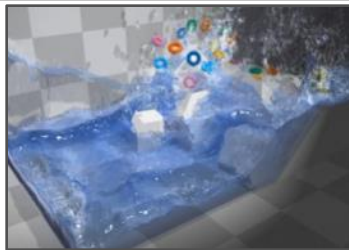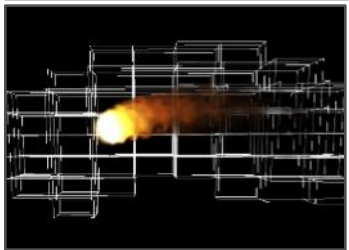# ULTRA-FAST DATA ACQUISITION: GPUS
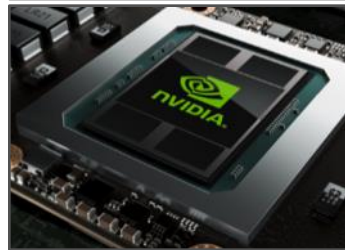
Peter Messmer, Sr. Manager HPC Vis/DevTech
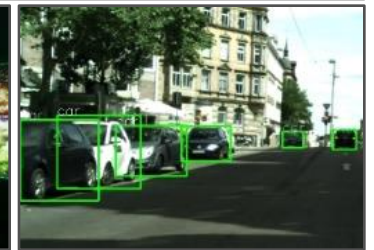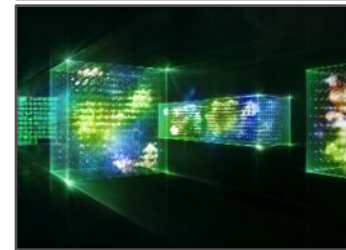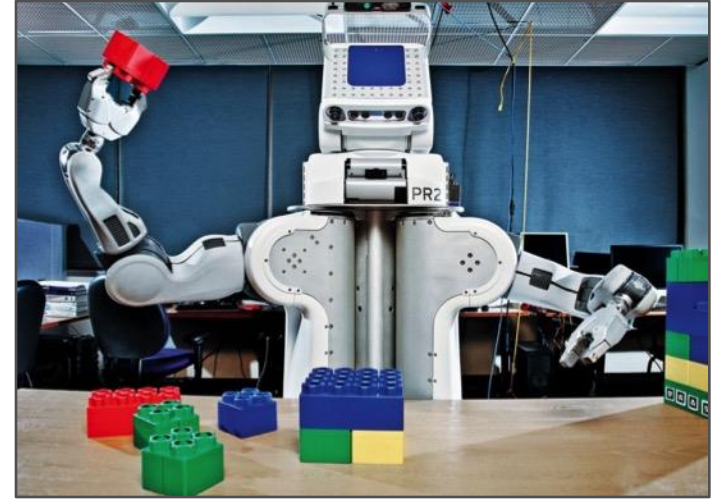
pmessmer@nvidia.com

# NVIDIA - AI COMPUTING COMPANY



Computer Graphics



GPU Computing



Artificial Intelligence

# TESLA PLATFORM
## Leading Data Center Platform for Accelerating HPC and AI

**APPLICATIONS**

| | | |
|---|---|---|
| Alibaba amazon Bai雷百度 ebay facebook flickr Google Microsoft Pinterest skype twitter yelp | Healthcare Manufacturing Finance  Automotive Retail Defense ... | ANSYS 3S SIMULIA GROMACS VASP +450 Applications |
| **INTERNET SERVICES** | **ENTERPRISE APPLICATIONS** | **HPC** |

**INDUSTRY FRAMEWORKS & TOOLS**

| | |
|---|---|
| Caffe2  Microsoft Cognitive Toolkit  mxnet  PYTORCH  TensorFlow  theano | allinea ArrayFire CONTINUUM ANALYTICS WOLFRAM ROGUE WAVE MathWorks |
| **FRAMEWORKS** | **ECOSYSTEM TOOLS** |

**NVIDIA SDK**

| | |
|---|---|
| cuDNN TensorRT NCCL cuBLAS cuSPARSE DeepStream SDK | CUDA C/C++ PGI OpenACC Directives for Accelerators |
| **DEEP LEARNING SDK** | **COMPUTEWORKS** |

**TESLA GPU & SYSTEMS**

| | | | | |
|---|---|---|---|---|
| TESLA GPU | NVIDIA DGX-1 | NVIDIA HGX-1 | DELL Hewlett Packard Enterprise IBM SYSTEM OEM | amazon web services Google Cloud Platform Microsoft Azure CLOUD |

Long-term support. E.g. for CT Scanners; ORNL Titan installed in 2012, still running

## Main functional features of the DAQ system

**Mandatory**
- Store the experimental data (raw or corrected)
- Provide access for off-line data analysis

**Nearly mandatory**
(feedback to users)
- Data visualisation
- Basic on-line processing:
  - Image/frame building
  - Support detector specific corrections (calibration)
  - Basic data extraction (ROI, alarms,..)

