W) x 900 (D) mm



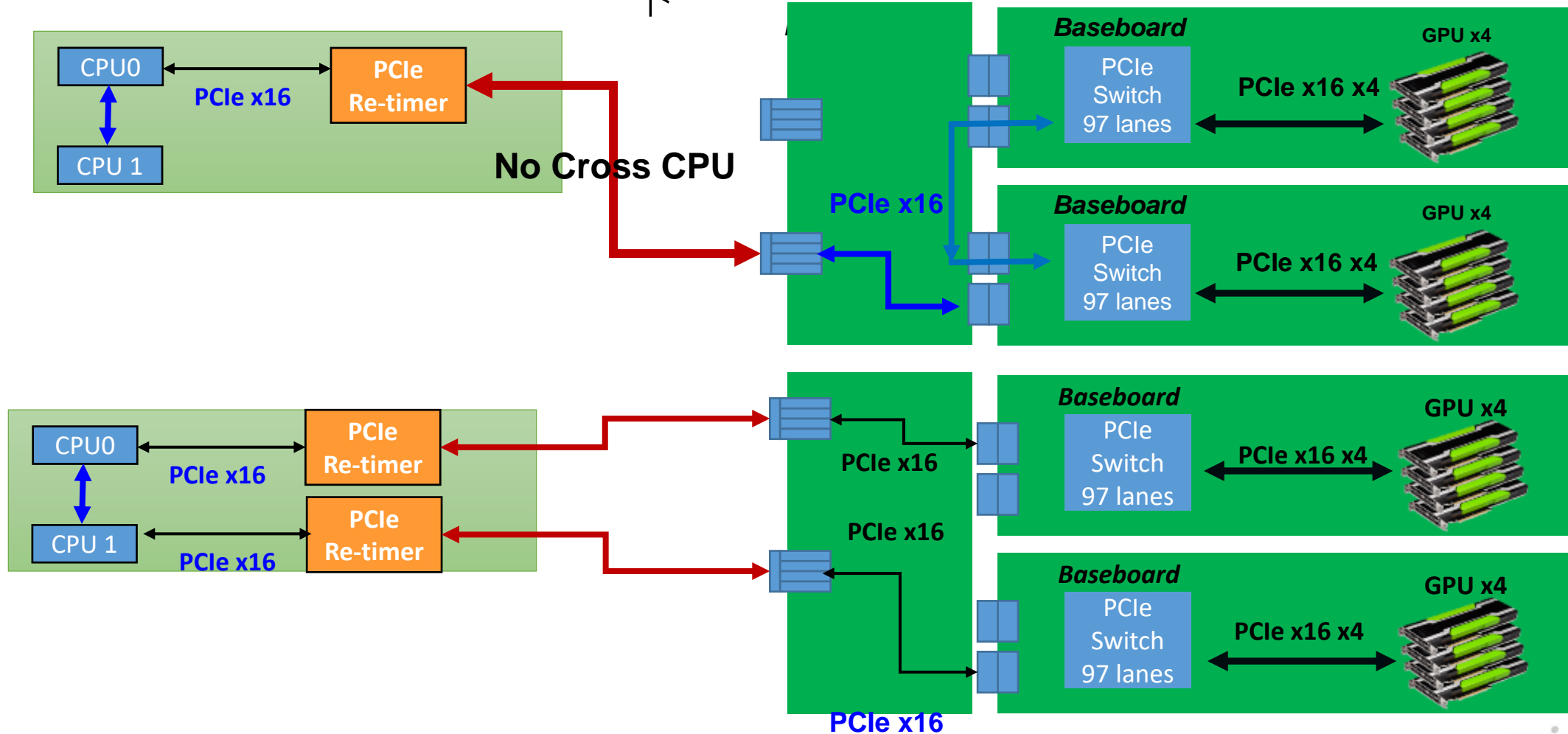
ソフトウェアスタック





アプリケーションのワークロード - 1:8

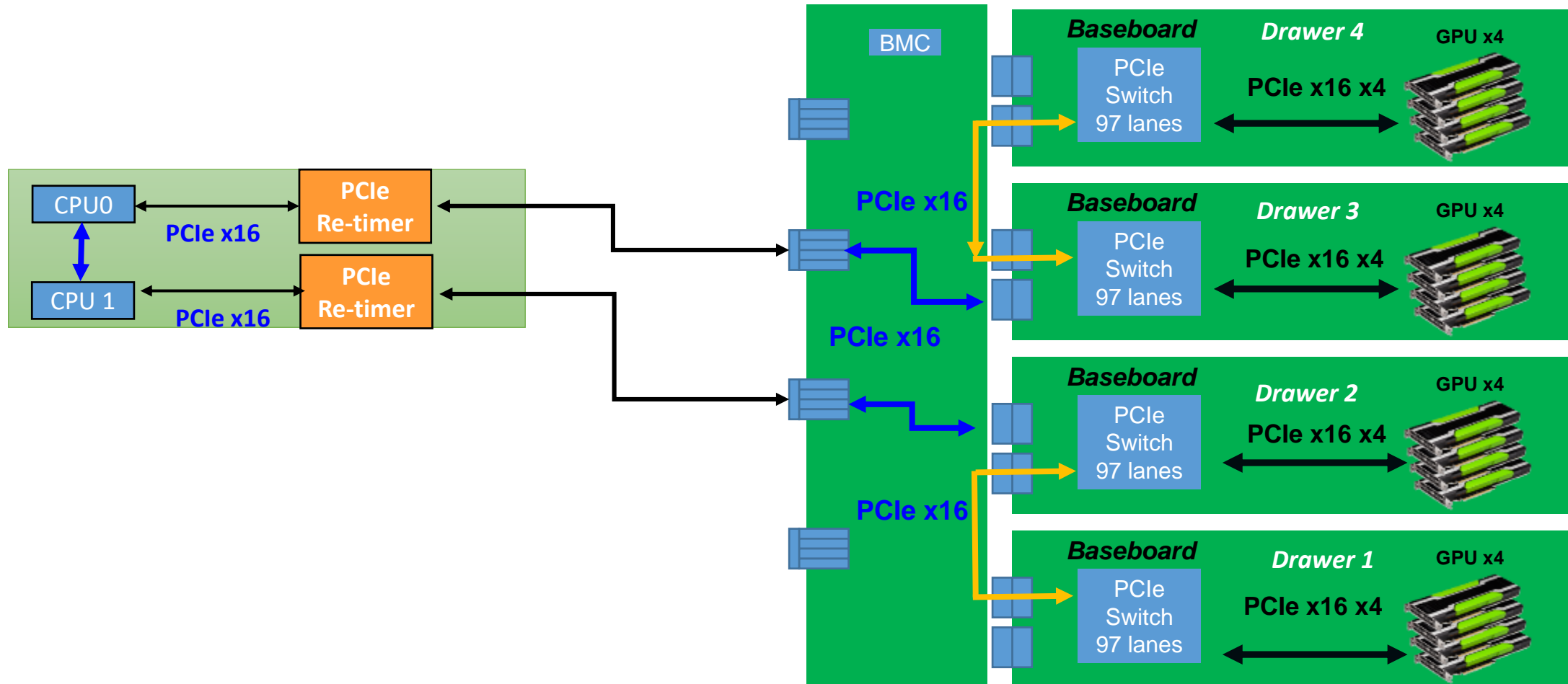
各ノードに専用の8つのGPUアクセラレータカード





アプリケーションのワークロード - 1:16

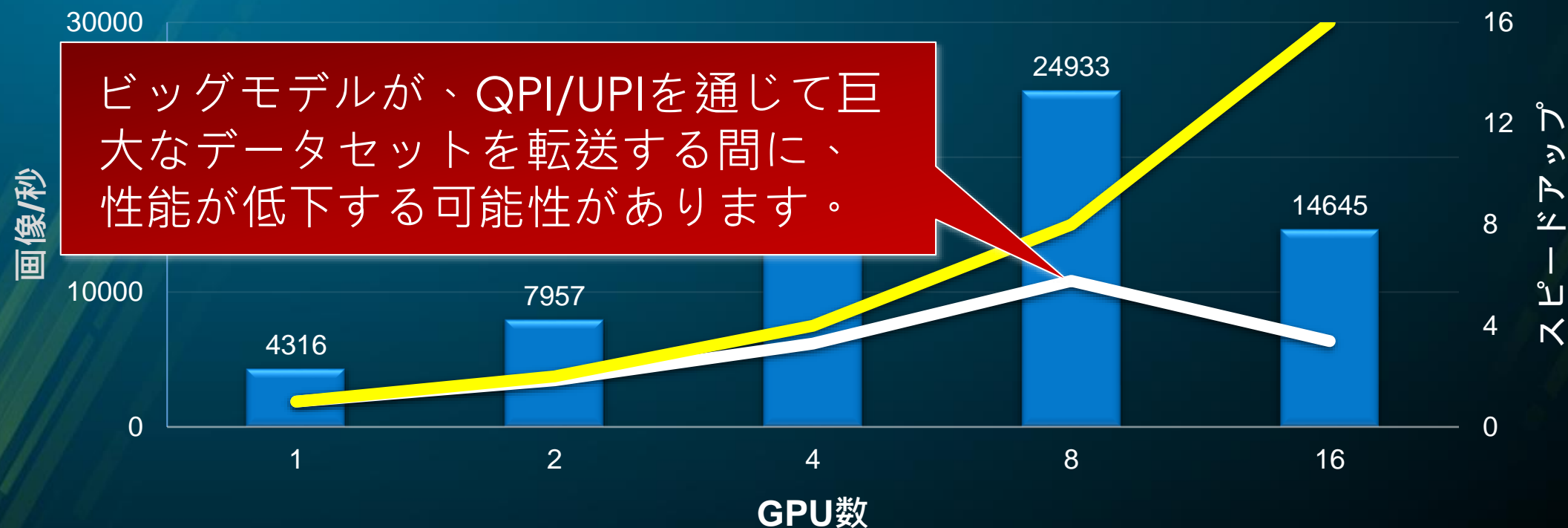
