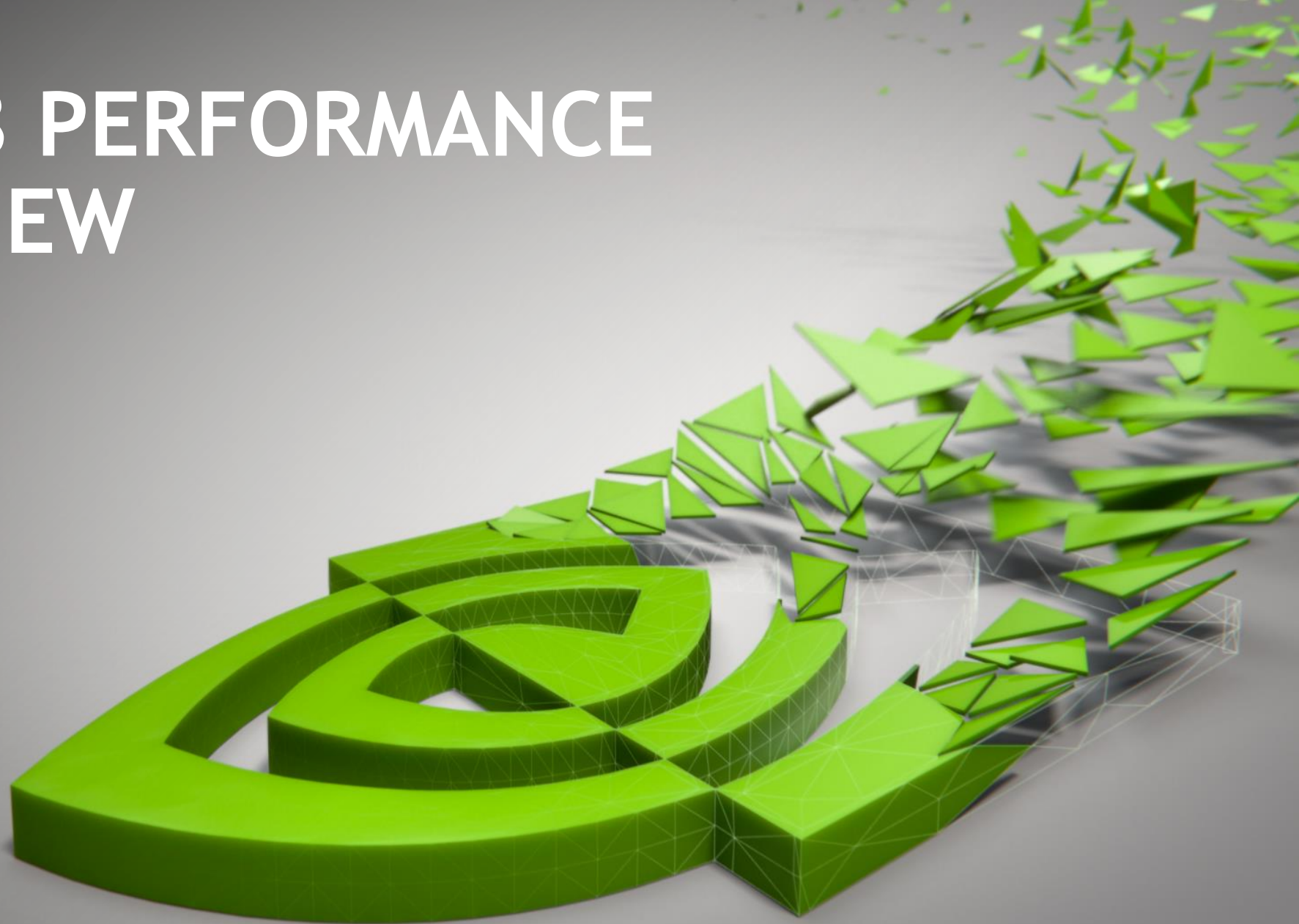# CUDA 8 PERFORMANCE OVERVIEW
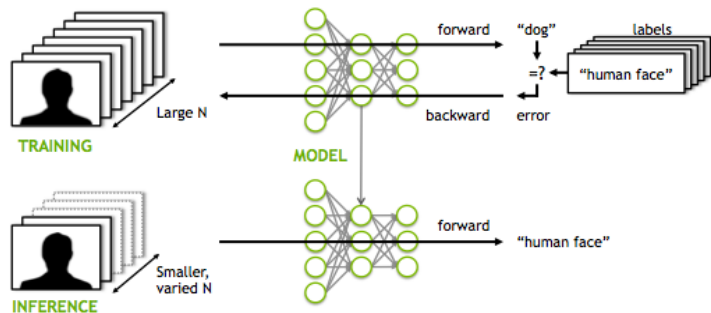
November 2016
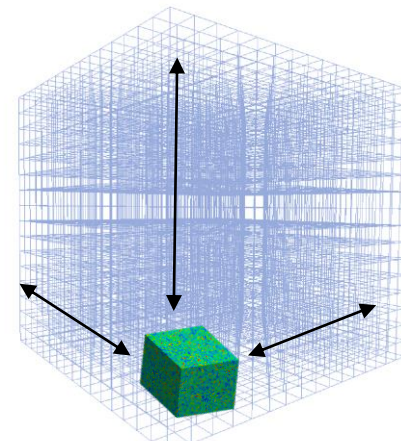
# CUDA 8 PERFORMANCE HIGHLIGHTS



