

**GTCx**

SEOUL | Oct.7, 2016

# Architecture & Programming for A Tightly Bound CPU-GPU World

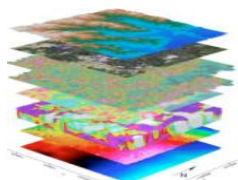
허욱 실장, Linux Server Solution Sales, IBM Korea

PRESENTED BY



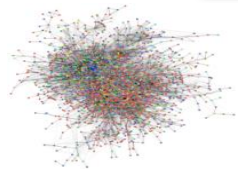
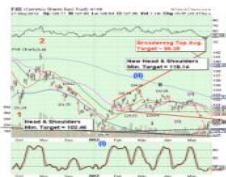
# Computing Acceleration with GPU

HPC & HPA needs new computing platform



Scientific Simulation


Financial Analysis

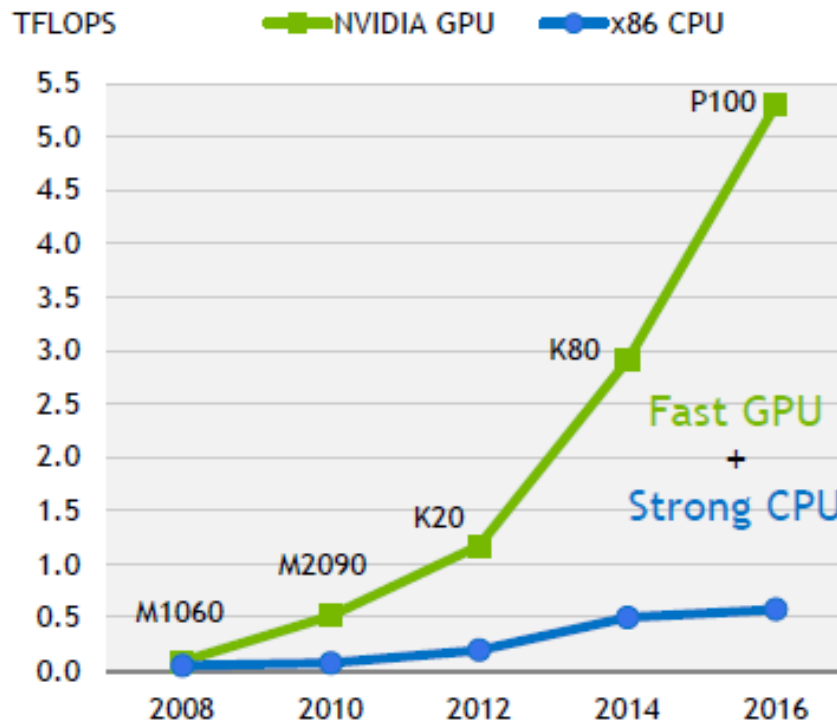


Graph Analytics

Deep Learning



More Data Volume  
  
 More Computing Power

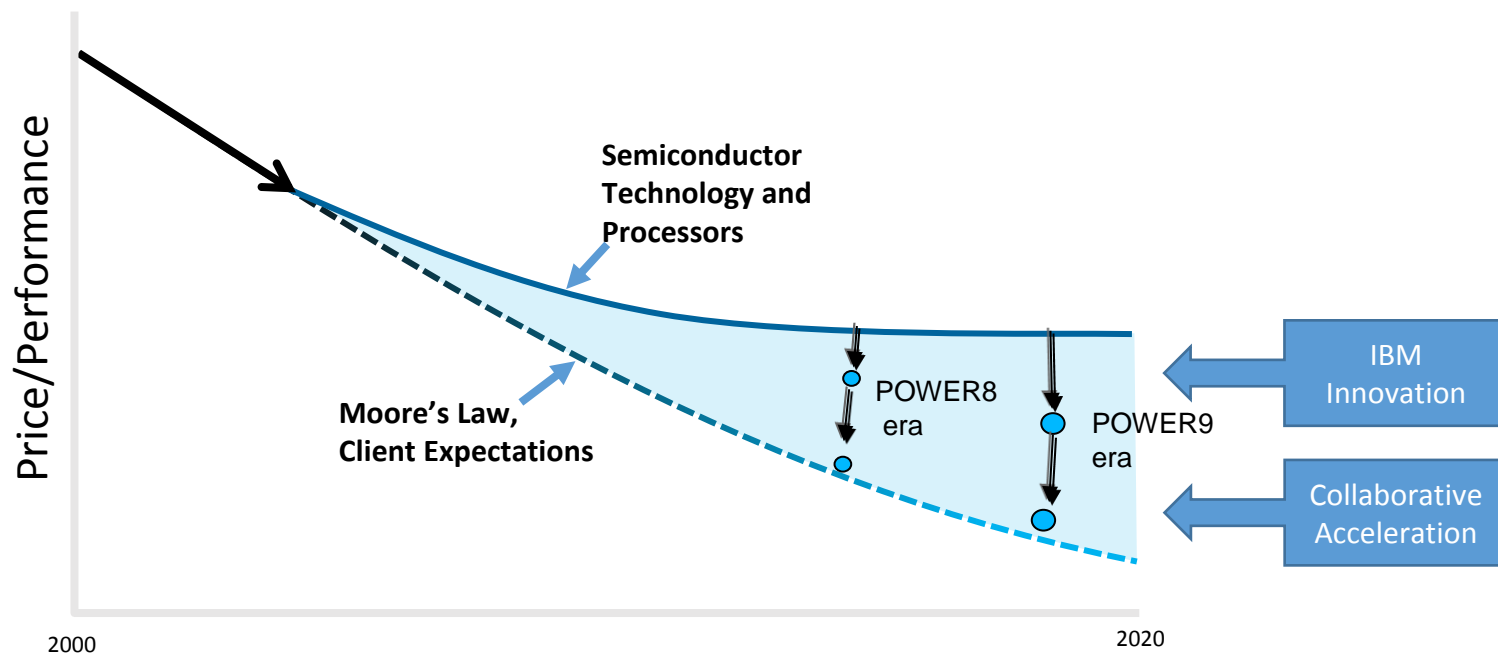


Manual Data Movement

Limits of PCI-E bus

# Innovation thru true Open Collaboration

## OpenPOWER: Open Architecture for HPC & HPA



**Open Source.  
Open Collaboration.  
Open the Throttle.**

