# OCP
## SUMMIT

March 20-21
2018
San Jose, CA

OPEN
Compute Project

# Facebook Flexible GPU Expander Big Basin Refresh

Whitney Zhao/HW Eng/Facebook Inc.

Xiaodong Wang/SW Eng/Facebook Inc.

OPEN
Compute Project

OCP
SUMMIT

# Agenda

Introduction

Architecture

Performance

Questions

# Agenda

Introduction ◆

Architecture ◇

Performance ◇

Questions ◇

# Impact

Facebook's commitment to developing AI & advancing ML



LANGUAGE TRANSLATION

FACE RECOGNITION

SEARCH

ADS

NEWS FEED

SIGMA

LUMOS

# Goal

- Open, full contribution to OCP
- Disaggregation/Modularity
- Serviceability



2016: Big Sur



2017: Leopard + Big Basin
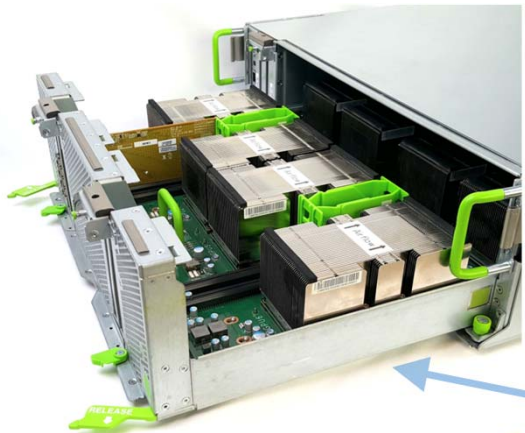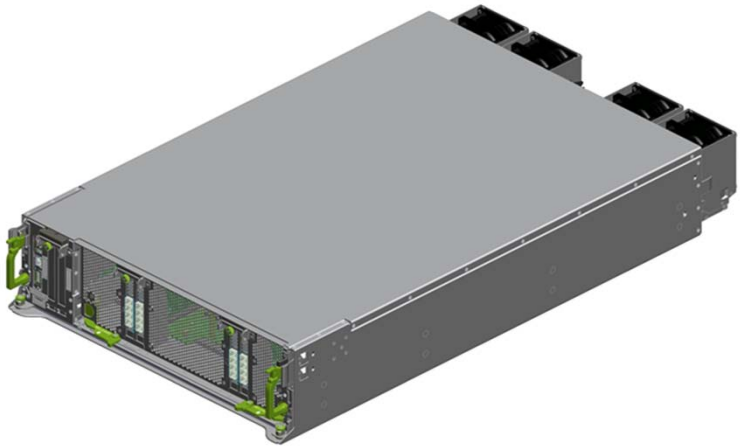2018: Tioga Pass + Big Basin V2
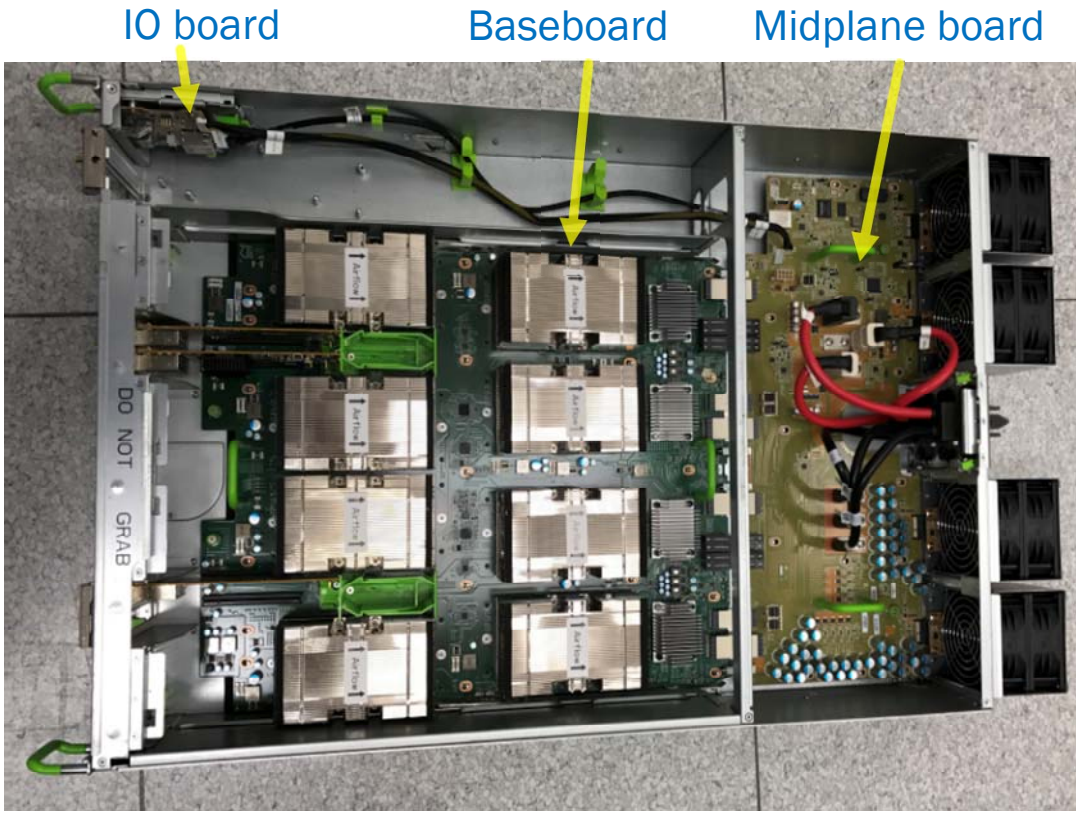
# Big Basin V2 Overview

