

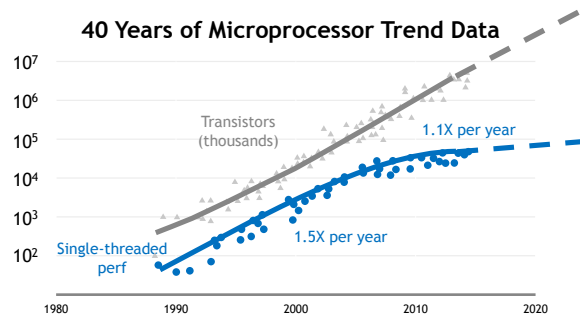
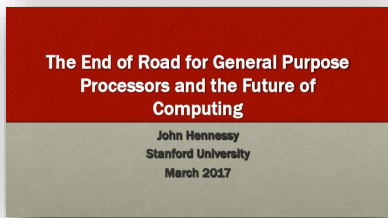
# ENDURING DIFFERENTIATION

Timothy Lanfear



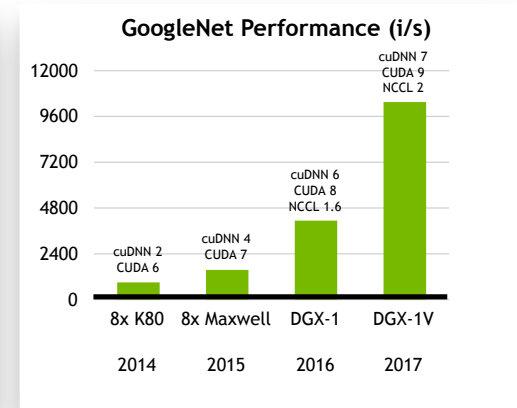
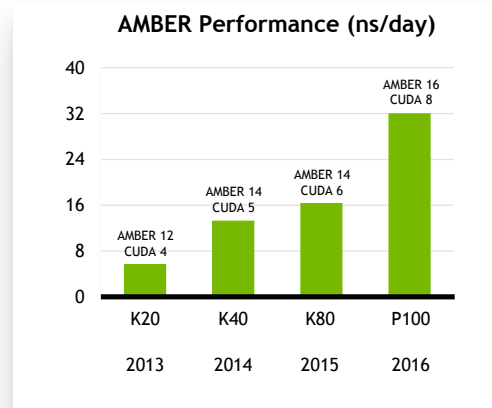
WHERE ARE WE?

## LIFE AFTER DENNARD SCALING



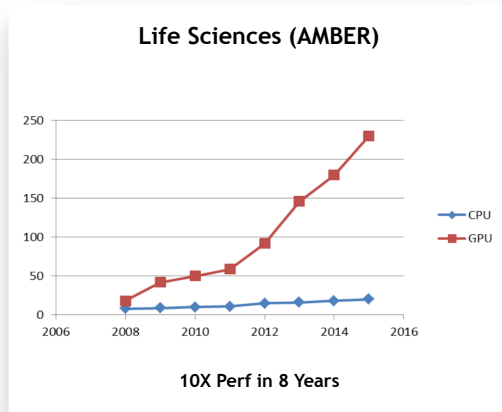
Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten. New plot and data collected for 2010-2015 by K. Rupp.