Join the OpenPOWER Developer Challenge.



Processor IP  
Licensing

Open  
Interfaces

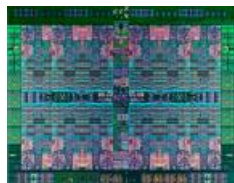
Systems  
& Software

# 1<sup>st</sup> Breakthrough in GPU computing

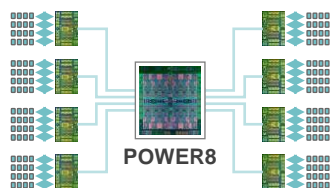
Combination of best of breed technologies thru open collaboration

POWER8 w/ NVLink

Tesla P100



Faster Cores



POWER8

More Cache size & Memory Bandwidth



Faster Data Communication

+



# IBM Power S822LC for HPC: ‘Minsky’

## 1<sup>st</sup> Custom-built GPU Accelerator Server w/ NVLink



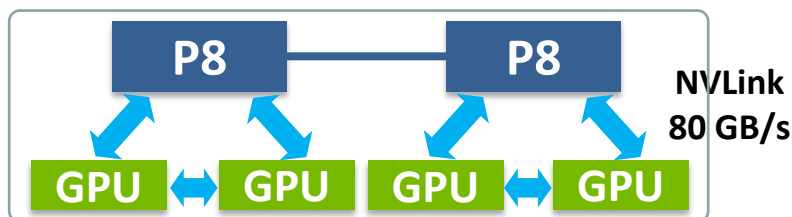
- 2U Form factor, 2 \* POWER8 & 2 or 4 \* P100, Up to 1TB mem
- High-Speed NVLink Connections between CPUs & GPUs and among GPUs
- 1<sup>st</sup> available system w/ P100 GPUs