- 3 OU chassis
- Open Rack v2 compatible
- 8x Nvidia Tesla V100 GPUs; NVLink capable
- 300W TDP for each Tesla V100 GPU
- Facebook 2S Server Tioga Pass as Head node

# A deeper look into Big Basin
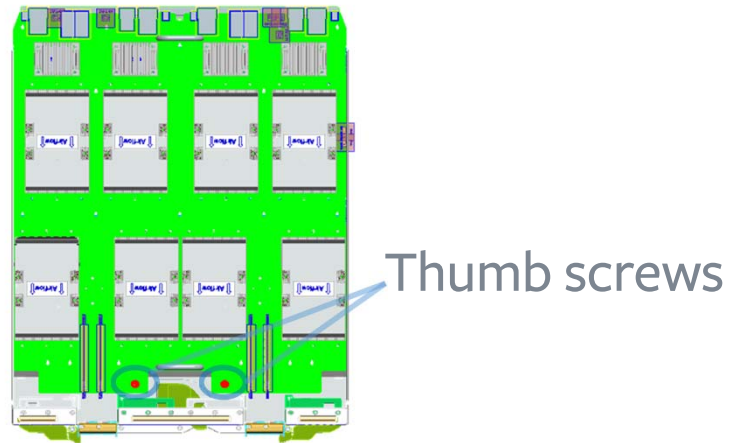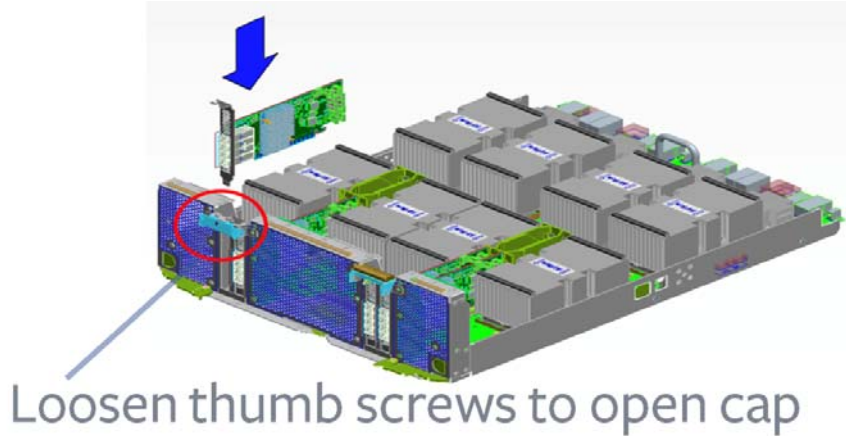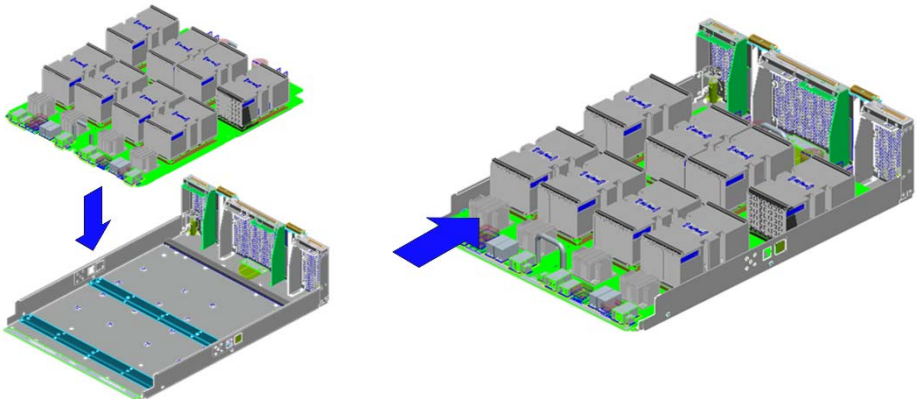