## GPU-ACCELERATED PERFORMANCE

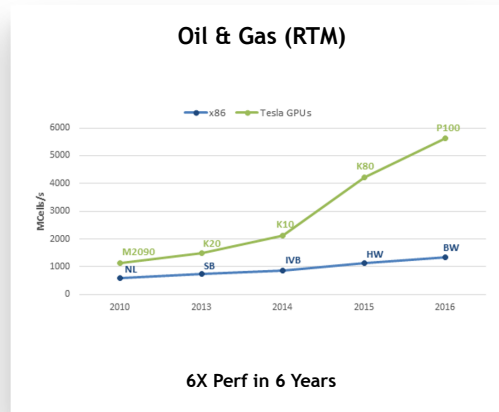


## TESLA PLATFORM ADVANTAGE

Delivered value grows over time



Amber performance: Nano Seconds Per Day delivered on 1xServer with GPUs and CPUs

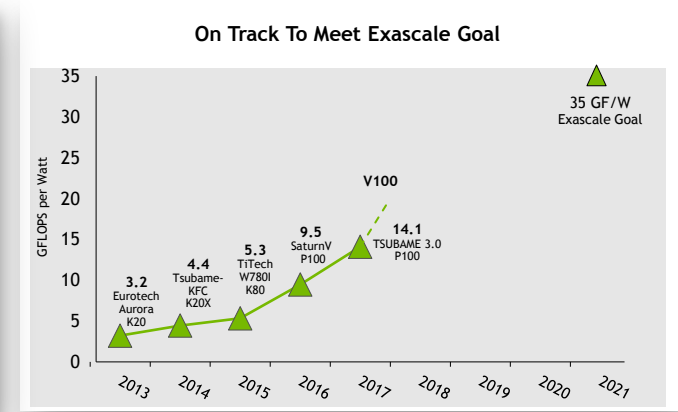


5 NVIDIA

## GPU-ACCELERATED EFFICIENCY

13/13 Greenest Supercomputers  
Powered by Tesla P100

TSUBAME 3.0  
Kukai  
AIST AI Cloud  
RAIDEN GPU subsystem  
Piz Daint  
Wilkes-2  
GOSAT-2 (RCF2)  
DGX Saturn V  
Reedbush-H  
JADE  
Facebook Cluster  
Cedar  
DAVIDE



## HOW ARE WE DOING THIS?

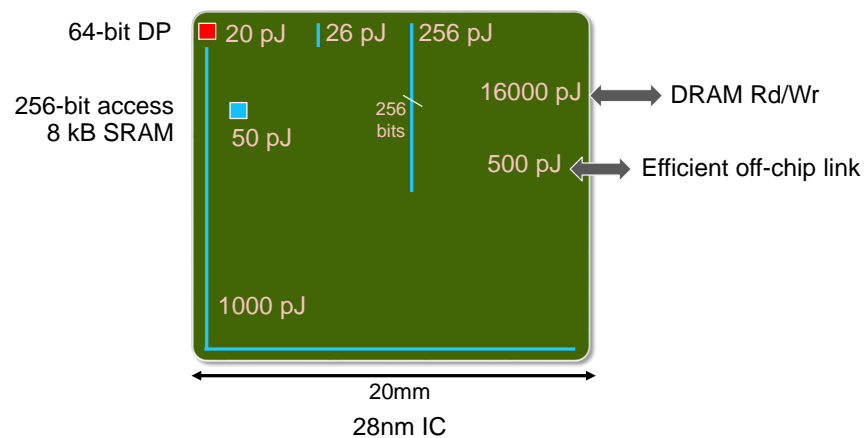
And, is our differentiation sustainable?

- What are the most important dimensions of our differentiation?
- Why are GPUs so much more efficient than CPUs?
- How can we continue scaling performance/efficiency as Moore's Law fades?
- Why can't competitors replicate GPU efficiency, performance, scaling, etc., with lots of weak CPU cores? (e.g., Intel KNC/KNL/KNM)
- How is optimizing GPUs for AI affecting their suitability for HPC?

7 NVIDIA

## ENERGY EFFICIENCY

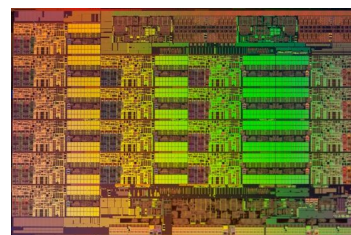
## COMPUTATION VERSUS COMMUNICATIONS



9 NVIDIA

## CPU 126 pJ/flop (SP)

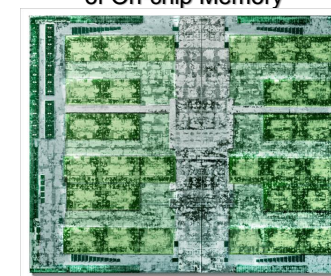
Optimized for Latency  
Deep Cache Hierarchy