**CORAL**  
COLLABORATION  
OAK RIDGE • ARGONNE • LIVERMORE

OAK  
RIDGE  
National Laboratory

Lawrence Livermore  
National Laboratory

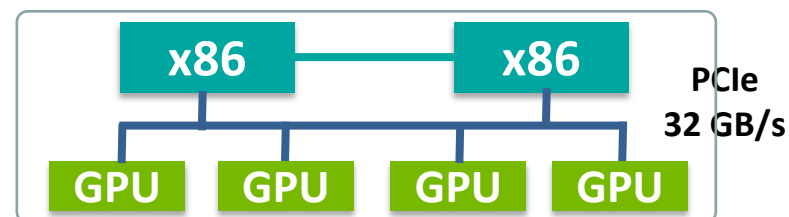
2.5x Faster CPU-GPU Data  
Communication **via NVLink**



“Minsky”

POWER8 NVLink Server

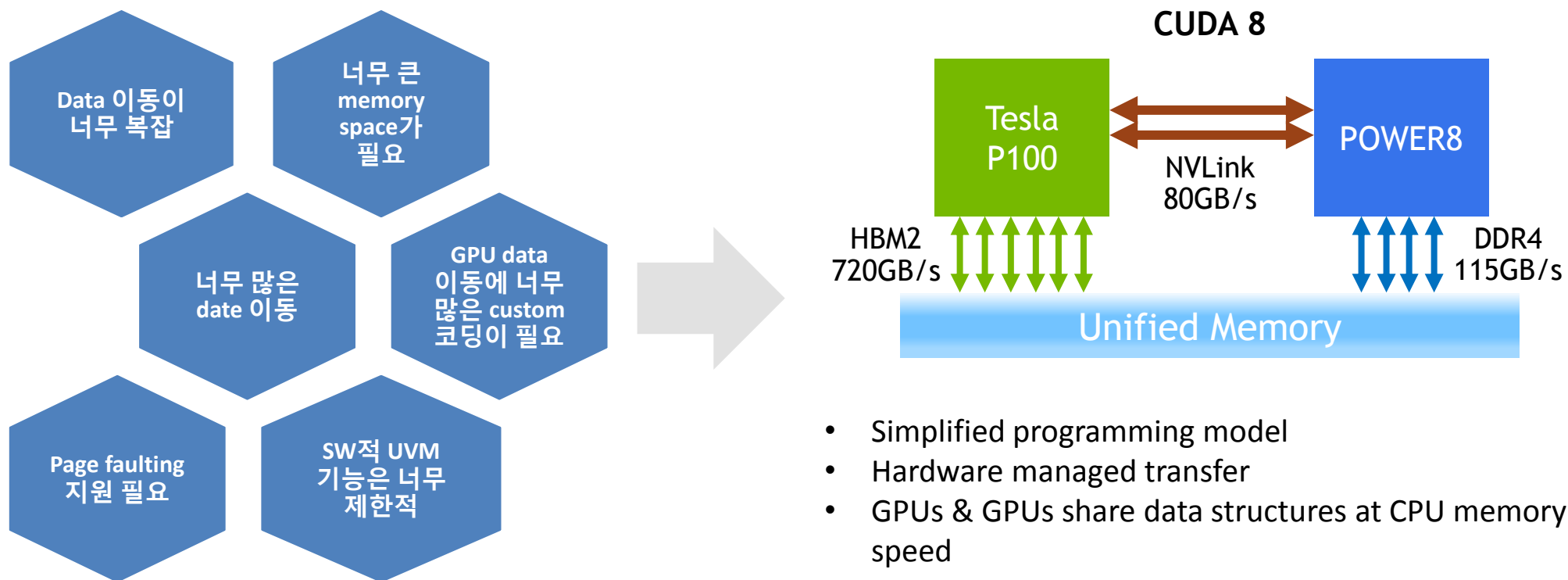
**No NVLink** between CPU & GPU for  
x86 Servers: PCIe Bottleneck