1.5-2X higher performance out-of-the-box



Solve larger problems than possible before with Unified Memory



Mixed precision support (FP16, INT8)
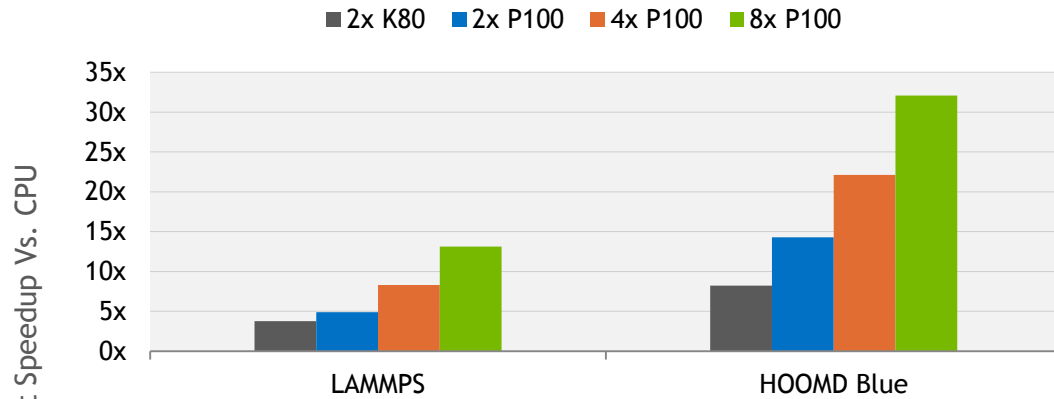
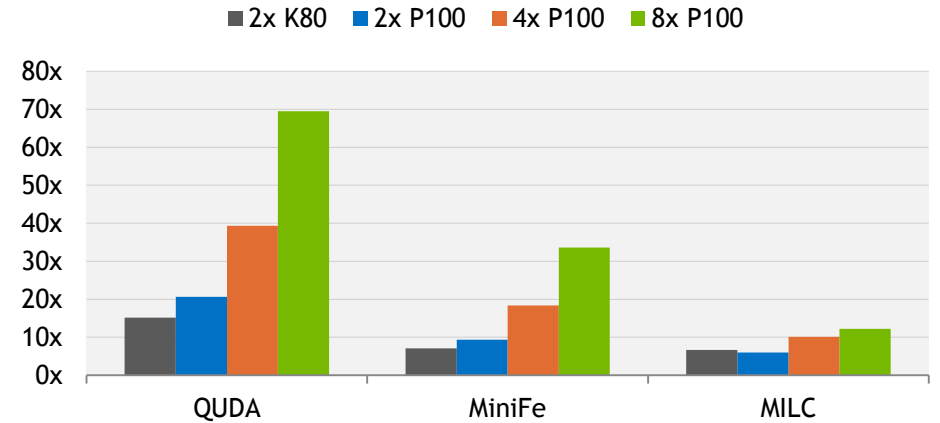**SOCIAL NETWORK ANALYSIS**



