



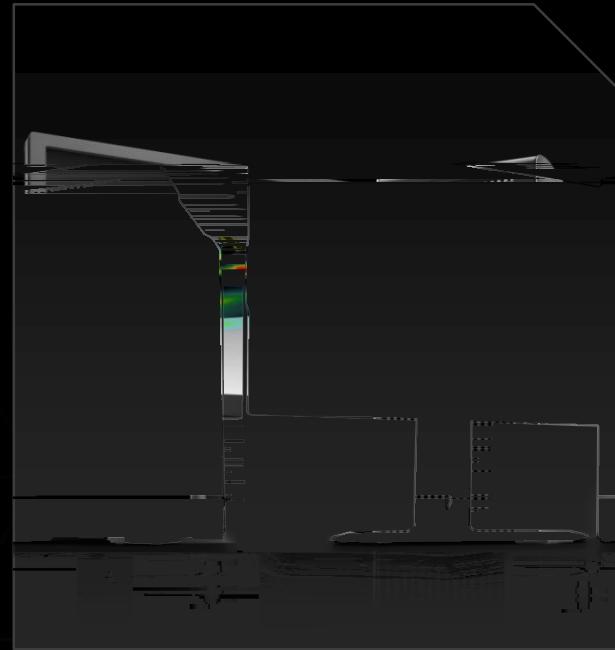
NVIDIA GPU

17/4/2015





GeForce | GRID



Quadro | Tesla | GRID

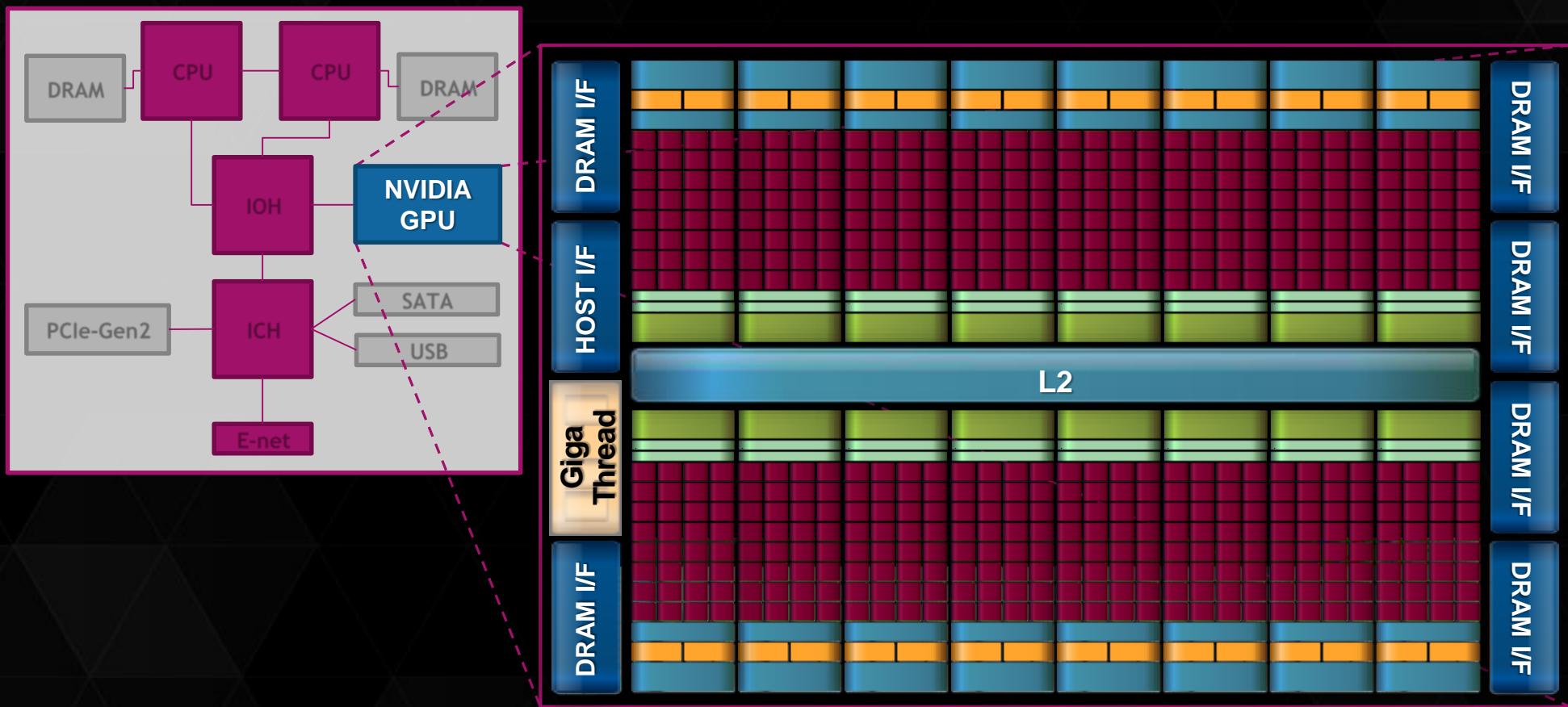


AUTO
Tegra

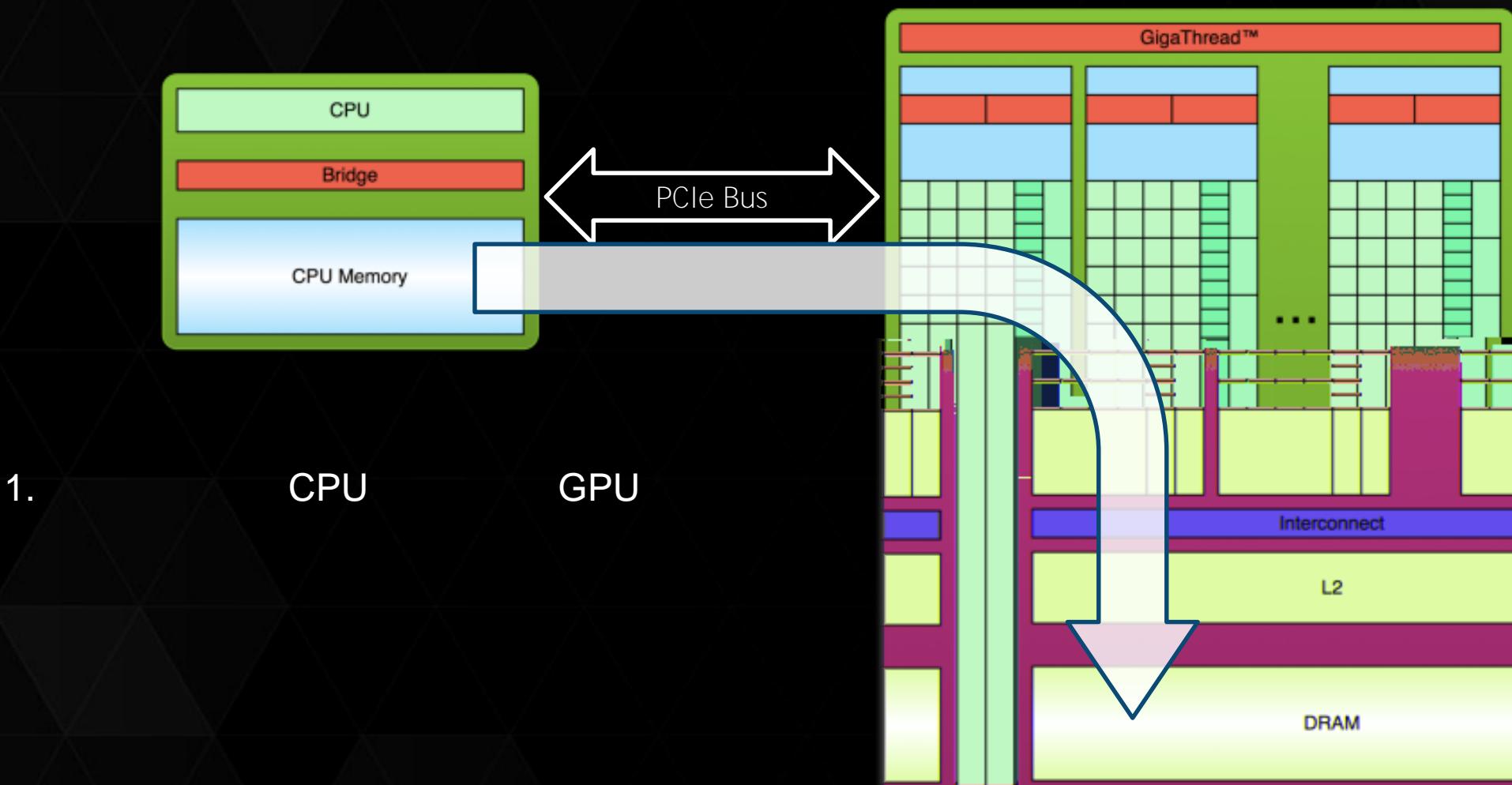
GPU

GPU

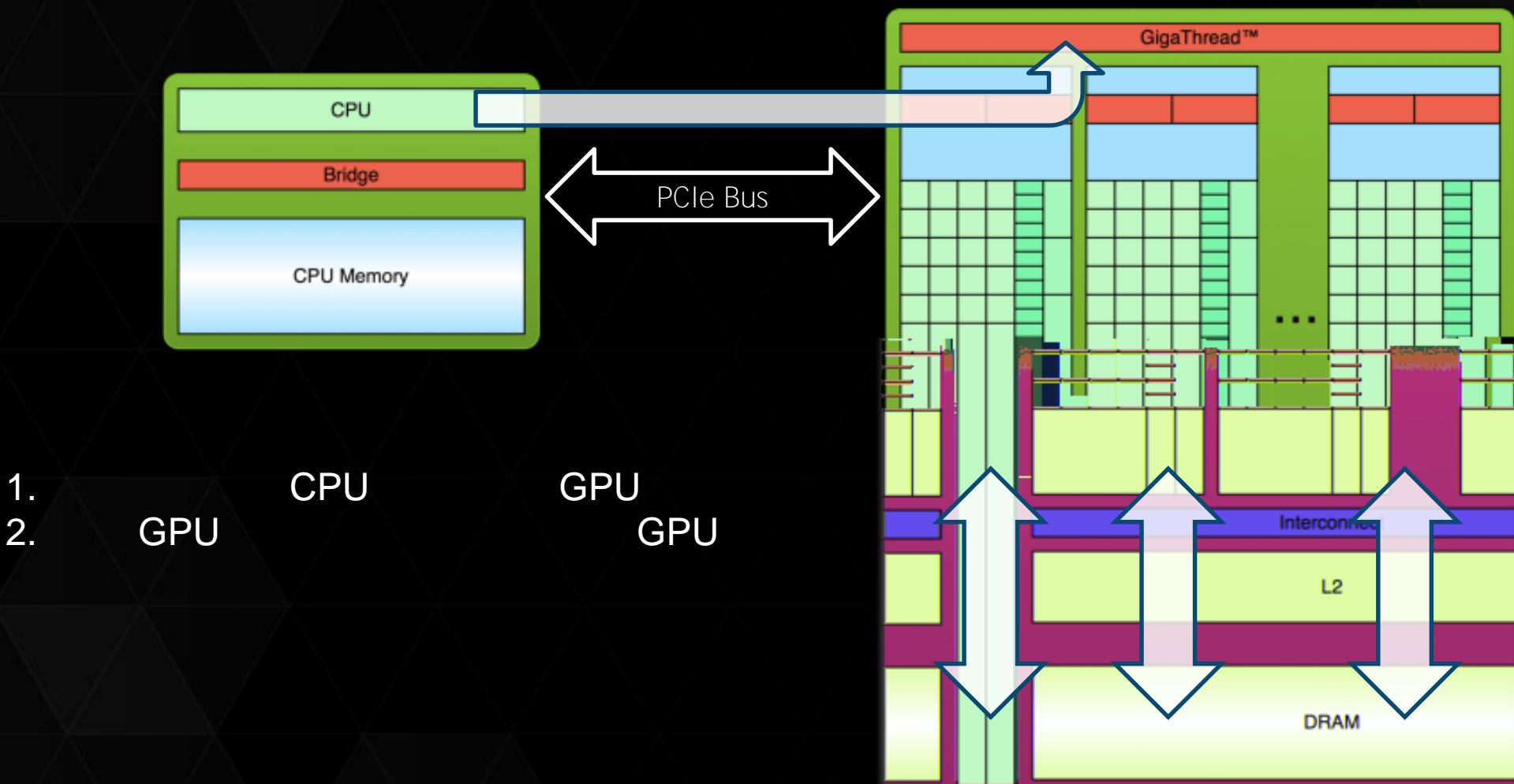
x86 + GPU



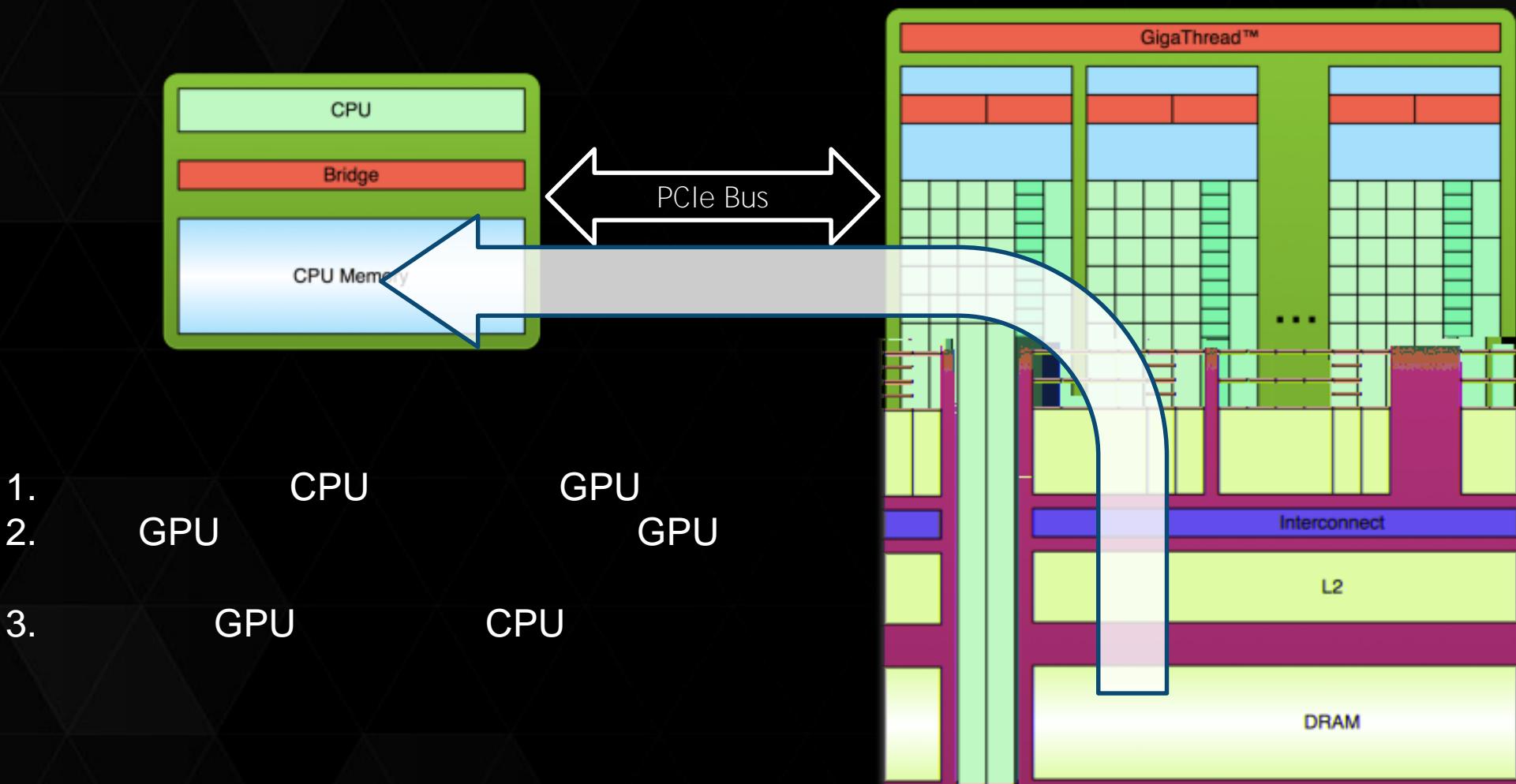
CPU + GPU