**Highly advisable**
- Very low latency partial data processing
- Advanced quasi-on-line data reduction and processing

**Slide courtesy of**
**Pablo Fajardo** *(fajardo@esrf.fr)*
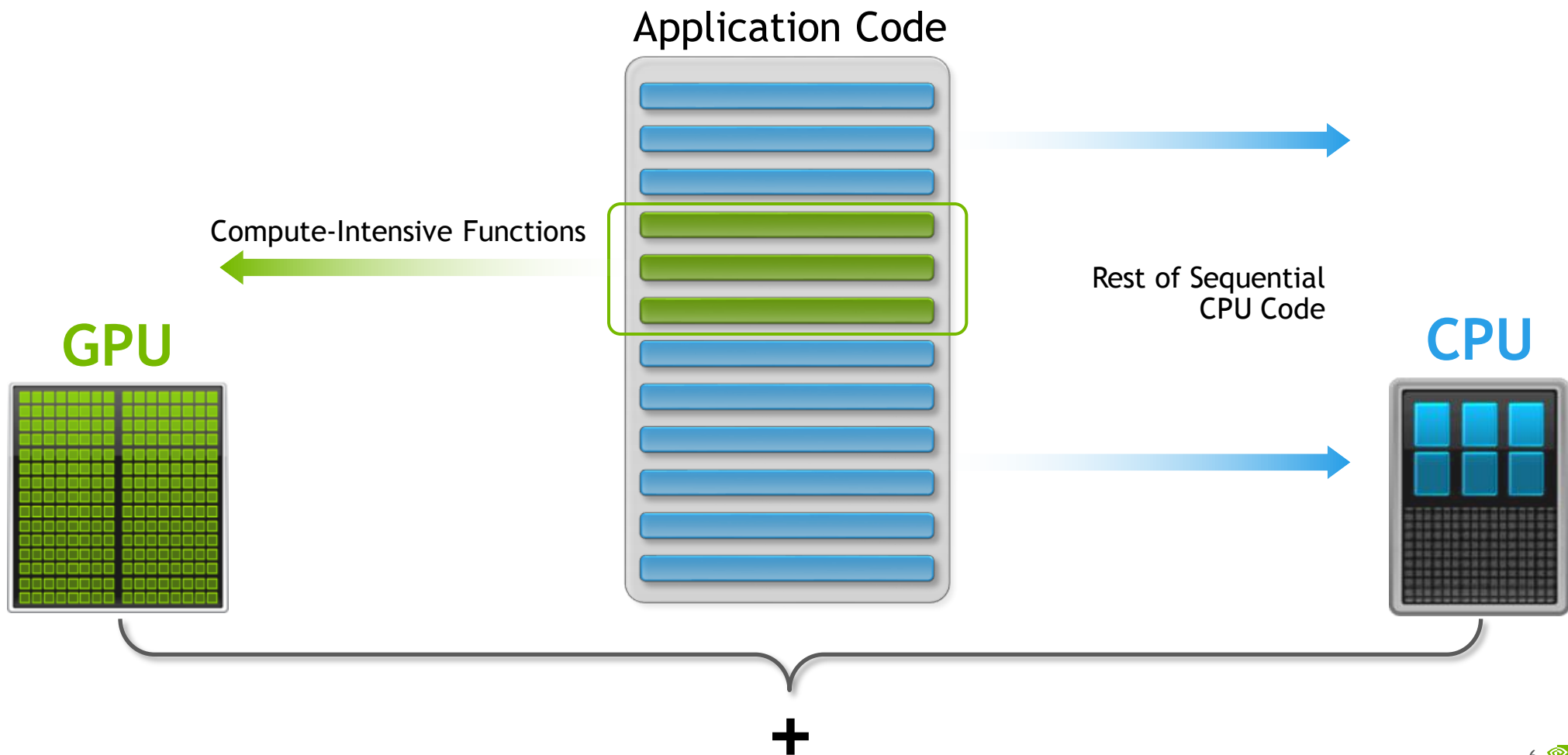**Detector & Electronics Group - ESRF**

The European Synchrotron | **ESRF**

# AGENDA

GPU overview

Getting data onto the GPU

Basic data processing

Advanced data processing

# HOW GPU ACCELERATION WORKS

Application Code

GPU

Compute-Intensive Functions

Rest of Sequential
CPU Code

CPU

+

NVIDIA.

# HETEROGENEOUS ARCHITECTURES
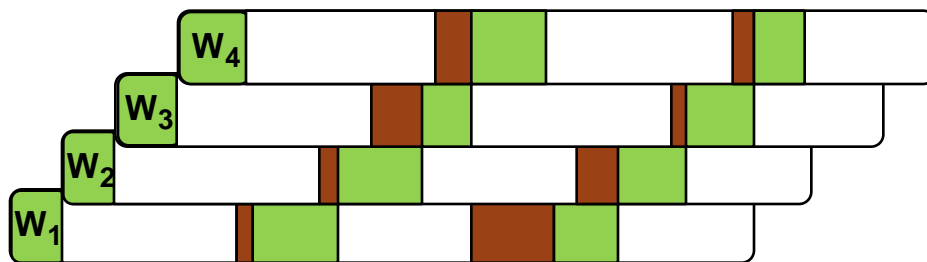
# LOW LATENCY OF HIGH THROUGHPUT?

CPU architecture must **minimize latency** within each thread

GPU architecture **hides latency** with computation from other threads (warps)
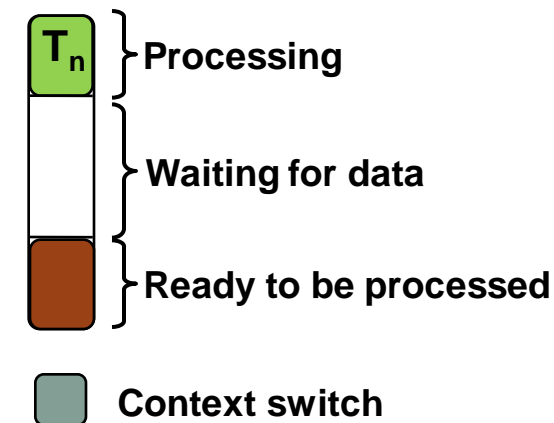
**CPU core – Low Latency Processor**