IO board

Baseboard

Midplane board

Baseboard on sliding tray

# Serviceability

- Quick repairs at data center

- Telemetries accessible from head node

- Provisioning Big Basin with its head node is not much different from provisioning existing servers; these servers come with additional GPUs.
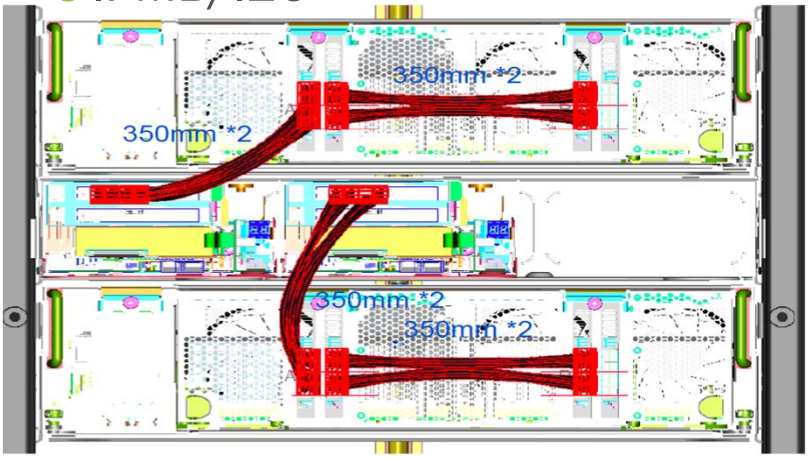