Broadwell E5 v4  
14 nm

## GPU 28 pJ/flop (SP)

Optimized for Throughput  
Explicit Management of On-chip Memory

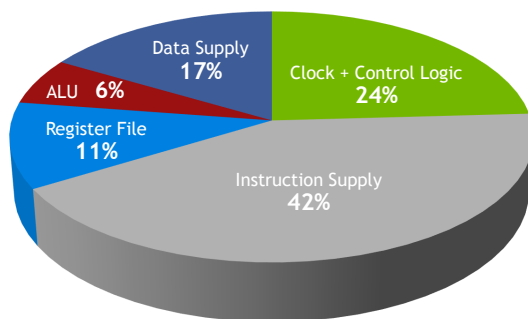


Pascal  
16 nm

10 NVIDIA

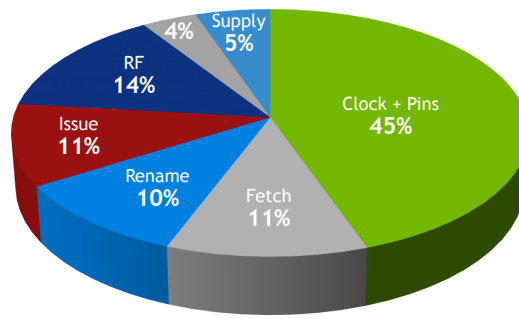
## HOW IS POWER SPENT IN A CPU?

### In Order, Embedded



Dally [2008] (Embedded in-order CPU)

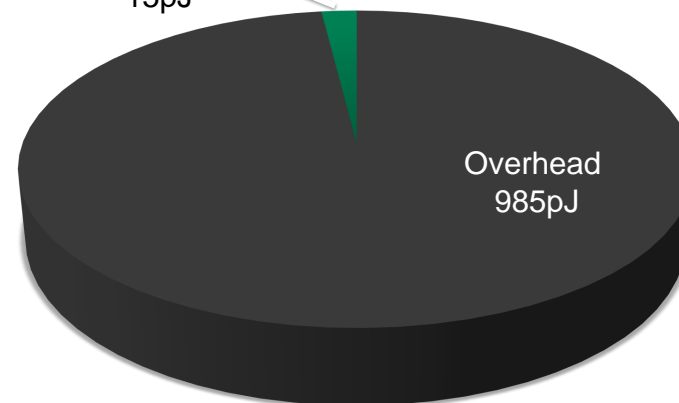
### Out of Order, High Performance



Natarajan [2003] (Alpha 21264)