**GPU Stream Multiprocessor – High Throughput Processor**
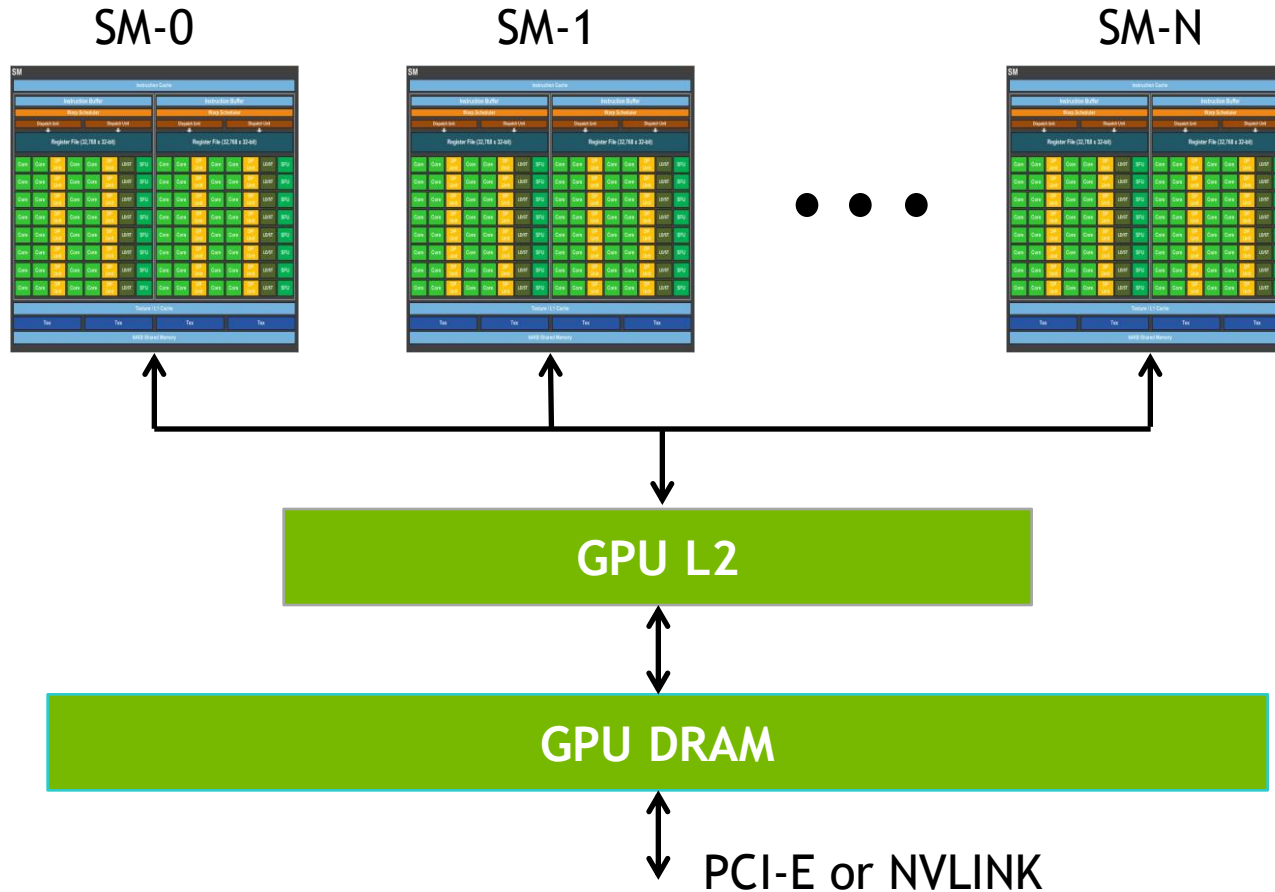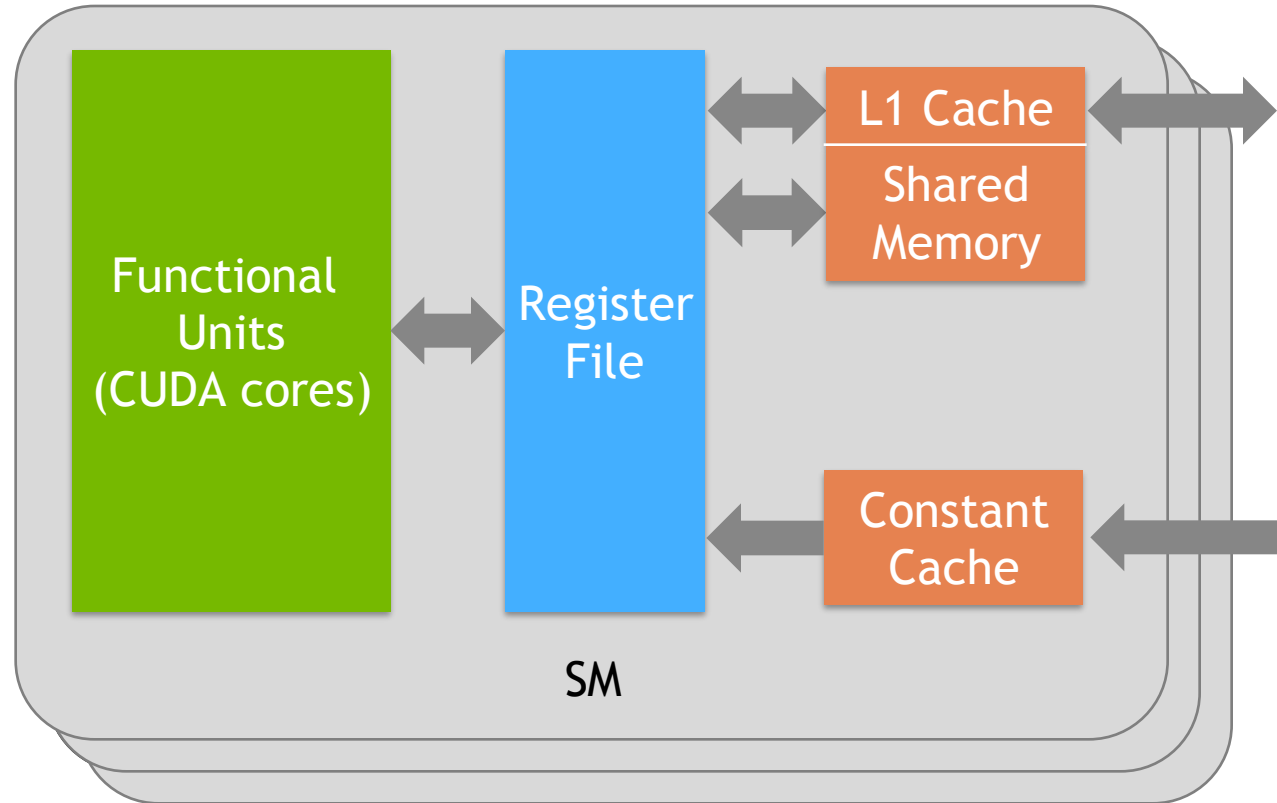
**Computation Thread/Warp**

$T_n$ Processing

Waiting for data

Ready to be processed

Context switch

NVIDIA.