x86 Servers with PCIe

# Better Developer Productivity

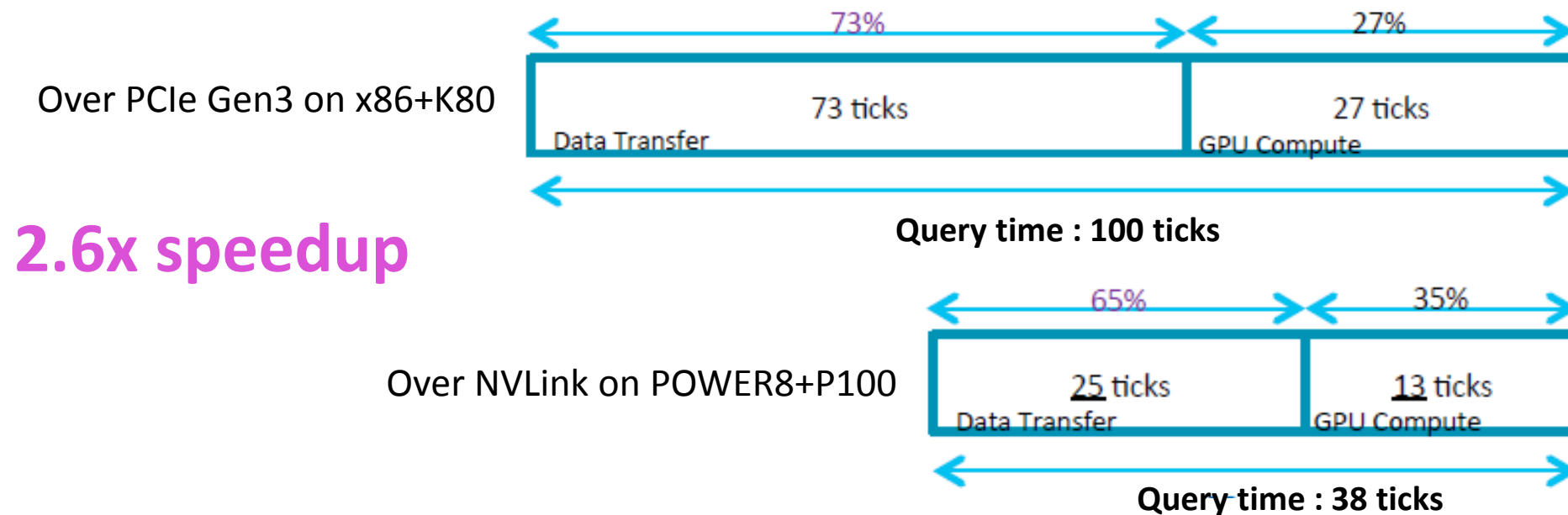
## NVLink w/ Unified Memory



- Simplified programming model
- Hardware managed transfer
- GPUs & GPUs share data structures at CPU memory speed