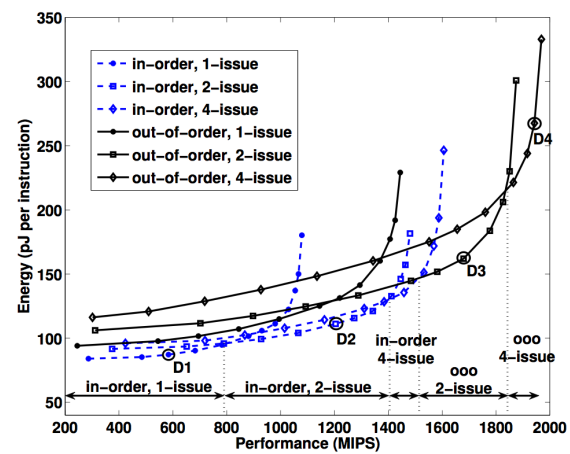
11 NVIDIA

Payload  
Arithmetic  
15pJ

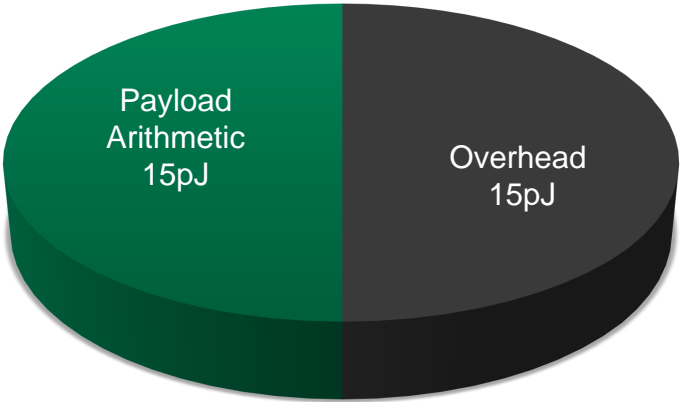


12 NVIDIA

# SIMPLER CORES = ENERGY EFFICIENCY

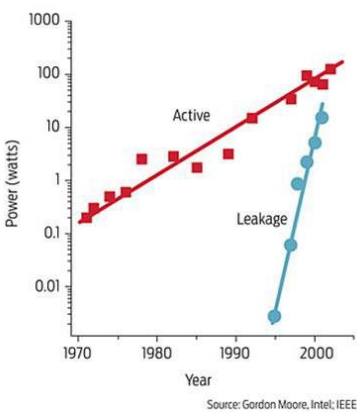


Source: Azizi [PhD 2010]



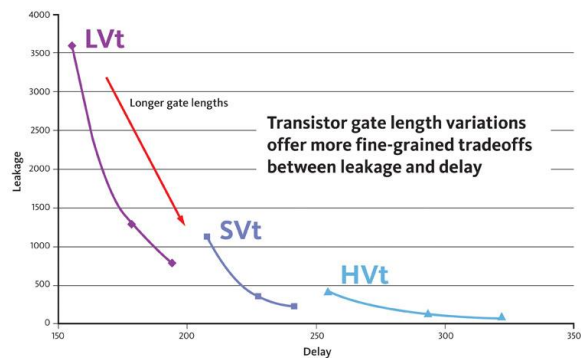
# THROUGHPUT PROCESSORS

# RISE OF LEAKAGE



Source: Gordon Moore, Intel, IEEE

## FREQUENCY VS. LEAKAGE

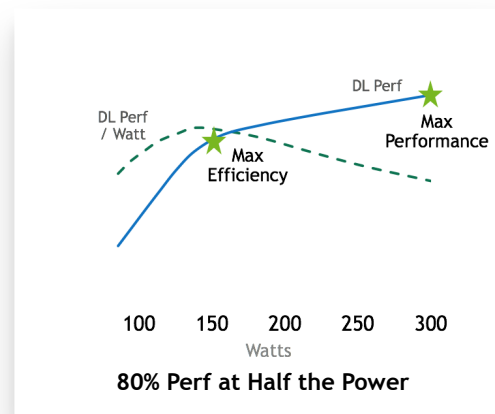


Source: Gordon Moore, Intel; IEEE

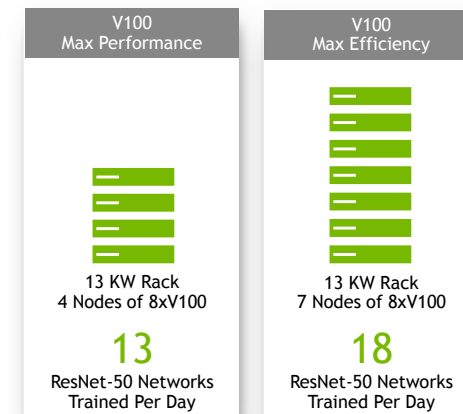
17 NVIDIA

## OPTIMIZED FOR DATACENTER EFFICIENCY

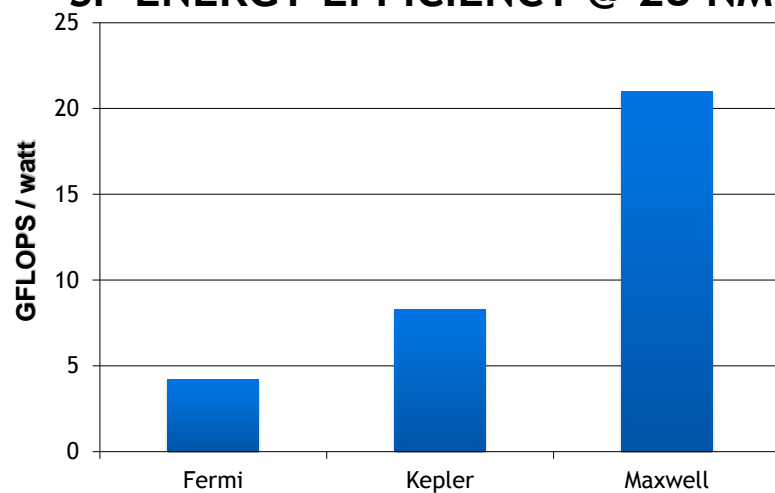
40% More Performance in a Rack



ResNet-50 Training, Max Efficiency run with V100@160W | V100 performance measured on pre-production hardware. 19 NVIDIA