# GPU ARCHITECTURE

SM-0            SM-1                    SM-N



**GPU L2**

**GPU DRAM**

PCI-E or NVLINK

# GPU SM ARCHITECTURE
## Volta SM

| | GV100 |
|---|---|
| FP32 Cores | 64 |
| FP64 Cores | 32 |
| Tensor Cores | 8 |
| Register File | 256 KB |
| Shared Memory | up to 96 KB |

Functional Units (CUDA cores)

Register File

L1 Cache

Shared Memory

Constant Cache

SM

80 SMs on Tesla V100

# NVIDIA TESLA V100

- 21B transistors
  815 mm$^2$, 12nm FFN

- 80 SM
  5120 CUDA Cores
  640 Tensor Cores

- 7.8 FP64 TFLOPS

- 15.6 FP32 TFLOPS

- 125 Tensor TFLOPS

- 16 GB HBM2
  900 GB/s memory bandwidth
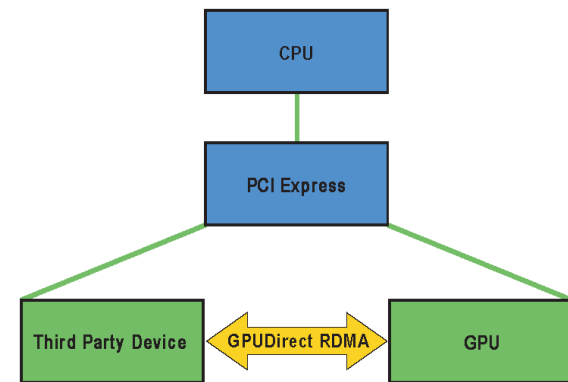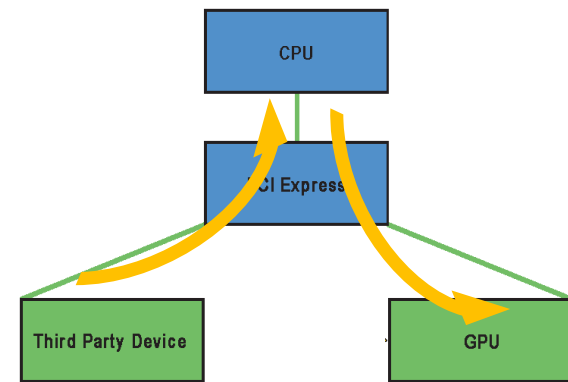
- 300 GB/s NVLink bandwidth



*full GV100 chip contains 84 SMs

# GETTING DATA ONTO THE GPU

- Transfer via host-memory

  - Simple if PCIe not saturated

  - 3$^{rd}$ Party Device/GPU not necessarily on same PCI root

  - Use asynchronous mem copies

- GPUDirect RDMA

  - Pin 3$^{rd}$ party device physical BAR addresses against GPU userspace addresses (3$^{rd}$ party device driver mod)

  - Devices on same PCI complex

  - Tesla/Quadro only

    Docu:    http://docs.nvidia.com/cuda/gpudirect-rdma
    Sample:  https://github.com/NVIDIA/gdrcopy



NVIDIA.

# TESLA FAMILY
## GPU comparison (boost clocks)

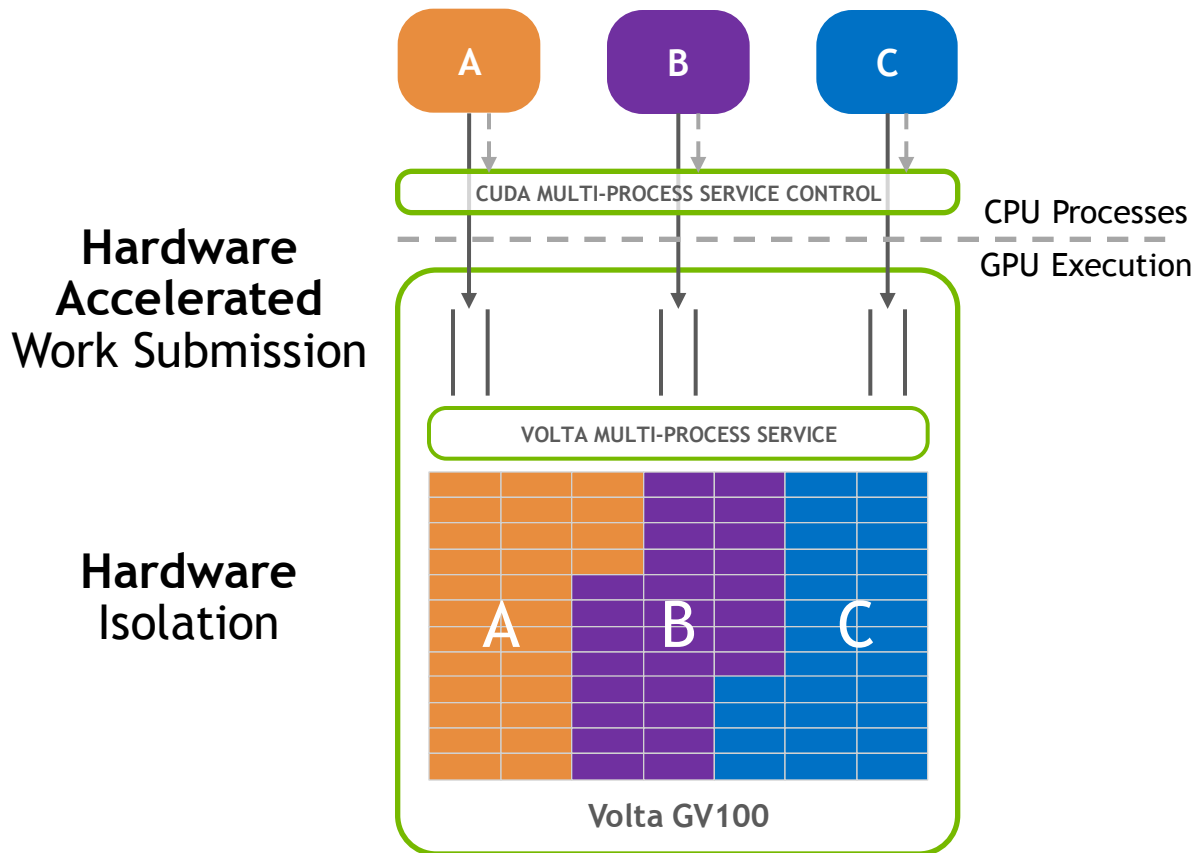|  | Tesla K40 | Tesla P100 | Tesla V100 |
|---|---|---|---|
| Peak FP32 (TFLOP/s) | 5.04 | 10.6 | 15 |
| Peak FP64 (TFLOP/s) | 1.68 | 5.3 | 7.5 |
| Peak Tensor Core (TFLOP/s) | N/A | N/A | 120 |
| Memory Size (GB) | 12 | 16 | 16 |
| Memory Bandwidth (GB/s) | 288 | 732 | 900 |