# Benefit of NVLink between CPU & GPU

Real application example: Kinetica GPU DB

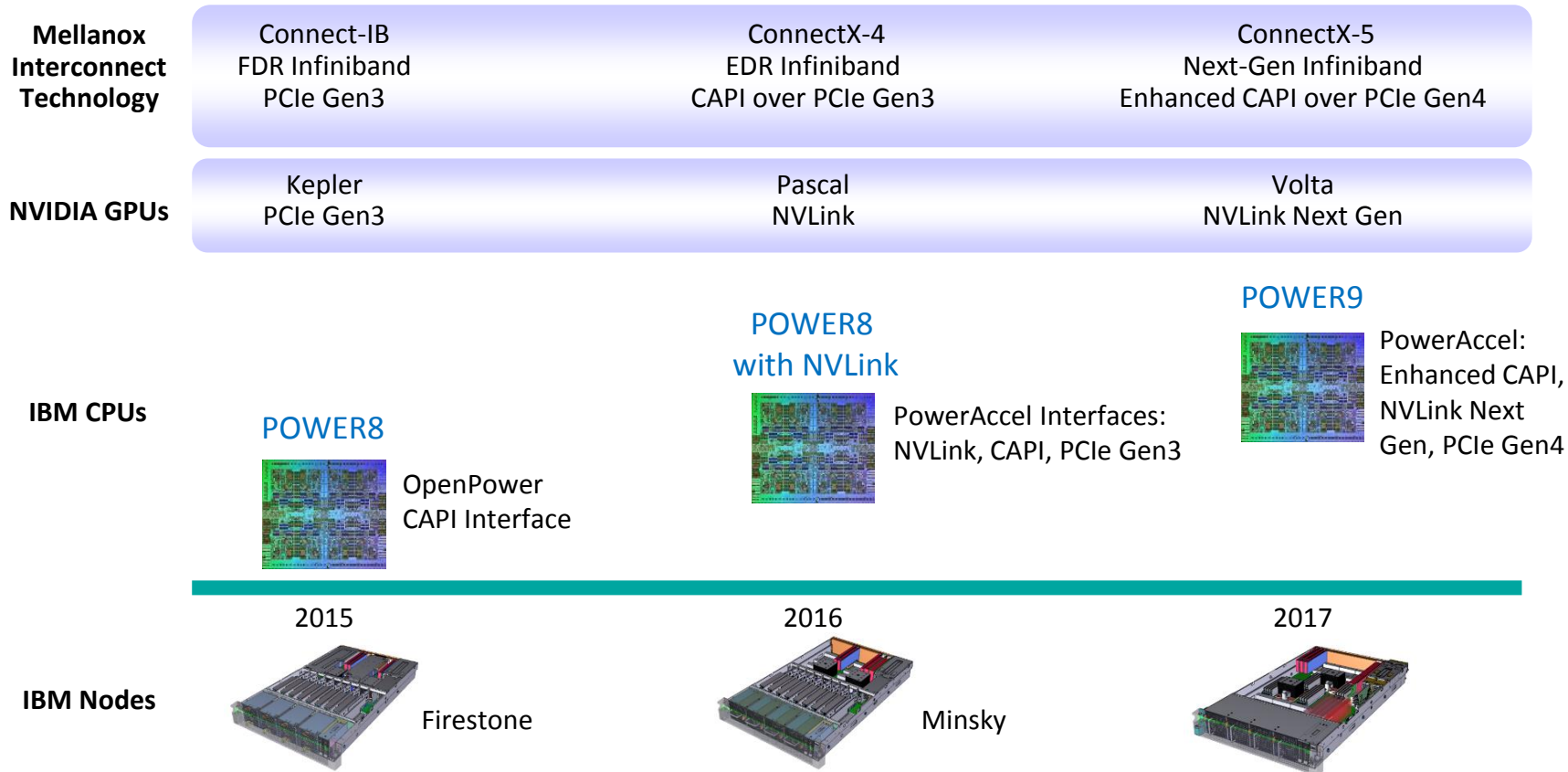


Total Reduction: 62 ticks

- Data Transfer Reduction: 48 ticks → **77%** contribution
- GPU Compute Reduction: 14 ticks → **23%** contribution

# Evolution w/ long-term collaboration

## Future Roadmap for the best GPU accelerator system





**GTCx**

SEOUL | Oct.7, 2016

# THANK YOU

JOIN THE CONVERSATION

**#GTCxKorea2016** **f**

PRESENTED BY