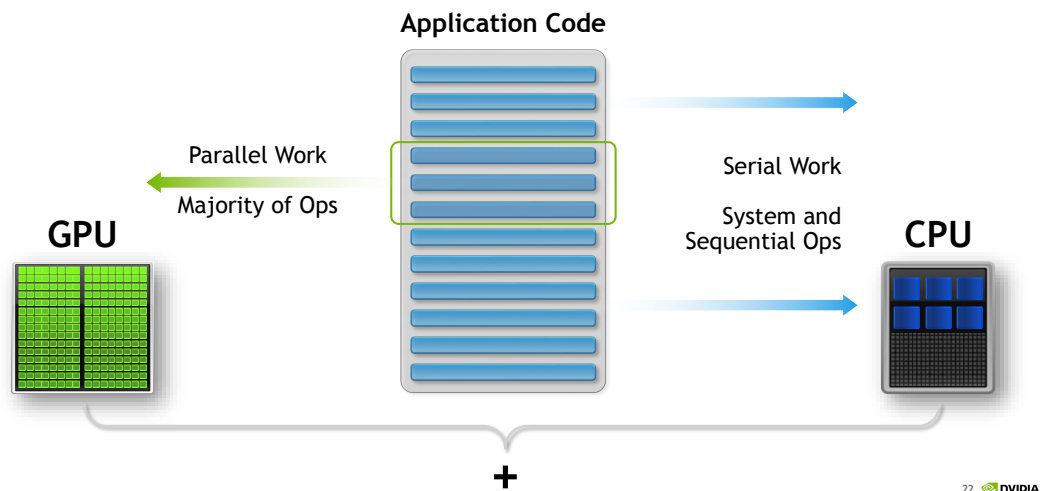
## SP ENERGY EFFICIENCY @ 28 NM



20 NVIDIA

## HETEROGENEOUS COMPUTING

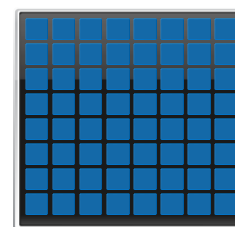
## OPTIMIZING SERIAL/PARALLEL EXECUTION



## TWO TYPES OF ACCELERATORS

**Many-Weak-Cores (MWC) Model**  
Single CPU Core for Both Serial & Parallel Work

**Xeon Phi (And Others)**  
Many Weak Serial Cores

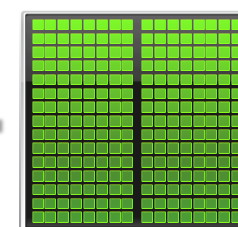


**Heterogeneous Computing Model**  
Complementary Processors Work Together

**CPU**  
Optimized for Serial Tasks



**GPU Accelerator**  
Optimized for Parallel Tasks