NVIDIA.

# LATENCY HIDING



Need sufficient work!  > 120000 threads in flight at any time

# VOLTA MULTI-PROCESS SERVICE

- GPUs getting wide

- Some problems too small to fill entire GPU

- Share GPU amongst multiple CPU processes

# GPU FOR SIGNAL PROCESSING

Bit Twiddling

- Integer INT32 instruction throughput: same as floating point (FP32) throughput *

- Half precision FP16 throughput 2xFP32 (sign, 5 exp, 10 mant). 10 bit integers.

- FP16 matrix-product cores (Tensor Cores), implicit FP16 -> FP32

    http://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#arithmetic-instructions

Rearranging data

- Relaxed coalescence constraints (32B transaction size)

- Shared memory on SM for fast, irregular access

Visualizing

NVIDIA.

# HARDWARE ACCELERATED COMPRESSION ENGINES

## Video compression and more

Multiple encoder/decoder per chip

(3 enc, 1 dec on V100)
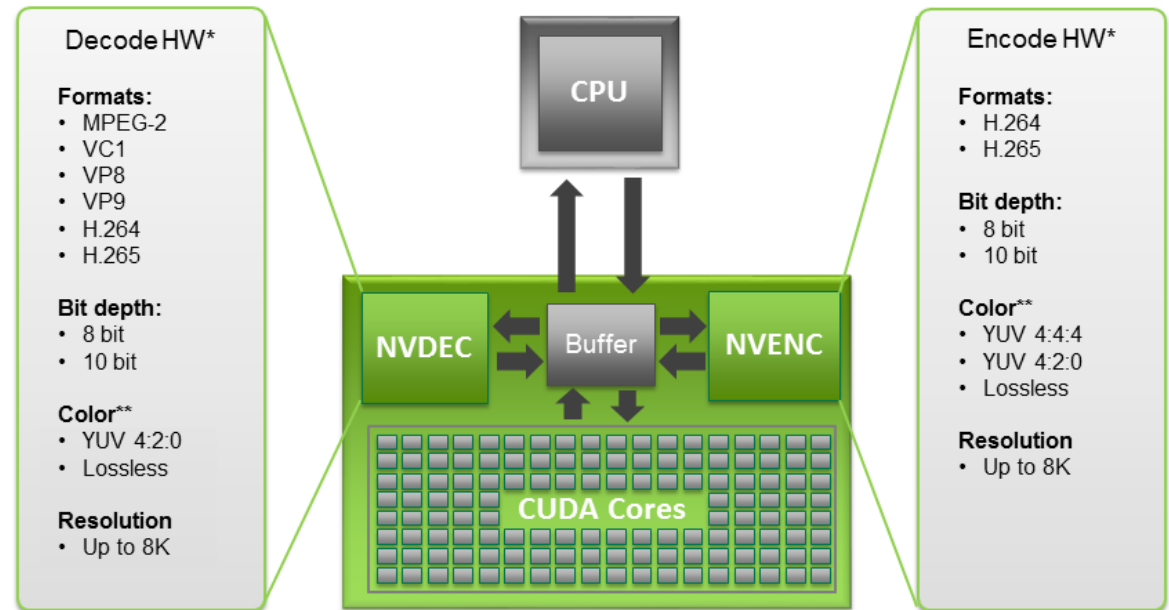
Transparently managed

Lots of formats supported, incl. lossless

I-frame, P-frame support

No CUDA/Rendering resources used

API: NVCODEC

https://developer.nvidia.com/nvidia-video-codec-sdk



**Decode HW***

**Formats:**
- MPEG-2
- VC1
- VP8
- VP9
- H.264
- H.265

**Bit depth:**
- 8 bit
- 10 bit

**Color****
- YUV 4:2:0
- Lossless

**Resolution**
- Up to 8K

**CPU**

**NVDEC**   **Buffer**   **NVENC**

**CUDA Cores**

**Encode HW***

**Formats:**
- H.264
- H.265

**Bit depth:**
- 8 bit
- 10 bit

**Color****
- YUV 4:4:4
- YUV 4:2:0
- Lossless

**Resolution**
- Up to 8K

NVIDIA.