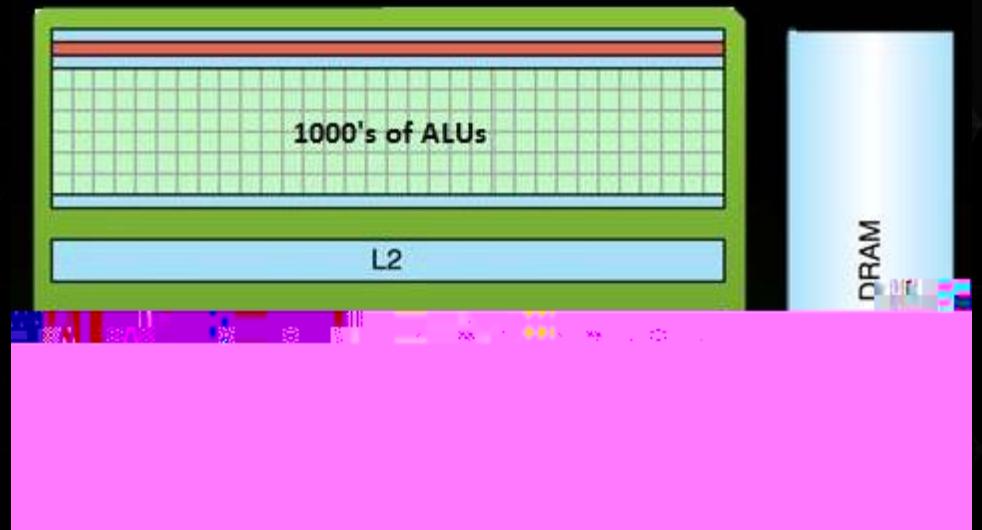
CPU + GPU



CPU + GPU



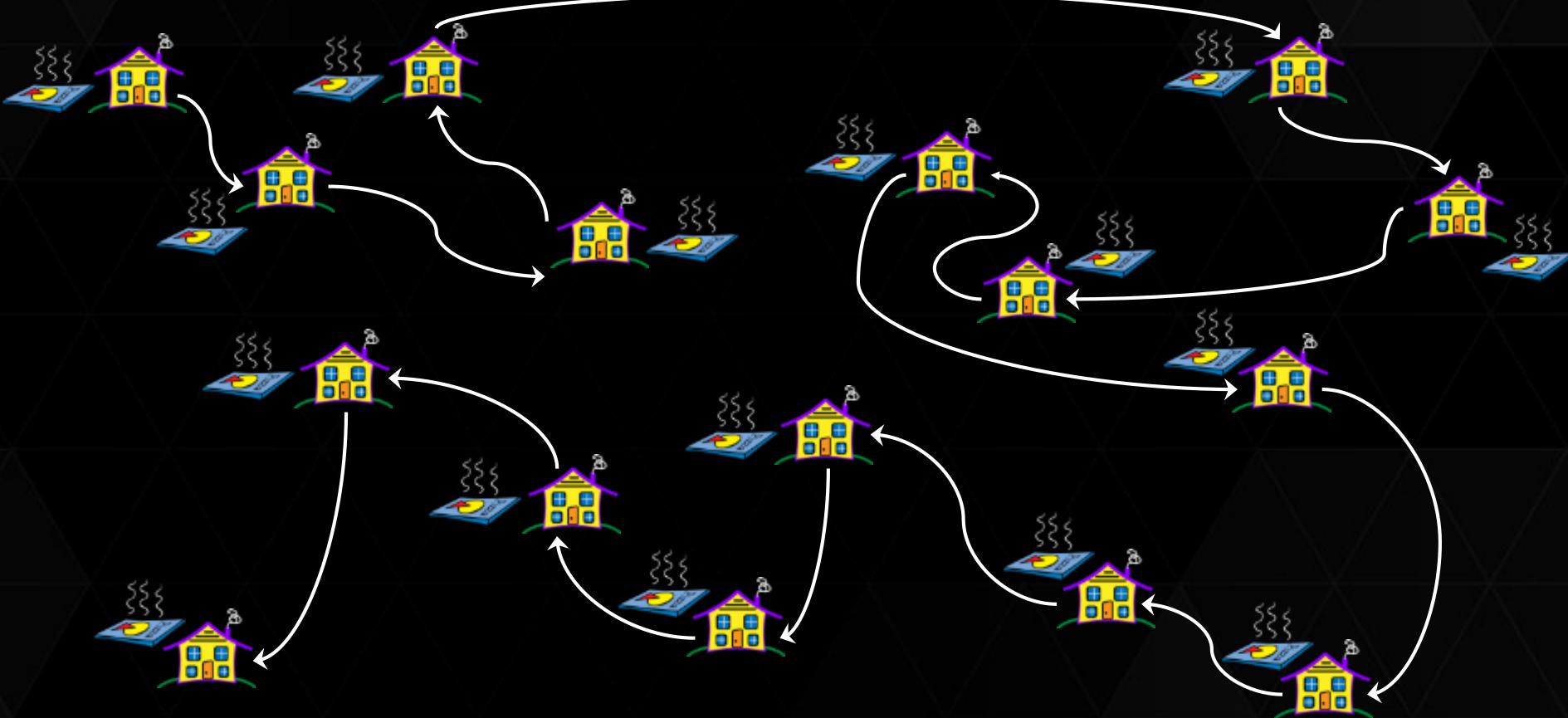
OR



CPU

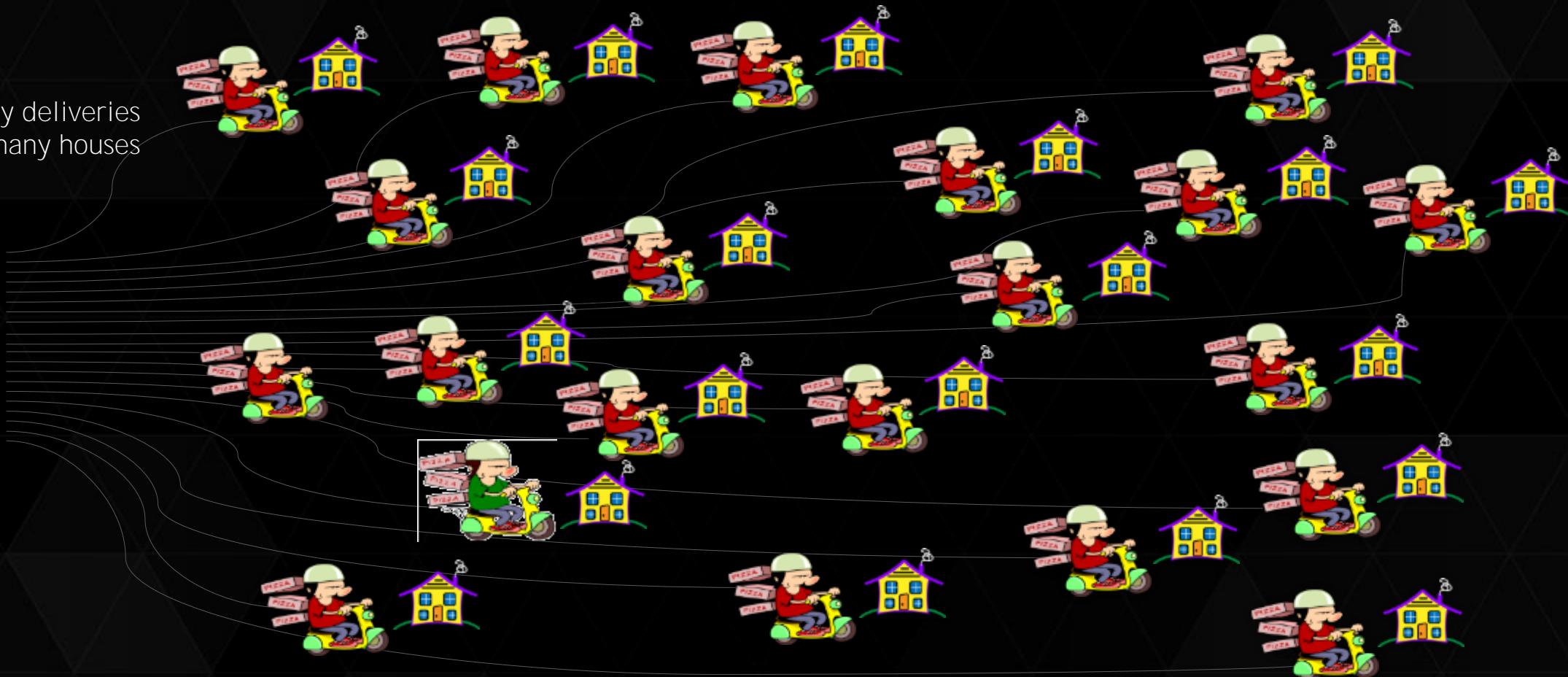


Delivery truck
delivers one pizza
and then moves to
next house

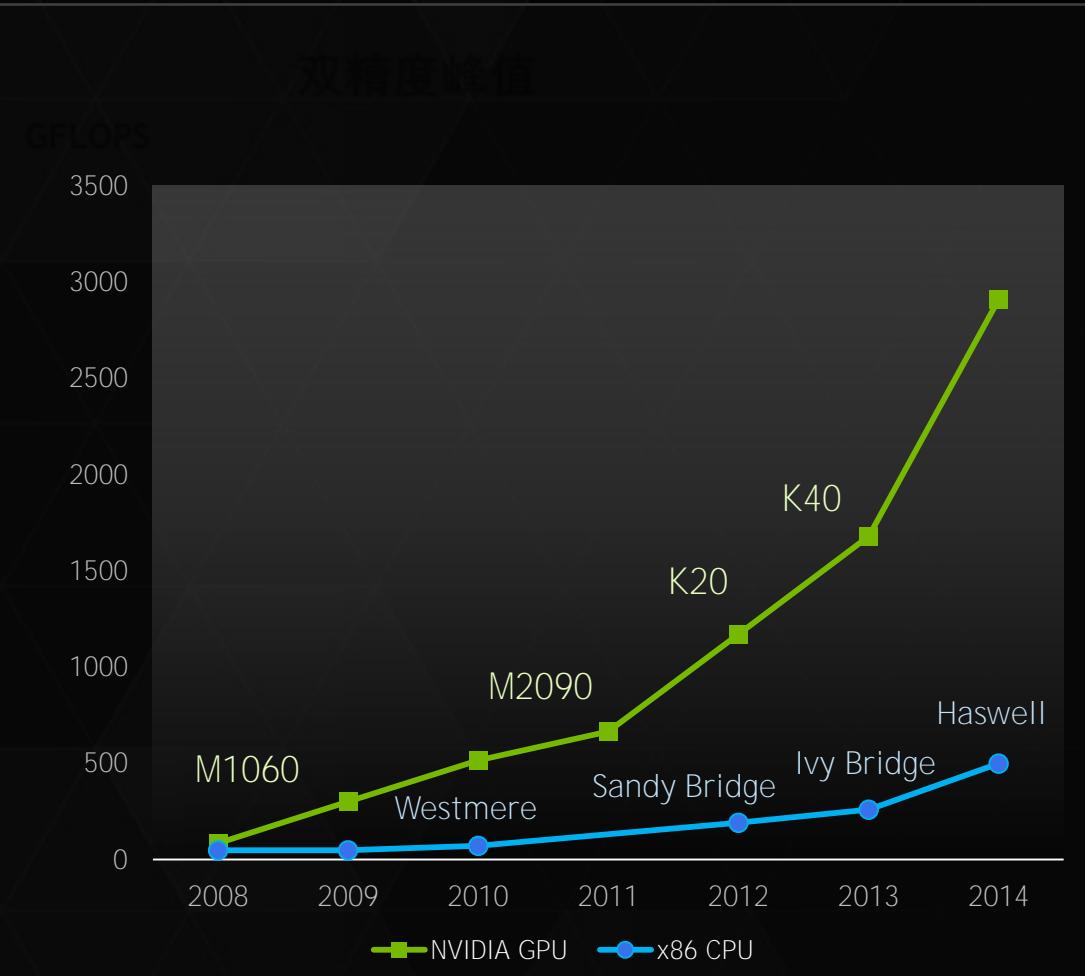


NVIDIA GPU

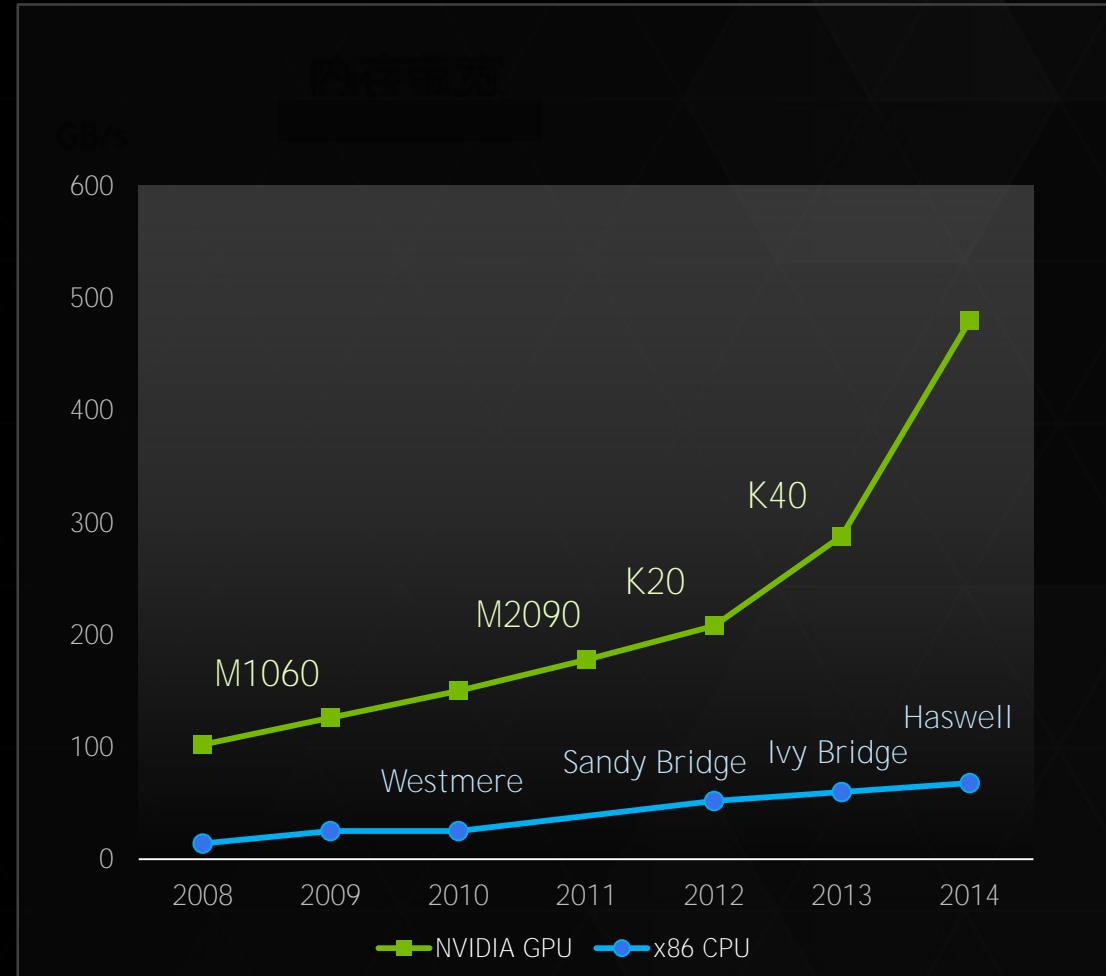