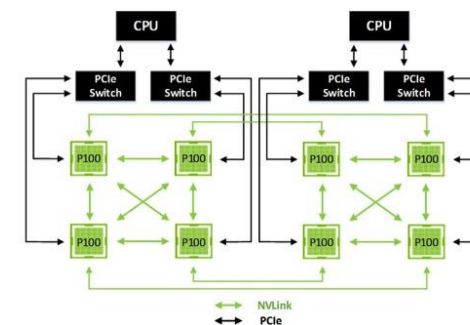


23

EXTENSIBILITY

## NVLINK: A MEMORY FABRIC, NOT A NETWORK

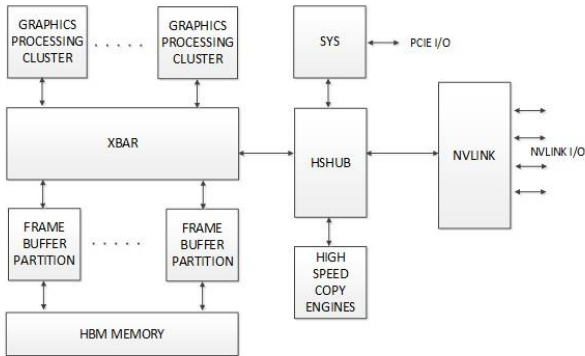
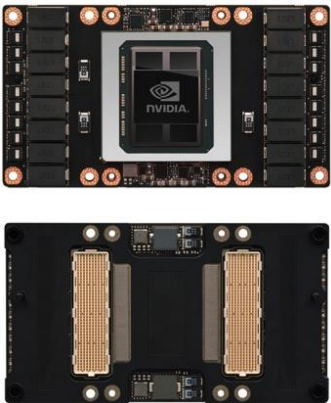
DGX-1: 8 NVLink-Connected GPUs



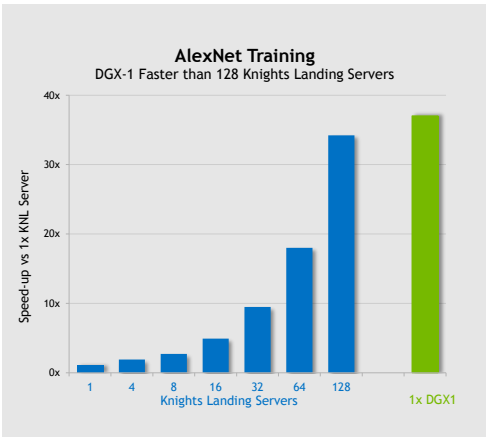
25

# LATENCY HIDING FOR LOAD/STORE/ATOMICS

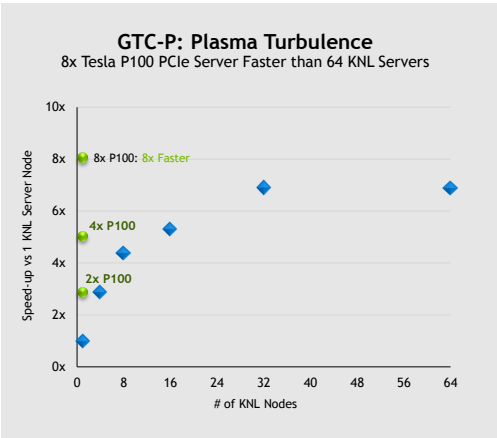
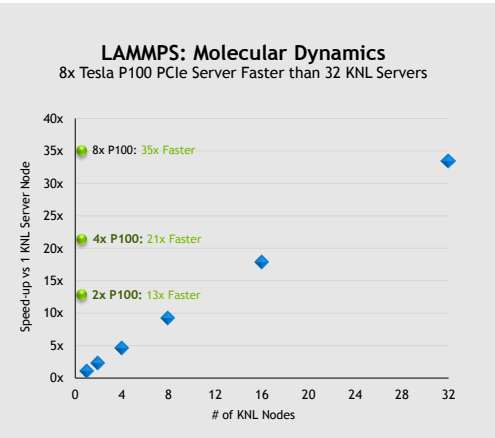
Where are the NICs? There are no NICs.



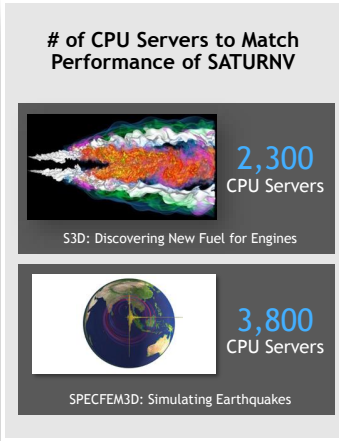
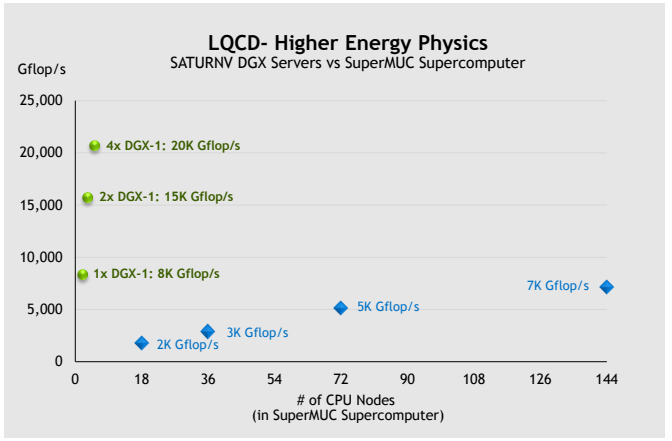
# STRONG SCALING