Wikimedia Commons

nvGRAPH: New graph analytics library

# HIGH PERFORMANCE COMPUTING:
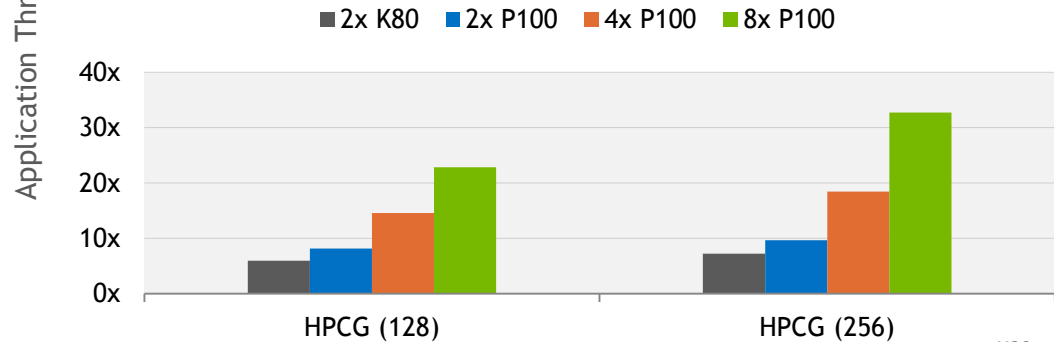## APPS WITH CUDA 8, OpenACC

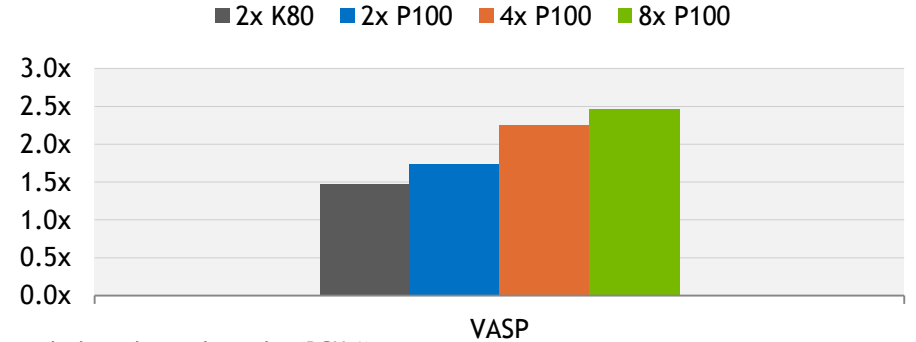# > 2X SPEEDUPS ON HPC APPS WITH P100

## Molecular Dynamics

■ 2x K80 ■ 2x P100 ■ 4x P100 ■ 8x P100

Application Throughput Speedup Vs. CPU

Chart showing LAMMPS and HOOMD Blue benchmarks, y-axis 0x to 35x.

## Physics

■ 2x K80 ■ 2x P100 ■ 4x P100 ■ 8x P100

Chart showing QUDA, MiniFe, MILC benchmarks, y-axis 0x to 80x.

## HPC Benchmarks

■ 2x K80 ■ 2x P100 ■ 4x P100 ■ 8x P100

Chart showing HPCG (128) and HPCG (256) benchmarks, y-axis 0x to 40x.