Many deliveries
to many houses



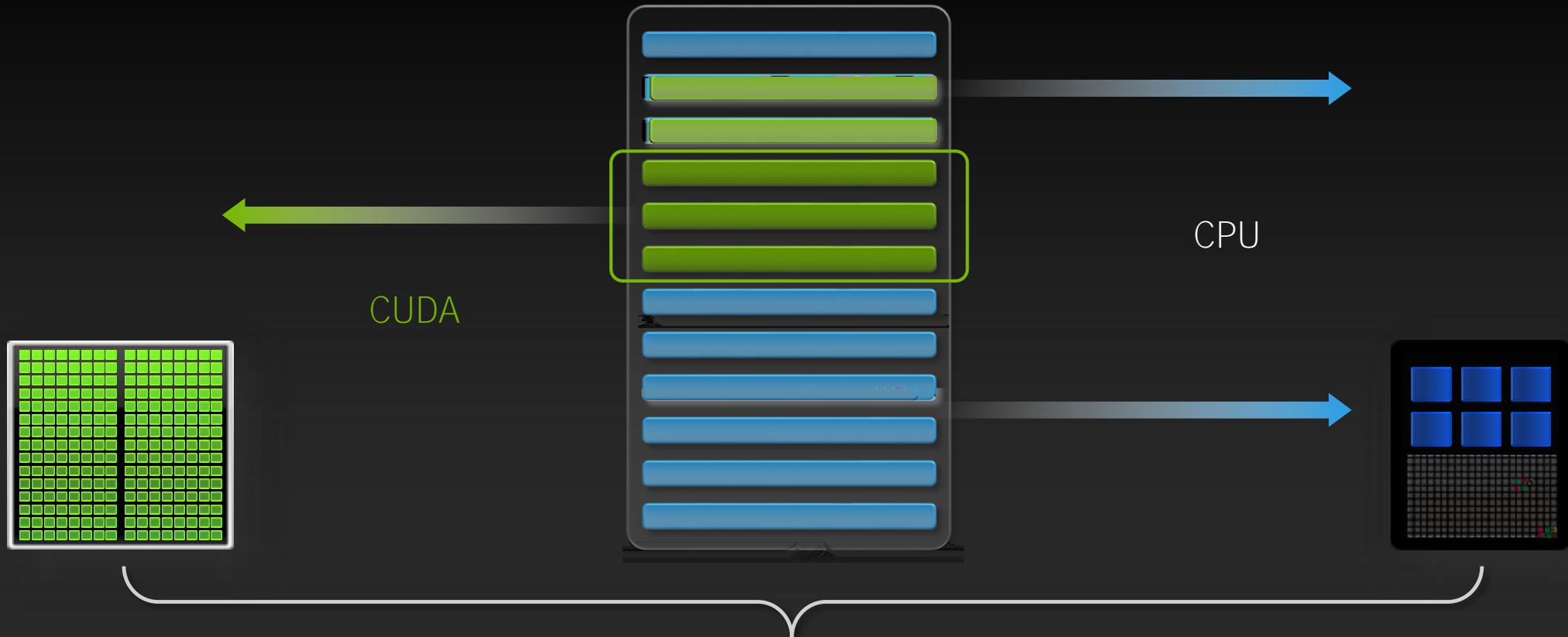
GPU

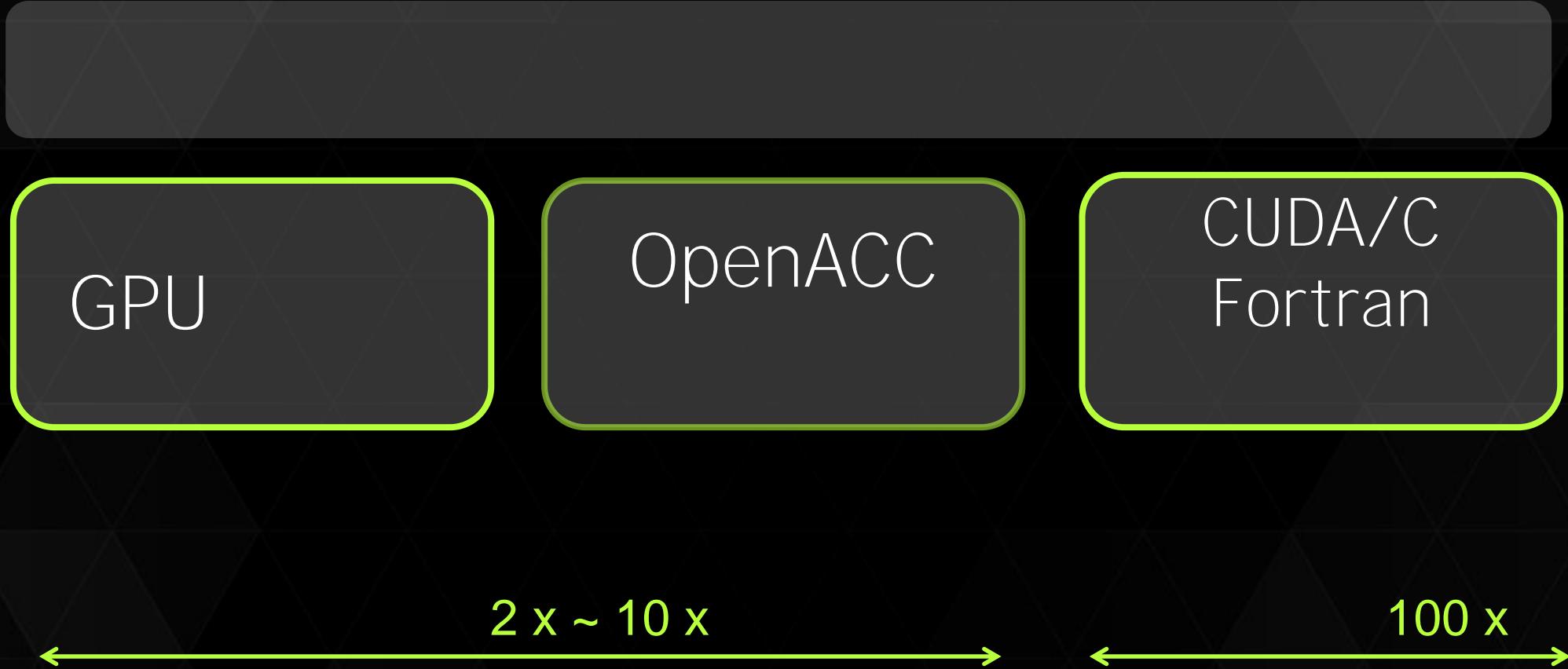


CPU



CPU + GPU =

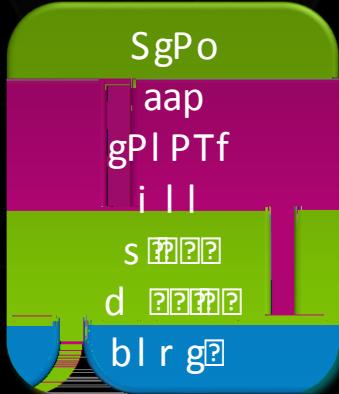
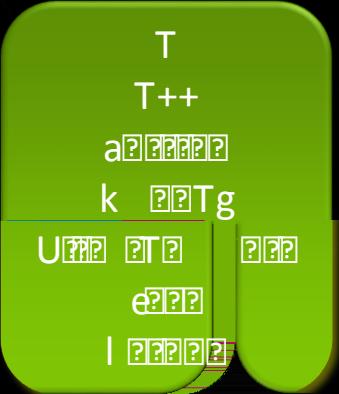
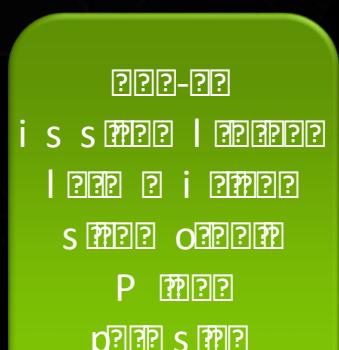




CPU + GPU



NVIDIA GPU



b l b l r T????? ????? & ??????