# STRONG SCALING



# STRONG SCALING



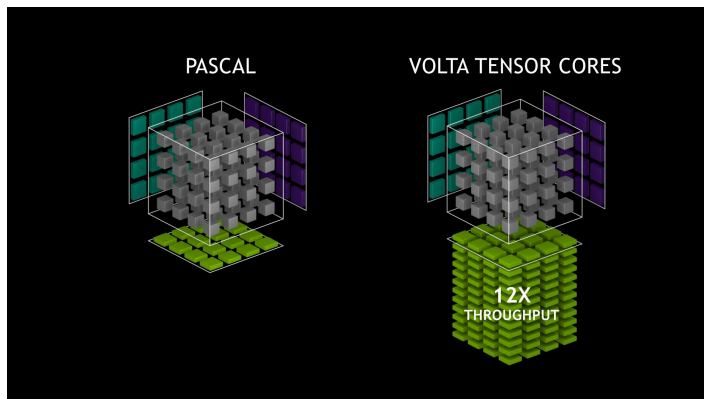
## NEW TENSOR CORE

New CUDA TensorOp instructions and data formats

4x4 matrix processing array

$$D_{FP32} = A_{FP16} \times B_{FP16} + C_{FP32}$$

Optimized for deep learning

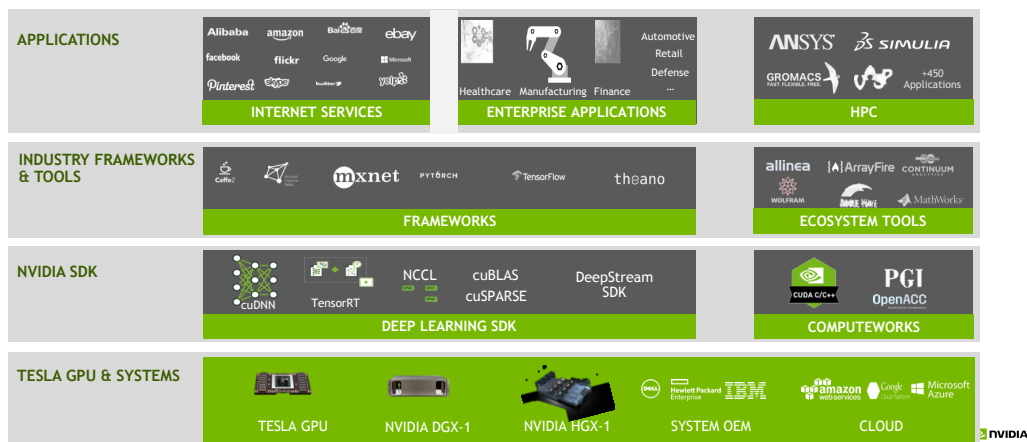


30 NVIDIA

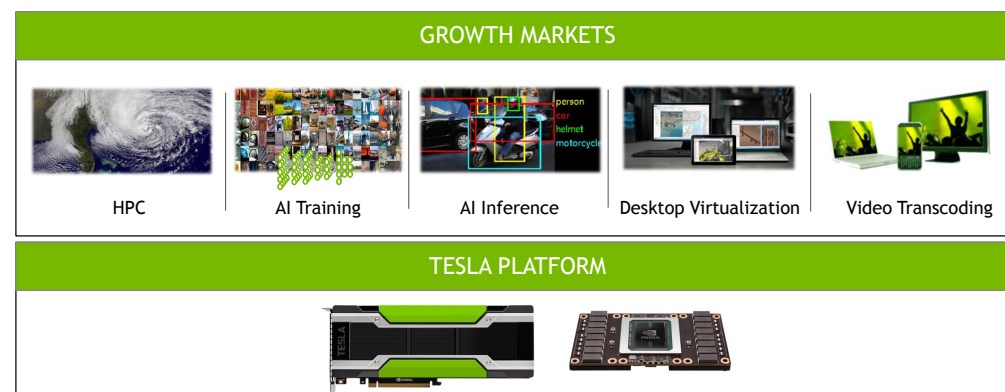
## TESLA PLATFORM

## TESLA IS A PLATFORM

World's Leading Data Center Platform for Accelerating HPC and AI



## MULTIPLE GROWTH MARKETS



33 NVIDIA



## CONCLUSION

## PASCAL TO VOLTA

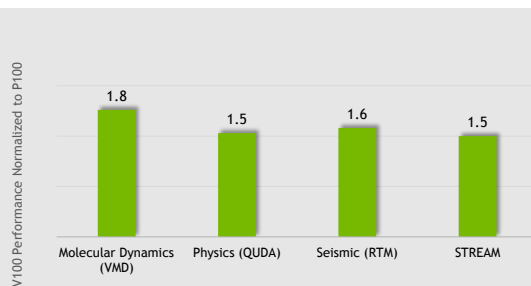
Architecture with Technology

- Area:  $\sim 600 \text{ mm}^2 \rightarrow \sim 800 \text{ mm}^2$  ( $\sim 33\%$  more area)
- Process:  $\sim$  small Pascal  $\rightarrow$  Volta improvement (a few percent)
- Clocks: similar dynamic range, power limited
- Memory BW (sustained): 50% improvement