# VIDEO FOR DATA COMPRESSION

## Not only for streaming
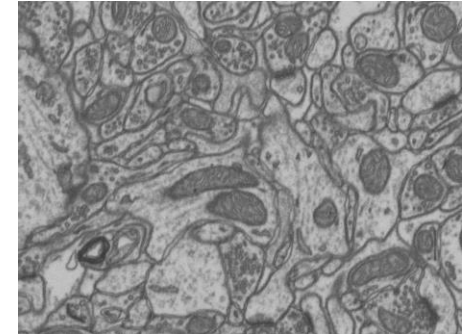
Measured data cubes exhibit lots of coherency
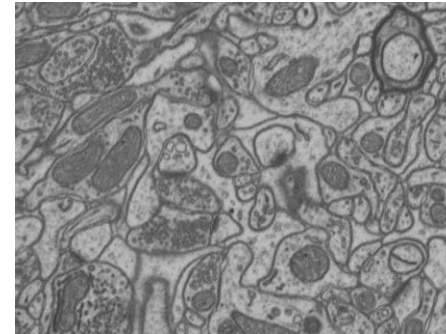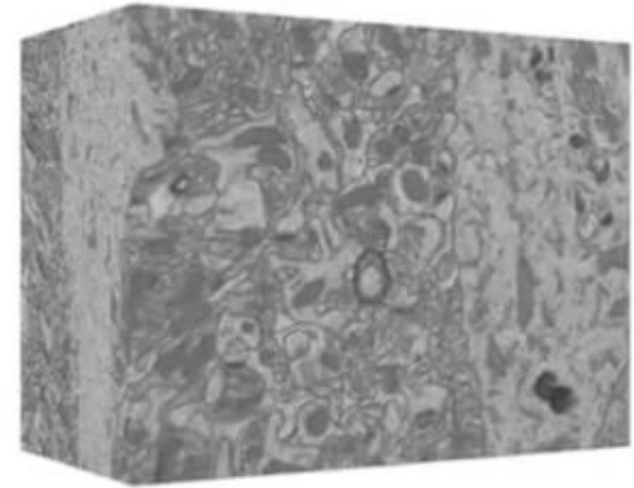
Huge compression possible

    3.2GB -> 51 MB (-> 21MB)

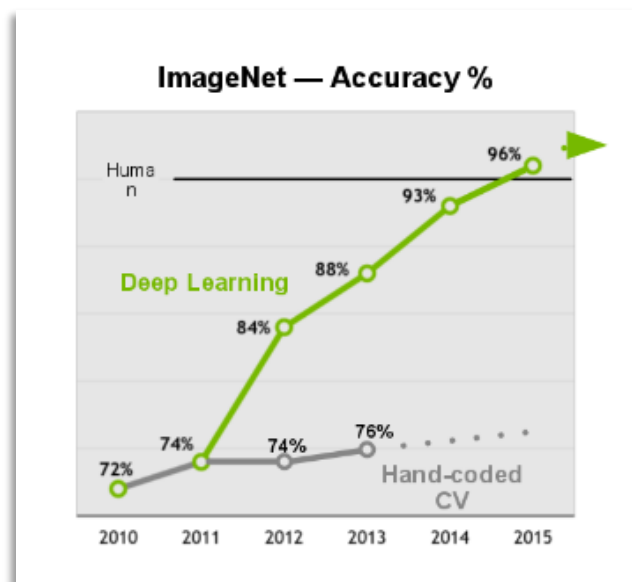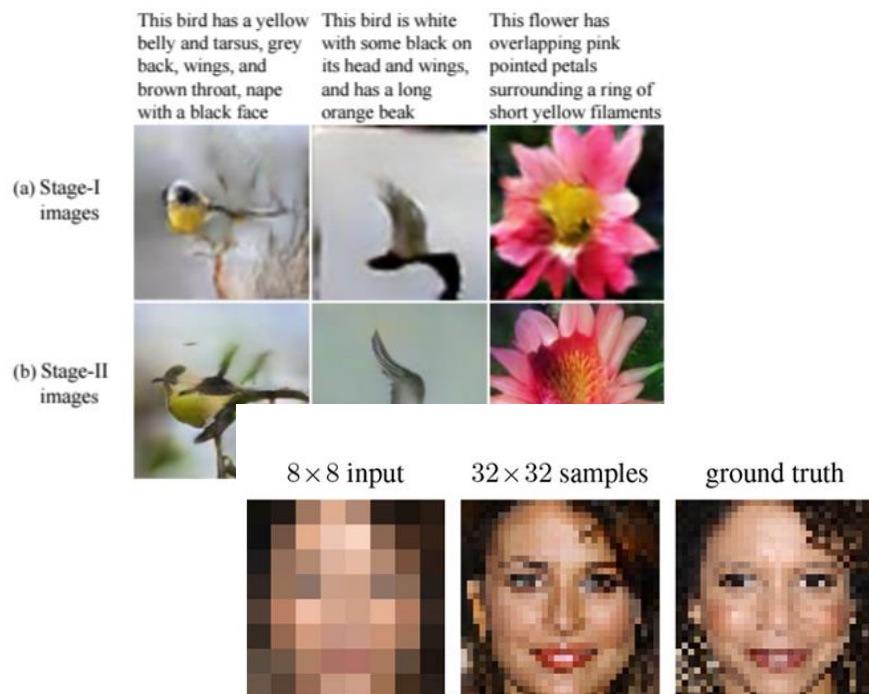Opportunity for large volume data

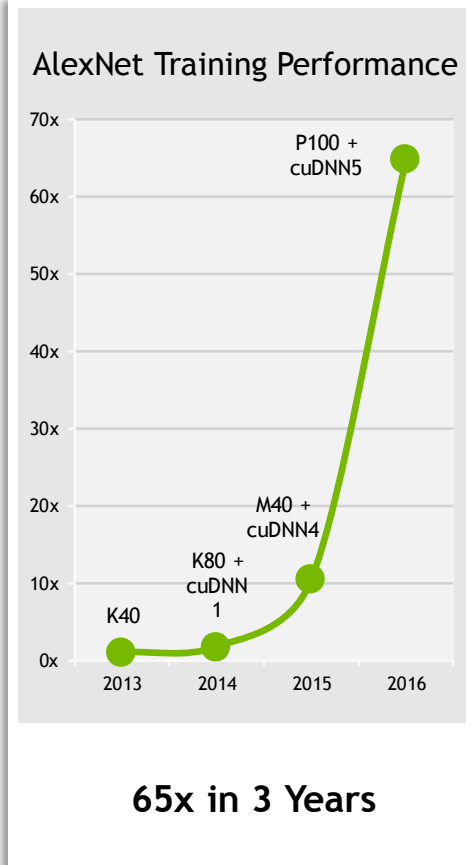GPU-GPU use case not yet supported

Get in touch for more details

Data and images: http://cvlab.epfl.ch/data/em

# WHY THE AI EXCITEMENT?

## GPUs as Enablers of Breakthrough Results

This bird has a yellow belly and tarsus, grey back, wings, and brown throat, nape with a black face

This bird is white with some black on its head and wings, and has a long orange beak

This flower has overlapping pink pointed petals surrounding a ring of short yellow filaments

(a) Stage-I images

(b) Stage-II images

8 × 8 input   32 × 32 samples   ground truth

Dahl et al. 2017

ImageNet — Accuracy %

Human

96%

93%

**Deep Learning**

88%

84%

74%       74%       76%

72%                           Hand-coded CV

2010  2011  2012  2013  2014  2015

AlexNet Training Performance

70x

P100 + cuDNN5

60x

50x

40x

30x

M40 + cuDNN4

20x

K80 + cuDNN1

10x

K40

0x

2013   2014   2015   2016

**65x in 3 Years**

We can generate photorealistic images from __textual__ descriptions and super-enhance blurry photos!