CUDA

>487,000,000

CUDA GPUs

>2,000,000

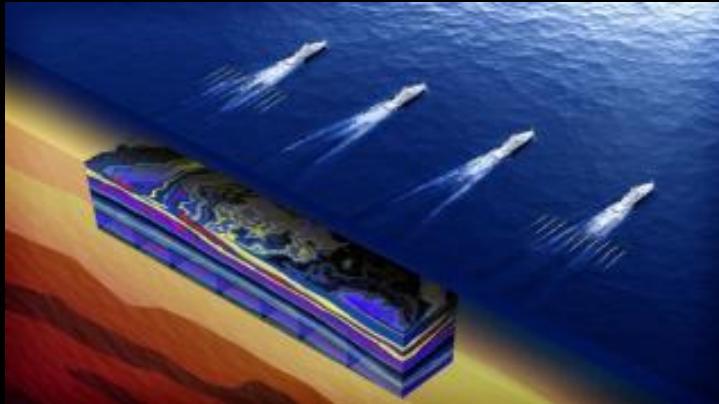
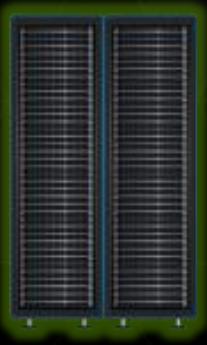
CUDA

>140,000

CUDA

>700

GPU



4x-6x



4x

12%



RTM

GPU



GPU

- ▶ CT
- ▶ X-RAY



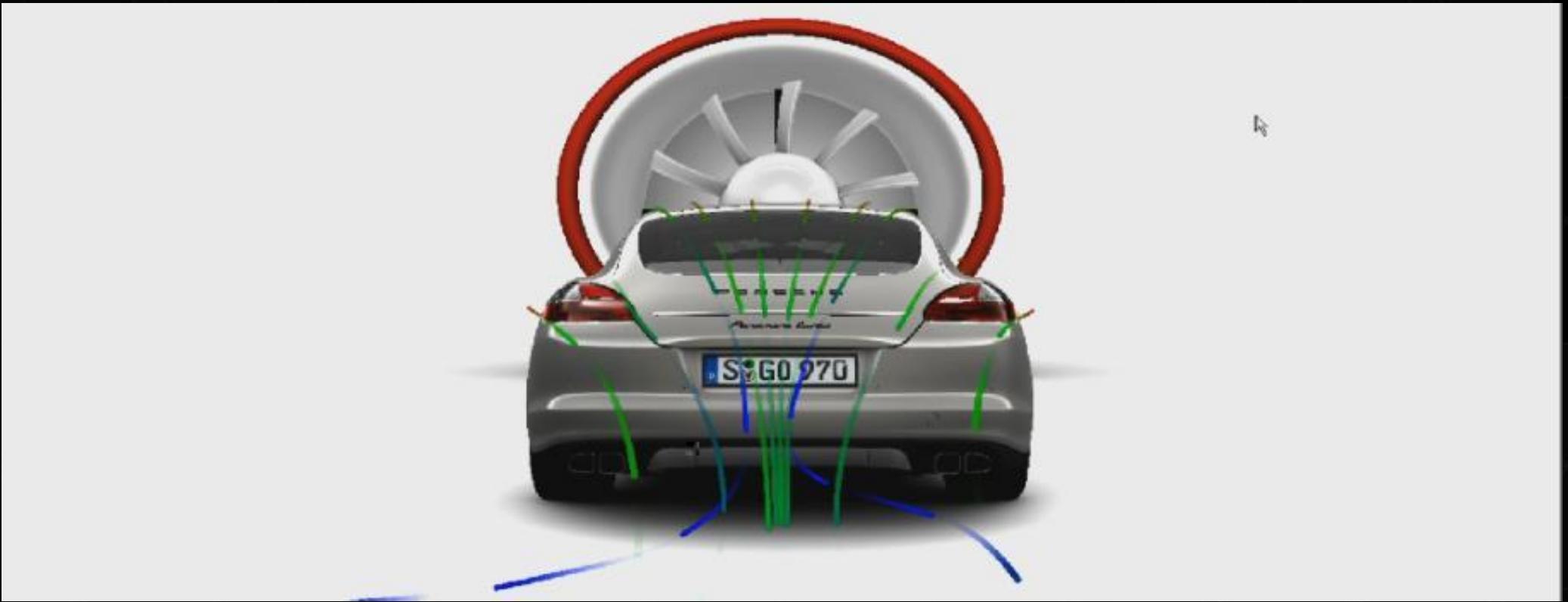
SIEMENS
medical

 GE Healthcare

GPU



GPU