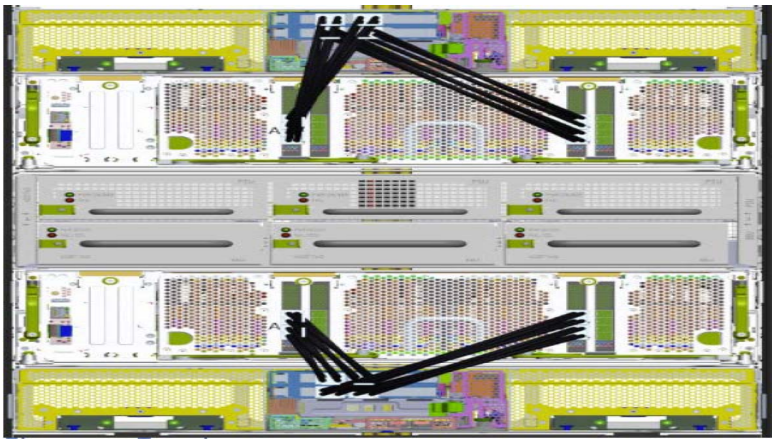Loosen thumb screws to open cap



Thumb screws

# Agenda

Introduction

**Architecture**

Performance

Questions

# Architecture (Headnode to Big Basin)

- MiniSAS HD cable(2 for each x16)
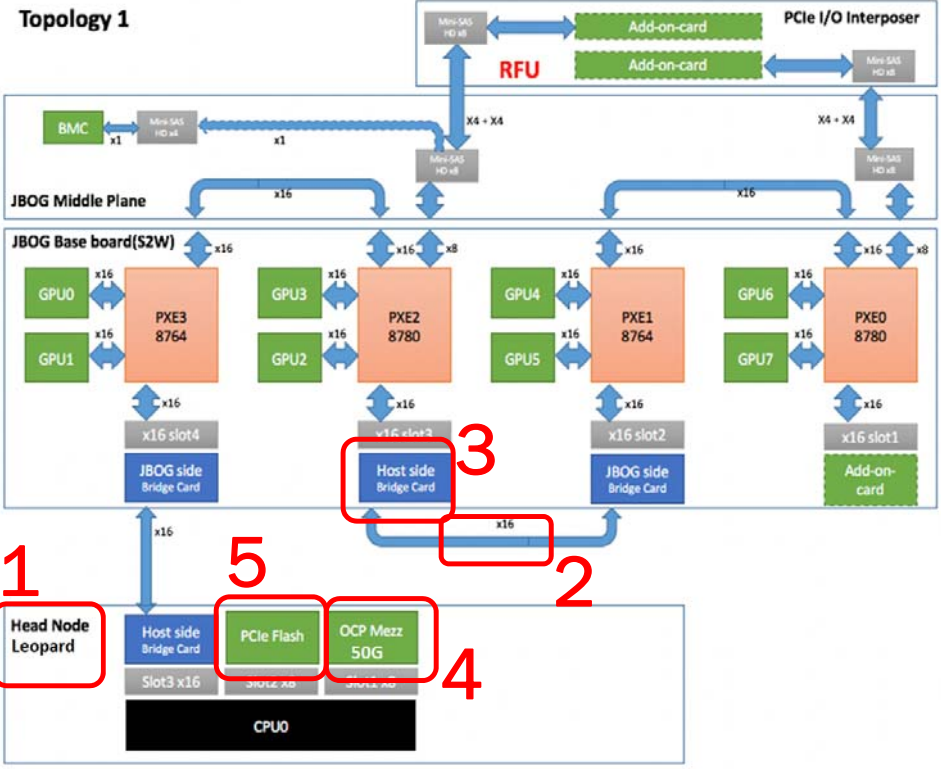  - Standard PCIe x16
  - Present Pin
  - USB2.0
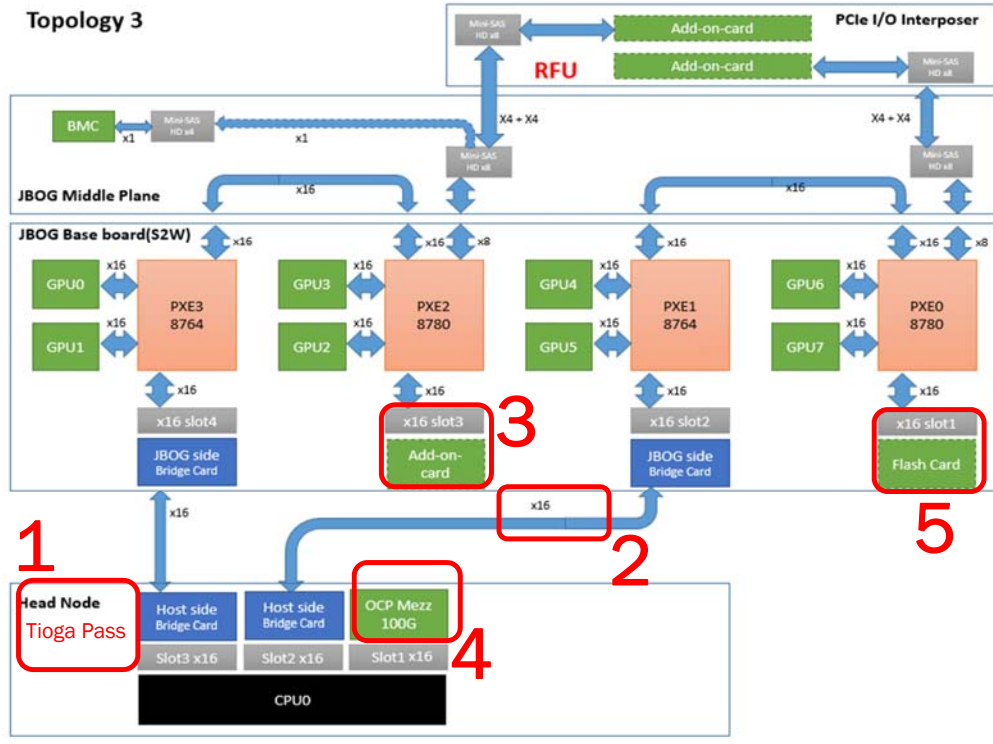  - IPMB/I2C