- Communications (NVLink): 160 GB/s  $\rightarrow$  300 GB/s (almost double!)
- AI (Tensor Cores):  $\sim 20 \text{ TFLOPS} \rightarrow 120 \text{ TFLOPS}$  ( $\sim 6\times$ )

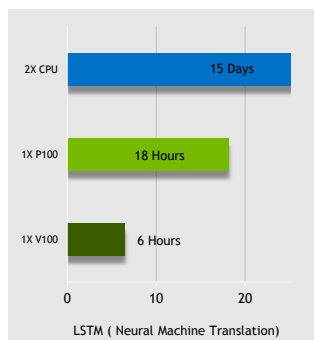
## REVOLUTIONARY PERFORMANCE FOR HPC AND AI

Single Platform For Data Science and Computation Science

1.5X HPC Performance In 1 Year



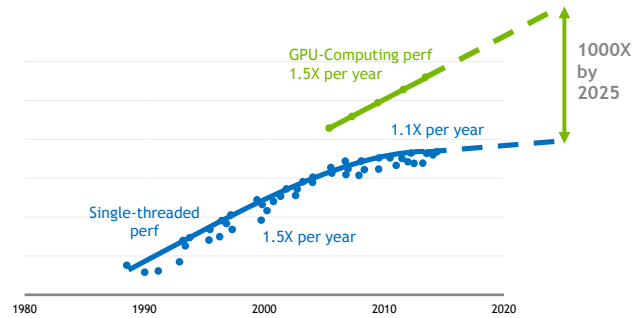
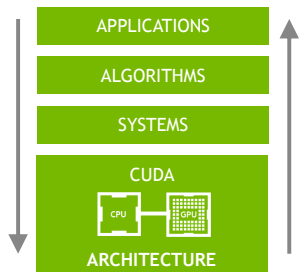
3X AI Performance in 1 Year



## GPU PERFORMANCE COMPARISON

	P100	V100	Ratio
Training acceleration	10 TOPS	120 TOPS	12x
Inference acceleration	21 TFLOPS	120 TOPS	6x
FP64/FP32	5/10 TFLOPS	7.5/15 TFLOPS	1.5x
HBM2 Bandwidth	720 GB/s	900 GB/s	1.2x
NVLink Bandwidth	160 GB/s	300 GB/s	1.9x
L2 Cache	4 MB	6 MB	1.5x
L1 Caches	1.3 MB	10 MB	7.7x

# GPU TRAJECTORY



Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten New plot and data collected for 2010-2015 by K. Rupp