GPU



GPU

/



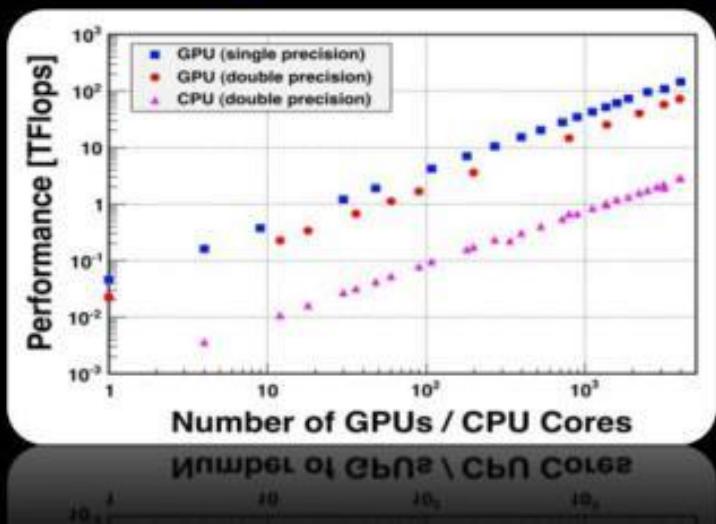
GPU

ASUCA TeraFlop Scaling (Weather Modeling)

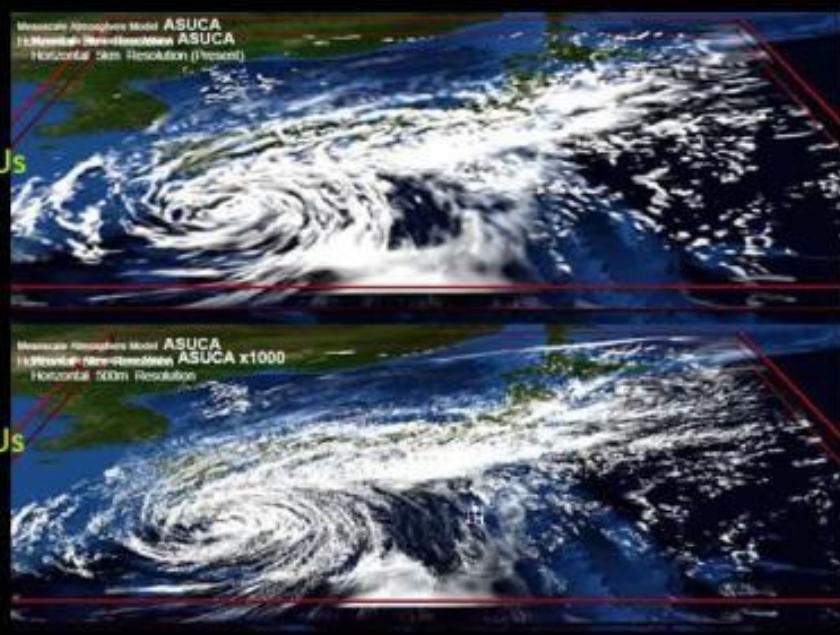
3990 Tesla M2050s

145.0 Tflops SP

76.1 Tflops DP



Before GPUs





- 40x application performance increase
- 80% lower costs
- Calculation results in minutes vs. overnight