各ノードに専用の16つのGPUアクセラレータカード



ベンチマーク

2倍のGPUを使用している間、コンピューティング能力は、約**1.7~1.8倍**になる可能性があります。

AlexNet (GPUダイレクト対応)



SV300G3+Dr. Know with V100

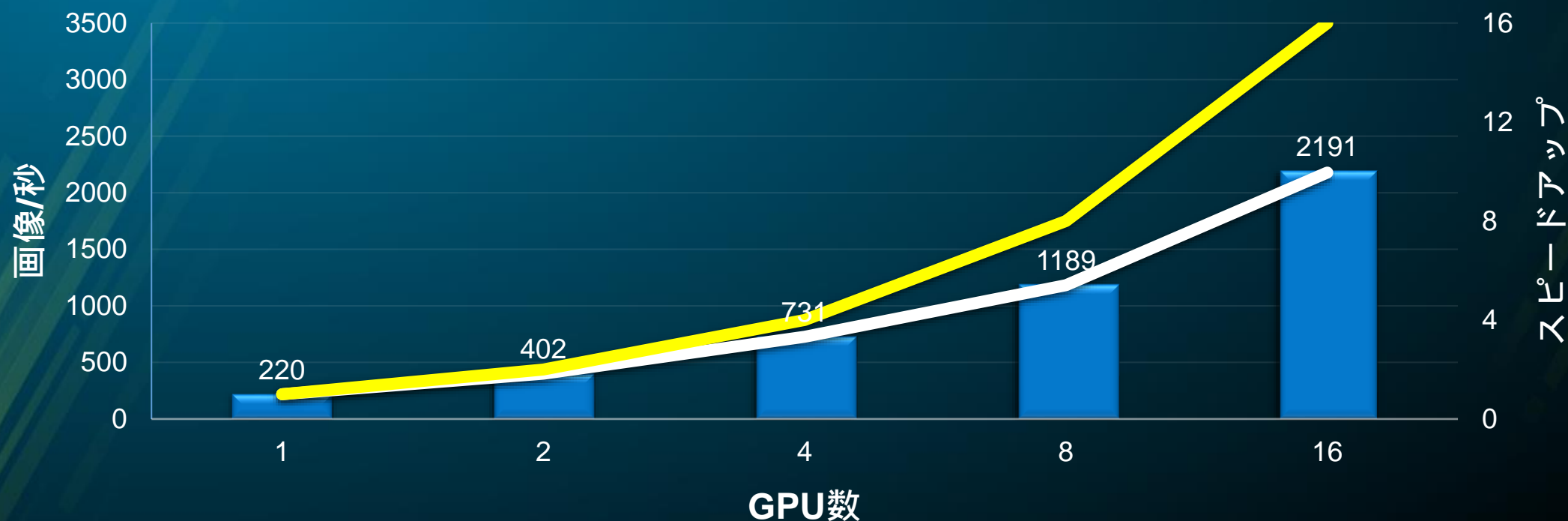
Speed up

Ideal Speed Up

ベンチマーク

CPUをまたぐ、良好な性能を実現したい場合は、より小さいモデルやパラメータサーバーモードを使用することをお勧めします。

Inception V3 (GPUダイレクト対応)



SV300G3+Dr. Know with V100 Speed up Ideal Speed Up



White Paper



Accelerating Your AI/Deep Learning Model Training with Multiple GPUs

White Paper



データセンターのワークロードに最適化されたITソリューション
と最高のTCOを提供します。