Achieve super-human accuracy in classification

And we are getting faster fast

# 1-SLIDE INTRO TO CONVOLUTIONAL NEURAL NETS

# AI INFERENCING IS EXPLODING

**2 Trillion**
Messages Per Day On LinkedIn

**PERSONALIZATION**

**500M**
Daily active users of iFlyTek

**SPEECH**

**140 Billion**
Words Per Day Translated by Google

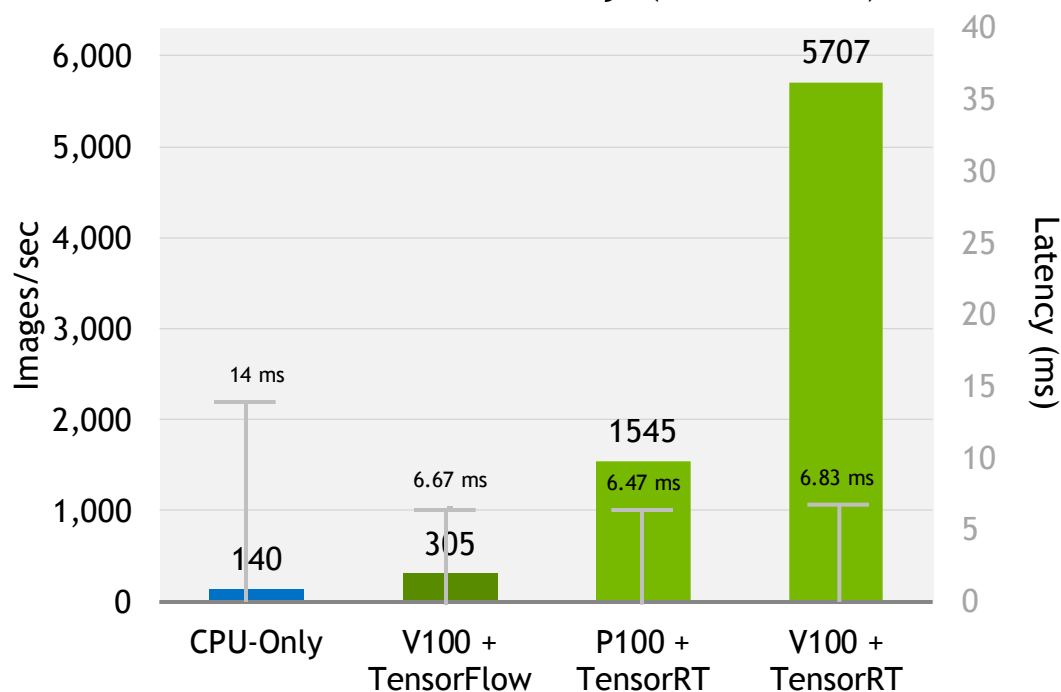**TRANSLATION**

**60 Billion**
Video frames/day uploaded on Youtube

**VIDEO**

# TENSORRT 3 PERFORMANCE

## 40x Faster CNNs on V100 vs. CPU-Only Under 7ms Latency (ResNet50)



Image classification network
224x224 pixels images
1000 categories
50 layers, mostly convolutions

Under 7ms latency from pixels to class

Inference throughput (images/sec) on ResNet50. **V100 + TensorRT**: NVIDIA TensorRT (FP16), batch size 39, Tesla V100-SXM2-16GB, E5-2690 v4@2.60GHz 3.5GHz Turbo (Broadwell) HT On. **P100 + TensorRT**: NVIDIA TensorRT (FP16), batch size 10, Tesla P100-PCIE-16GB, E5-2690 v4@2.60GHz 3.5GHz Turbo (Broadwell) HT On **V100 + TensorFlow**: Preview of volta optimized TensorFlow (FP16), batch size 2, Tesla V100-PCIE-16GB, E5-2690 v4@2.60GHz 3.5GHz Turbo (Broadwell) HT On. **CPU-Only:** Intel Xeon-D 1587 Broadwell-E CPU and Intel DL SDK. Score doubled to comprehend Intel's stated claim of 2x performance improvement on Skylake with AVX512.
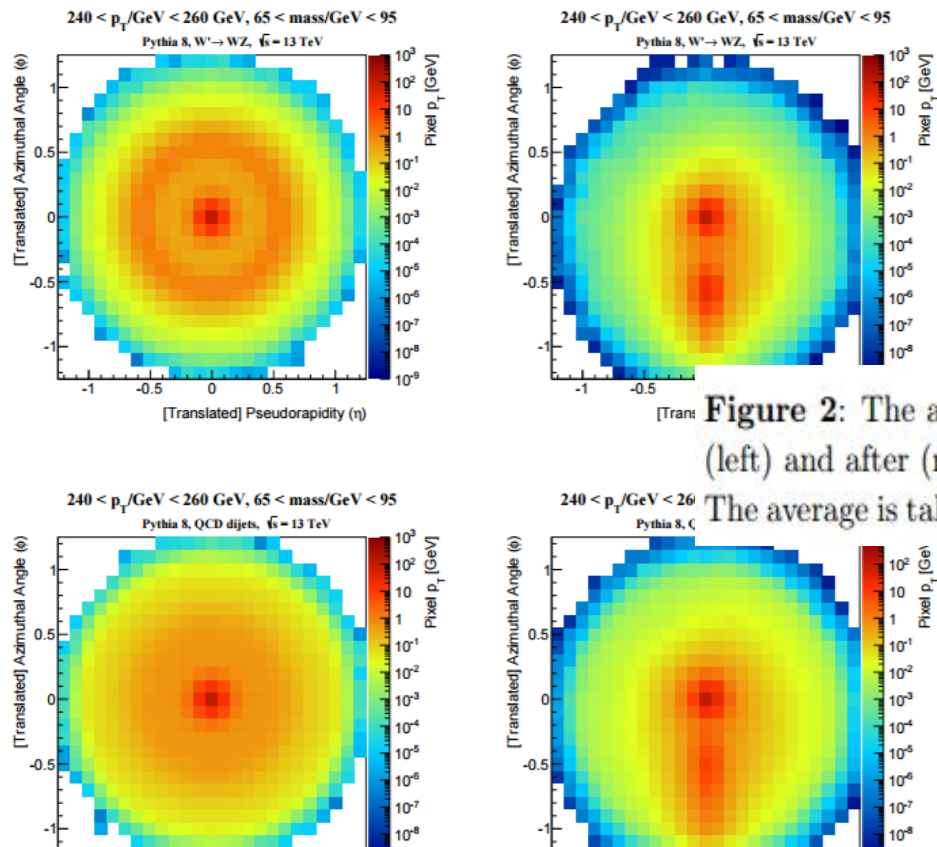
# NORMALIZING AND PRE-PROCESSING