Understanding market risk is key to avoiding surprises that impact profitability. J.P. Morgan depends on Tesla GPUs to analyze risk faster than their competitors and to arrive at better decisions through more frequent, more complex calculations. GPUs are



GPU



Typical scale of training data:



Datasets

- Image recognition: 100 millions
- OCR: 100 millions
- Speech: 10 billions
- CTR: 100 billions

Training time:

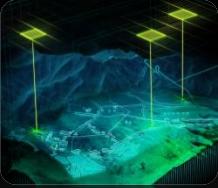
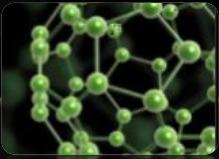
Weeks to Months
on GPU clusters

**Big data + Deep learning + HPC
= Success**

Projected training data to
grow 10x each year



HPC



Schlumberger
Excellence in Energy

BR
PETROBRAS

 **Eni**

 **Chevron**

 **Statoil**

 **HARVARD**
School of Engineering
and Applied Sciences

 **STANFORD**
UNIVERSITY

 **Georgia Tech**

 **ETH**
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

 **UNIVERSITY OF**
CAMBRIDGE

Raytheon

 Air Force
Research
Laboratory

 **NASA**

 Naval Research
Laboratory

 **CSCS**

 **NCSA**

 Tokyo Institute
of Technology

 **OAK RIDGE**
National Laboratory

 **Lawrence Livermore**
National Laboratory

J.P.Morgan

 **BARCLAYS**



 **BNP PARIBAS**

 **MUREX™**

 **Baidu** 百度

 salesforce
SOFTWARE

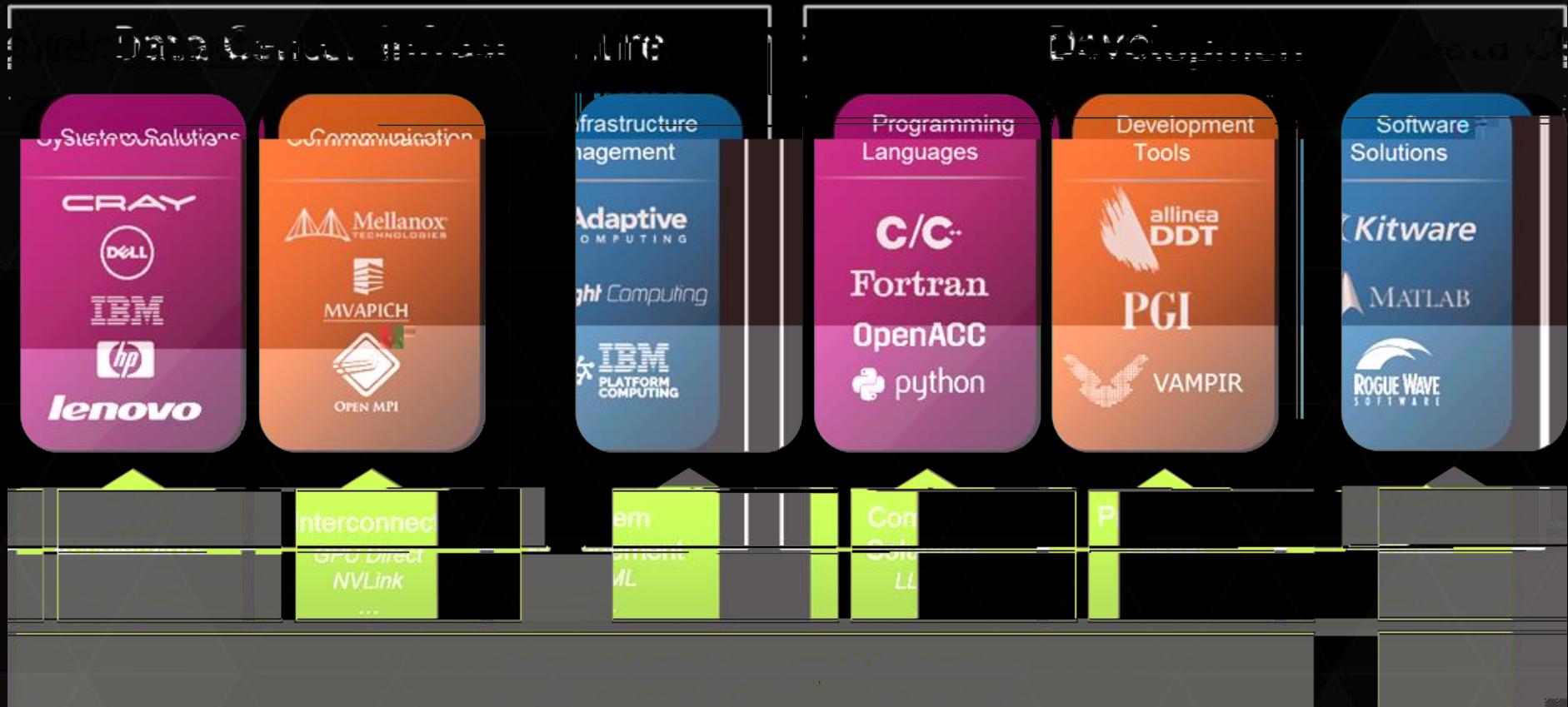
 **SHAZAM®**

 **amazon.com**

 **Yandex**

 **NVIDIA**

Tesla



2014,

%

Libraries



OpenACC

Programming
Languages



2015 Tesla GPU



Tesla K40

Best Single GPU Performance

Server, Workstation, Liquid Cooled

Higher Ed, Data Analytics, HPC Labs, Defense

Double Precision Workloads



Tesla K80

Maximize Throughput within a Server

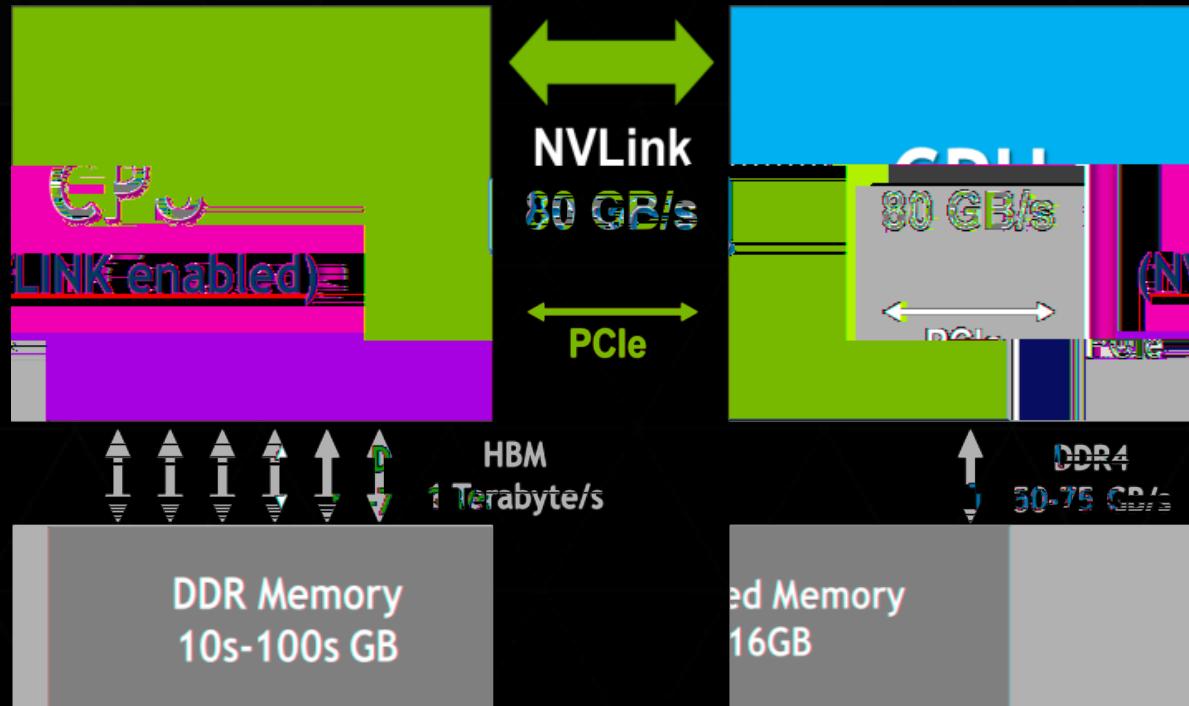
Seismic, Data Analytics, HPC Labs, Defense