Leopard + Big Basin(Tesla P100)

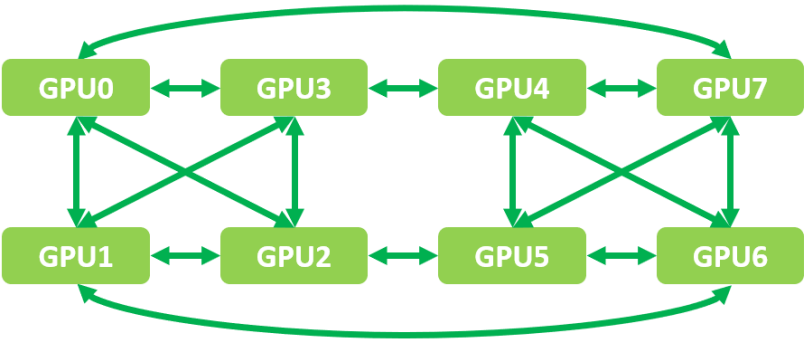Tioga Pass + Big Basin V2(Tesla V100)

# Architecture (PCIe)

# Architecture (NVLINK)



Big Basin W/Nvidia Tesla P100
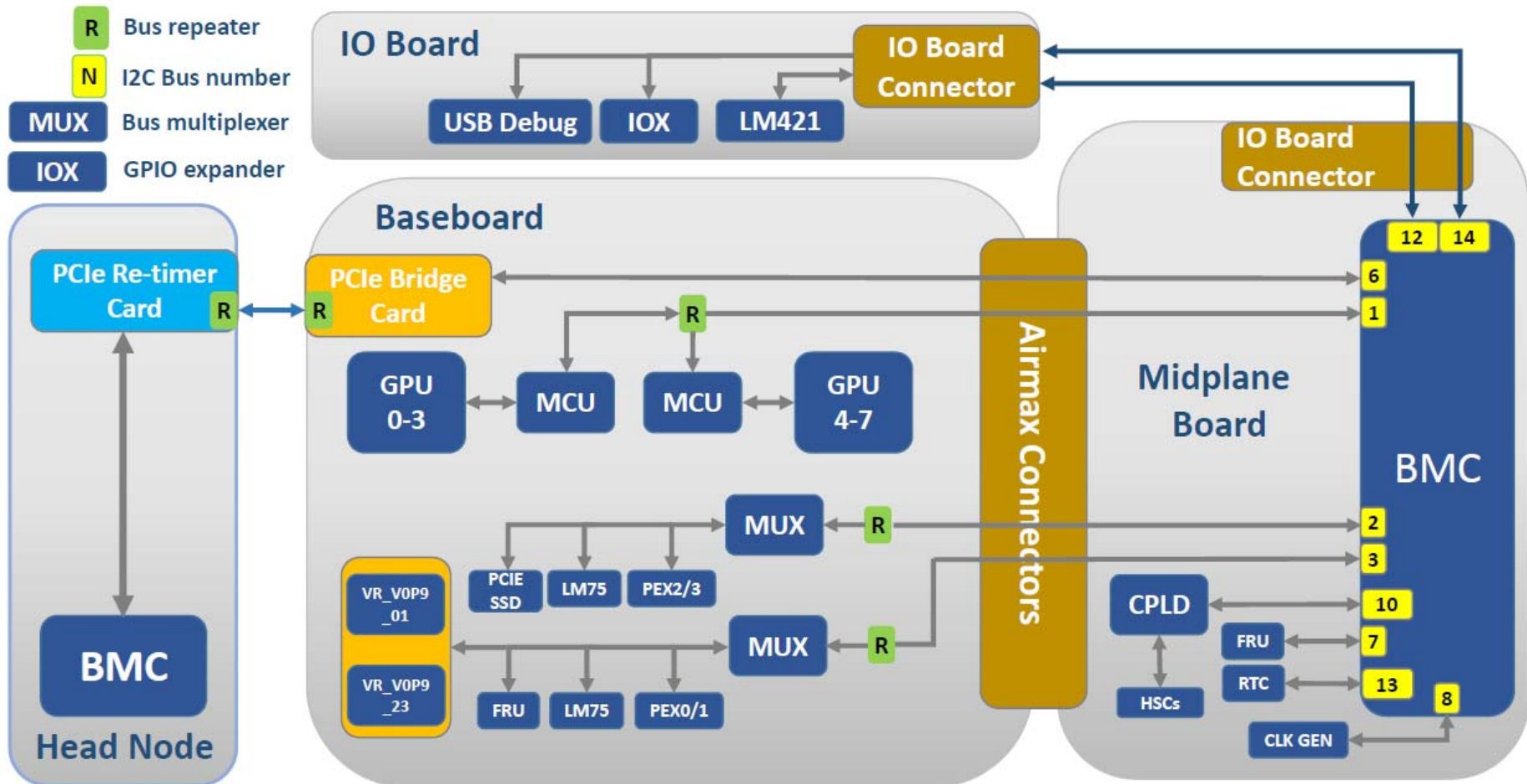
Big Basin V2 W/ Nvidia Tesla V100

# Architecture (IPMB/I2C/PMBUS)

# Agenda

Introduction ◇

Architecture ◇

**Performance** ◆

Questions ◇

# Performance

- Hardware Spec Improvement

- Application performance
  - Computer vision
    - Single-GPU
    - Multi-GPU scalability
    - TensorCore
  - Neural machine translation

# Performance

- Comparisons of GPU Hardware

| | Metrics | NVIDIA V100 | NVIIDA P100 | Improvement |
|---|---|---|---|---|
| **Performance** | FP-32 | 15 TFLOPS | 10.6 TFLOPS | 1.42x |
| | FP-16 | 30 TFLOPS | 21.2 TFLOPS | |
| | TensorCore | 125 TFLOPS | NA | Up to 5x |
| | Mem Bandwidth | 900 GB/s | 720 GB/s | 1.25x |
| | NVLink | 300 GB/s | 160 GB/s | 1.88x |
| **Power** | | 300 W | 300 W | |

# Performance

- Comparisons of GPU Hardware

- Head-node upgrade: Tioga Pass
  - New CPU architecture: Broadwell to Skylake
  - Double PCIe bandwidth
  - Upgraded 100G NIC

- CUDA 9 + cudnn 7: faster libraries, etc.