## DL trained on jet images vs. physically-motivated feature driven approaches



**Figure 2**: The average jet image for signal $W$ jets (top) and background QCD jets (bottom) before (left) and after (right) applying the rotation, re-pixelation, and inversion steps of the pre-processing. The average is taken over images of jets with $240\ \text{GeV} < p_T < 260\ \text{GeV}$ and $65\ \text{GeV} < \text{mass} < 95\ \text{GeV}$.

# Jetson TX1 Module

January 2016 availability

| | JETSON TX1 MODULE |
|---|---|
| GPU | 1 TFLOP/s 256-core Maxwell |
| CPU | 64-bit ARM A57 CPUs |
| Memory | 4 GB LPDDR4 | 25.6 GB/s |
| Video decode | 4K 60Hz |
| Video encode | 4K 30Hz |
| CSI | Up to 6 cameras | 1400 Mpix/s |
| Display | 2x DSI, 1x eDP 1.4, 1x DP 1.2/HDMI |
| Wifi | 802.11 2x2 ac |
| Networking | 1 Gigabit Ethernet |
| PCIE | Gen 2 1x1 + 1x4 |
| Storage | 16 GB eMMC, SDIO, SATA |
| Other | 3x UART, 3x SPI, 4x I2C, 4x I2S, GPIOs |

# JETSON EMBEDDED PLATFORM

A comprehensive platform for development and deployment of advanced embedded products

## Jetson SDK

L4T kernel     Platform drivers
Bootloader
Reference FS



## Jetson TX1 module

Unmatched performance
Production-ready
Long lifetime



## Library support

cuDNN     OpenVX
CUDA libs     OpenCV, ROS
VisionWorks     OpenGL



## Developer kit

Module
Breakout board
5MP camera



## Developer tools

Nsight     Debuggers
NVTX     Profilers
Jetpack installer



## Design collateral

Ref. design files     Reference diags
Design guides     Factory users guide
TRM



## Developer portal

Downloads     White papers
Training     Success stories
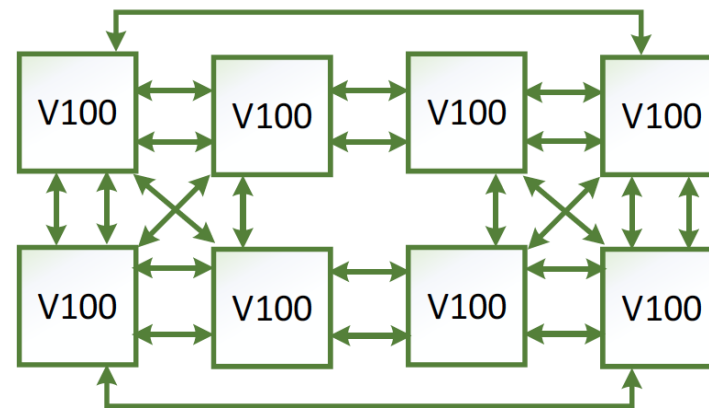Forum     Partner pages



## Ecosystem

Design partners     Higher Ed
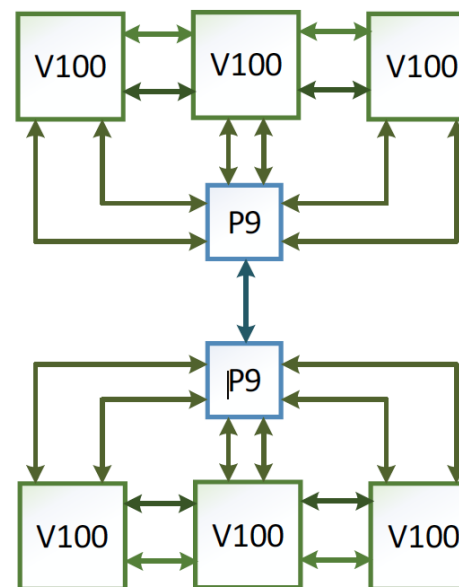IHVs     Hacker/maker
ISVs     OSS

# VOLTA NVLINK

- 6 NVLINKS @ 50 GB/s bidirectional

- Reduce number of lanes for lightly loaded link (Power savings)

- Coherence features for NVLINK enabled CPUs



Hybrid cube mesh
(eg. DGX1V)

POWER9 based node

NVIDIA.

# SUMMARY

- Getting data to GPU directly via RDMA-Direct

- Perform processing, even bit-twiddling, on the GPU. Can FP16 be used?

- GPUs are wide and have lots of memory. Run many things concurrently

- Compression via hardware video encoders.

- Advanced processing on the fly via Deep Learning methods

- Jetson for compute on the edge

- NVLInk for fast exchange of data between GPUs

NVIDIA.