Multi-GPU Accelerated Apps

Single and Double Precision Workloads

NVLINK

CPU GPU

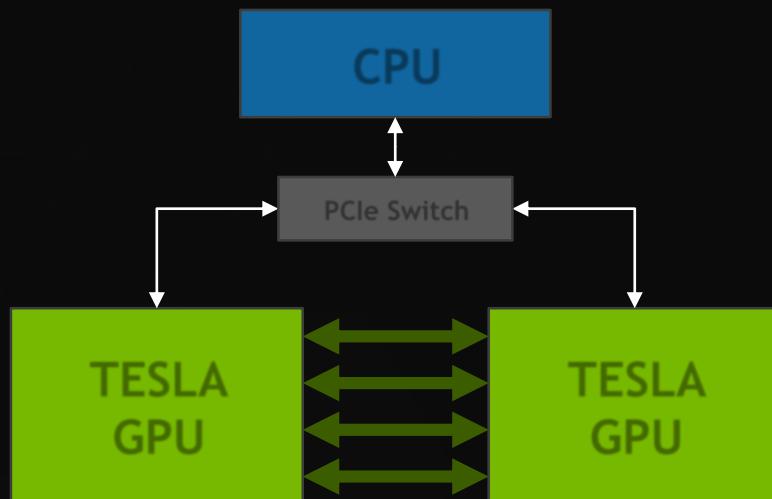


NVLINK 1.0
4 data links with 20GB/s each

X86

GPU NVLINK

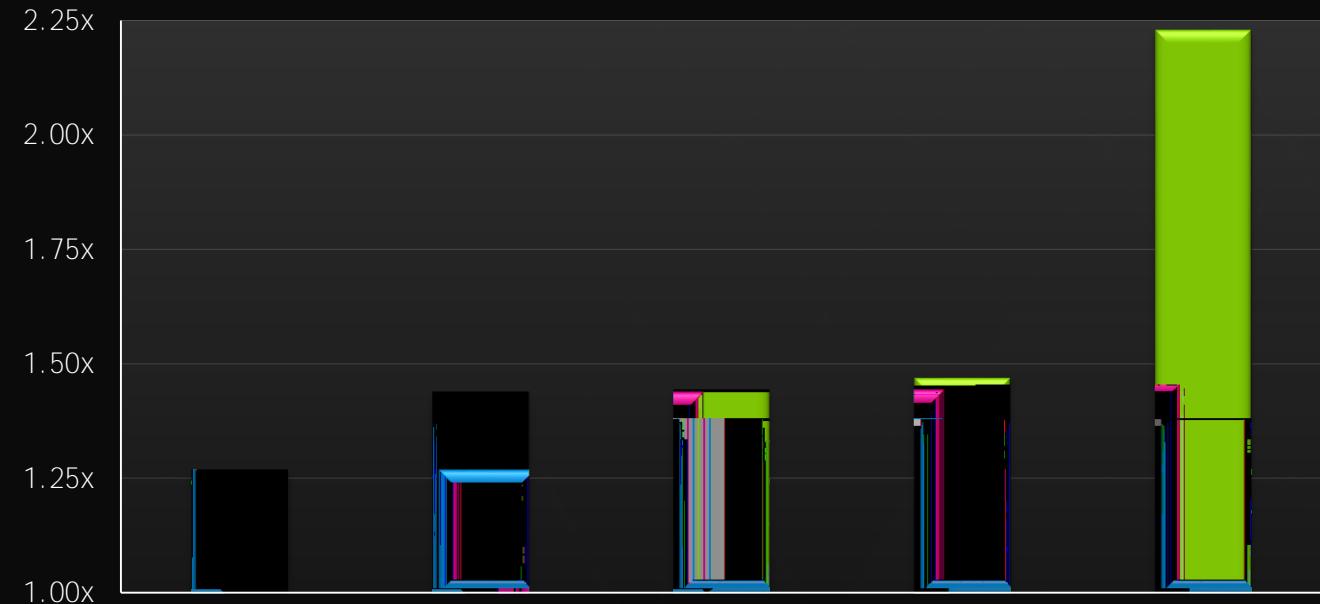
GPUs Interconnected with NVLink



5x Faster than
PCIe Gen3 x16

Over 2x Application Performance Speedup
When Next-Gen GPUs Connect via NVLink Versus PCIe

Speedup vs
PCIe based Server



To learn more: <http://www.nvidia.com/object/nvlink.html>

US to Build Two Flagship Supercomputers Powered by the Tesla Platform



100-300 PFLOPS Peak

10x in Scientific App Performance

IBM POWER9 CPU + NVIDIA Volta GPU

NVLink High Speed Interconnect

40 TFLOPS per Node, >3,400 Nodes

2017

Major Step Forward on the Path to Exascale



"Machine Learning" is in some sense a rebranding of AI.

The focus is now on more specific, often perceptual tasks, and there are many successes.

Today, some of the world's largest internet companies, as well as the foremost research institutions, are using GPUs for machine learning.