# Impact - Computer Vision

# Performance metrics in Computer Vision

- Computer Vision: resnet-50
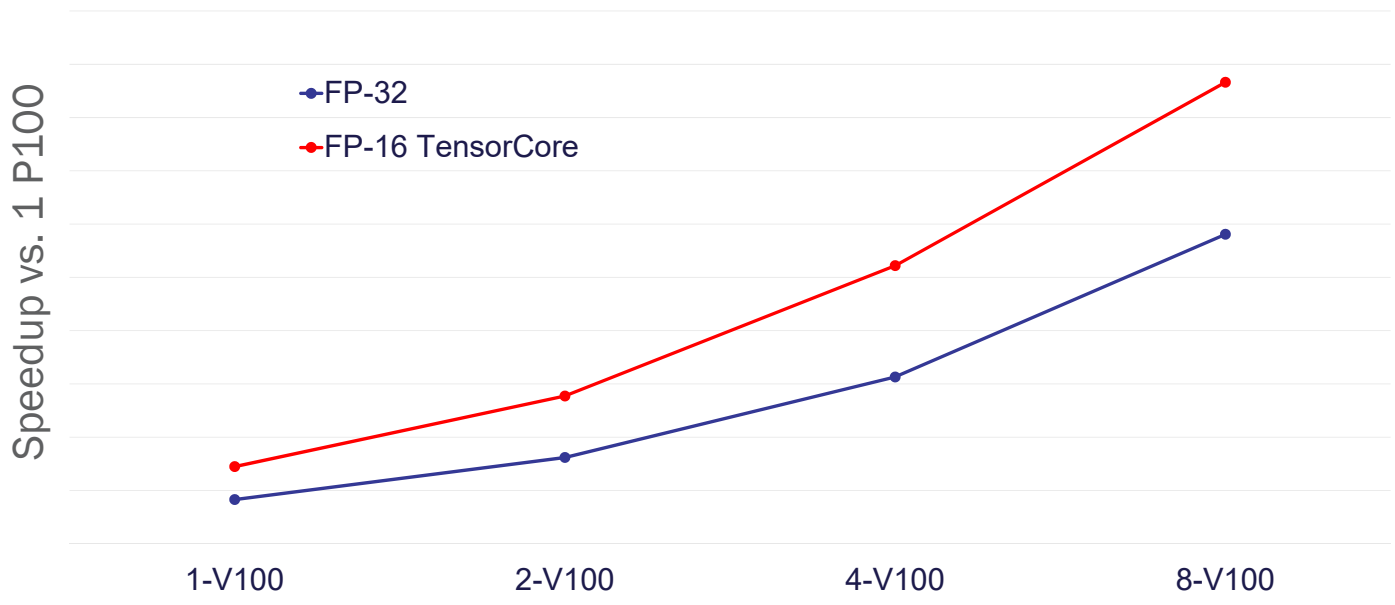  - 1-GPU training speed: use P100 + CUDA 8 as baseline

# Computer Vision Performance

- Computer Vision
  - Multi-GPU speedup vs. 1 P100

# Computer Vision Performance

- Computer Vision
  - High-bandwidth FP-16 TensorCore (WIP)

# Machine Translation

## Better Translation Quality



Phrase-based statistical approach



Neural network approach

# Machine Translation Performance

- Neural Machine translation

| | | |
|---|---|---|
| P100 + CUDA 8<br>Training Throughput as Baseline | V100 + CUDA 8<br>**1.45X** | V100 + CUDA 9<br>**2.2X** |

# Questions?

# OCP Marketplace

- http://www.opencompute.org/products/specsanddesign?keyword=Big+basin