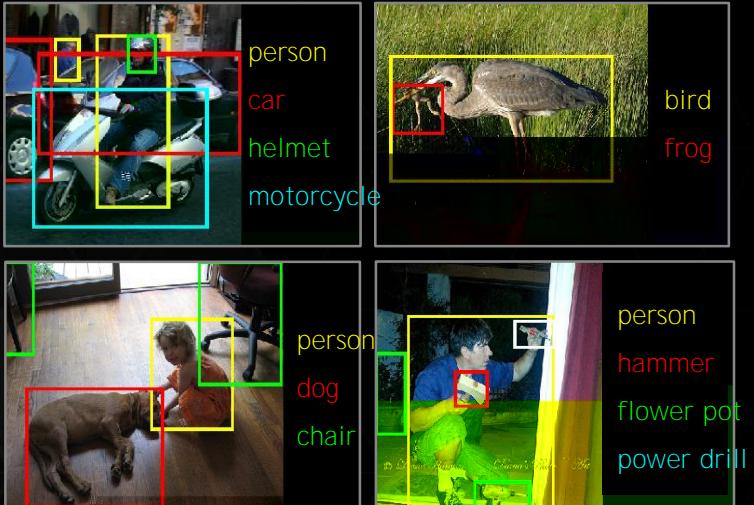
GPU -

1.2

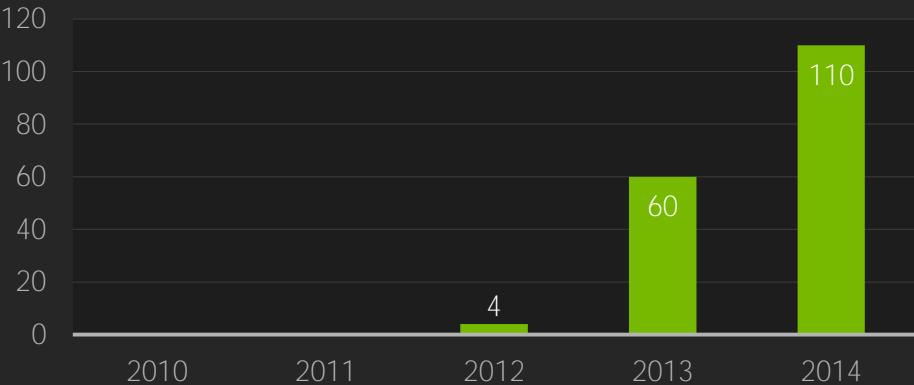
1000

Hosted by

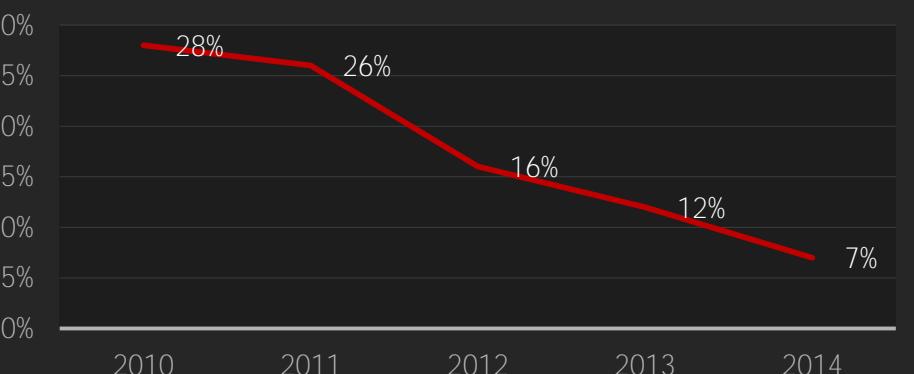
IMAGENET



GPU Entries



Classification Error Rates

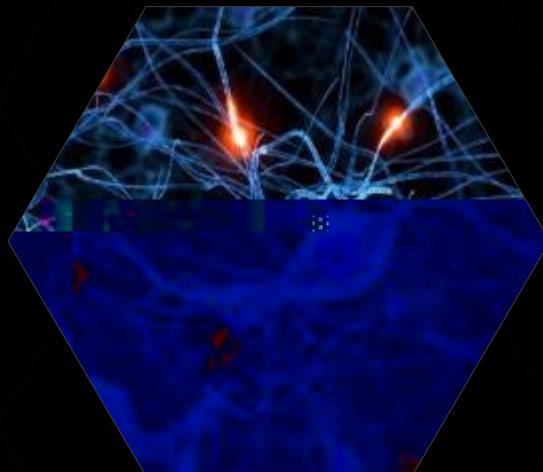


3

大数据



模型



GPU加速器



GPU

Early Adopters



Image Analytics
for Creative
Cloud



Speech/Image
Recognition



Image
Classification



Hadoop



Recommendation



Search Rankings

Use Cases

&

&

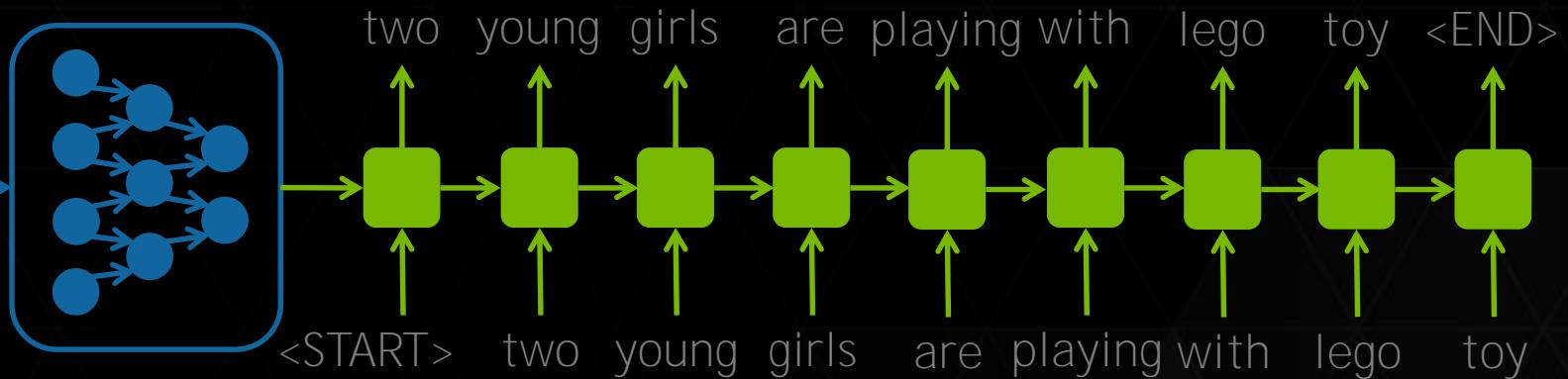
&

Talks @ GTC



ConvNets + RNN/LSTM

- ▶ Multiple papers (Stanford, Google, UC Berkley, and others) in Nov 14



Car and Lane detection via CNN from Baidu, Stanford, Twitter, TI

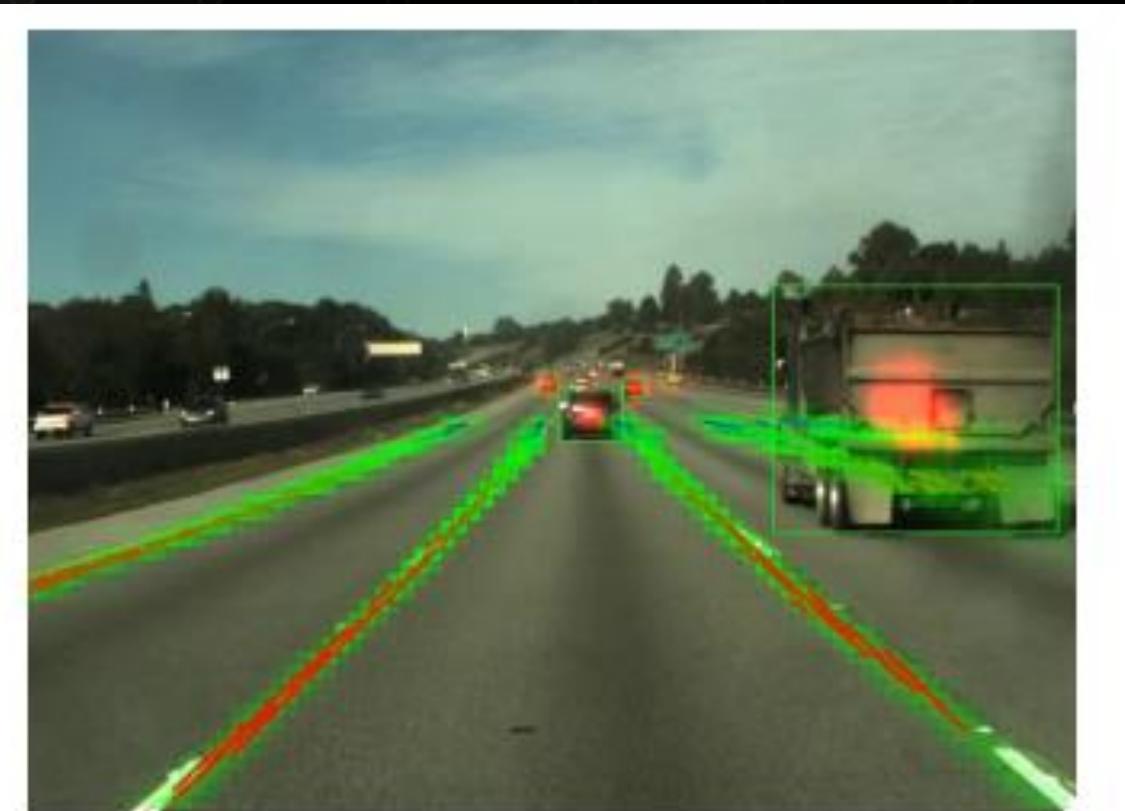


Fig. 1: Sample output from our neural network capable of lane and vehicle detection.

FUELING THE DEEP LEARNING REVOLUTION

March 17 – 20, 2015 | Silicon Valley | #GTC15

REGISTER NOW

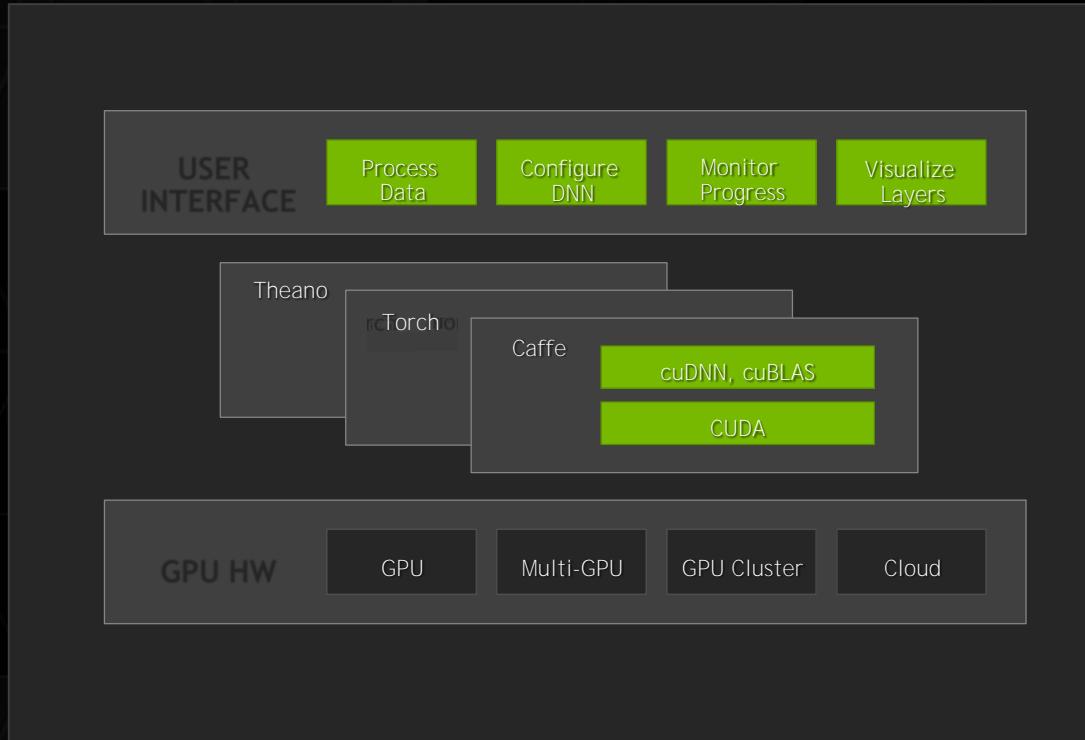
50+
Deep Learning Sessions

Adobe	Google
Alibaba	iFlytek, Ltd
Baidu	NUANCE
Carnegie Mellon	Stanford Univ
Facebook	UC Berkeley
Flickr / Yahoo	Univ of Toronto

Developer Labs

Caffe
Torch
Theano

DiGITS - GPU



DEEP GPU TRAINING SYSTEM
FOR DATA SCIENTISTS

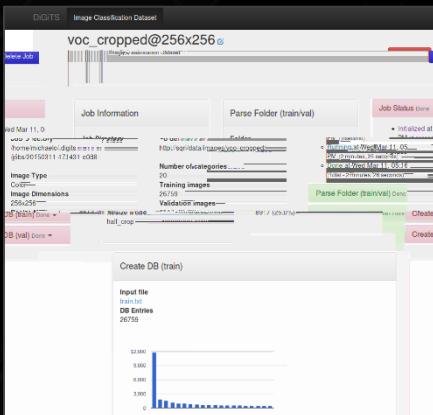
Design DNNs

Visualize activations

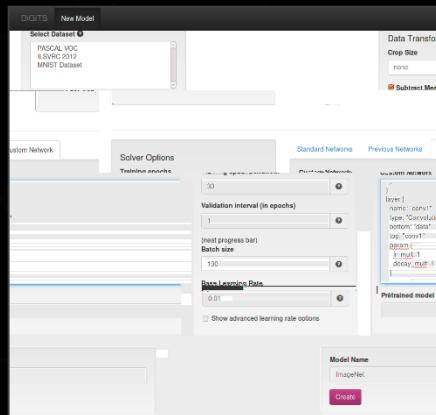
Manage multiple trainings

DIGITS

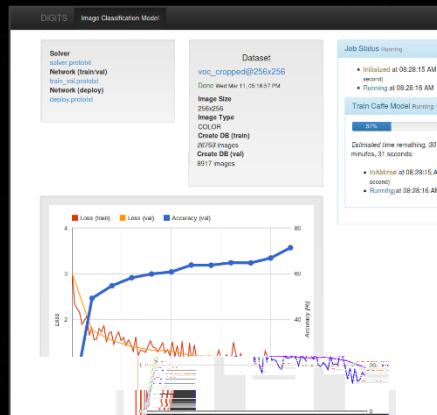
Process Data



Configure DNN



Monitor Progress



Visualize Layers



HLUO@NVIDIA.COM