## Quantum Chemistry

■ 2x K80 ■ 2x P100 ■ 4x P100 ■ 8x P100

Chart showing VASP benchmark, y-axis 0.0x to 3.0x.
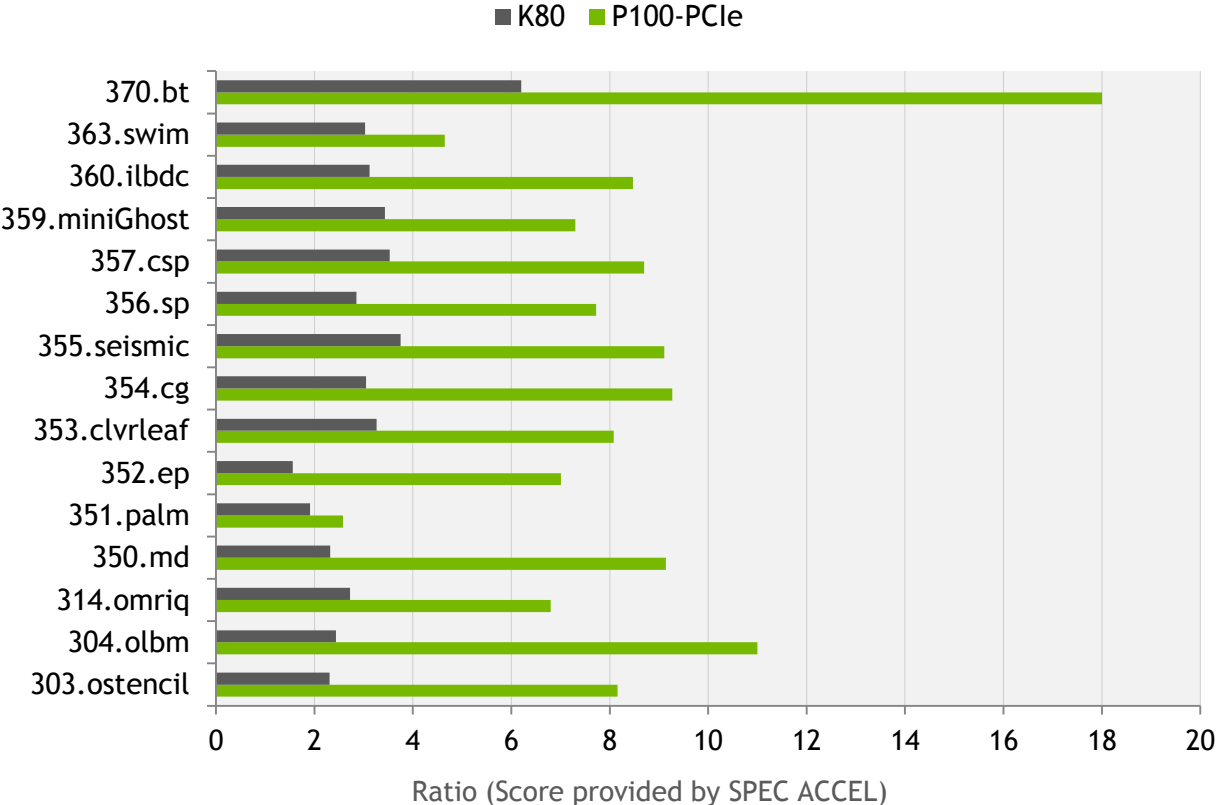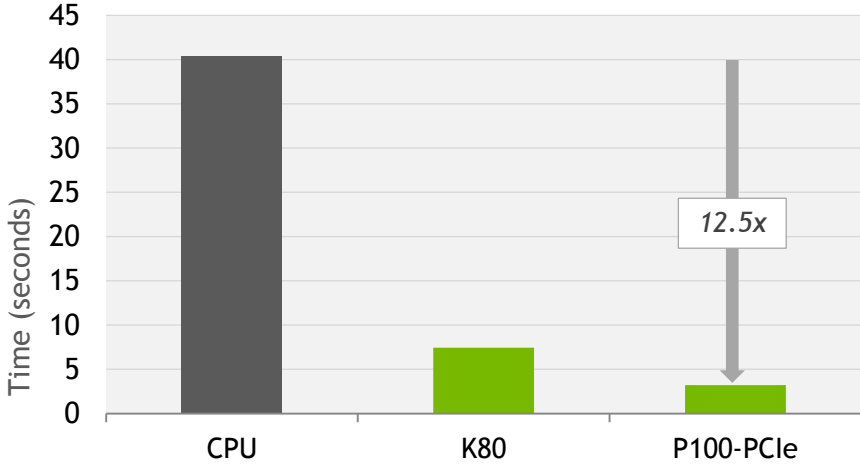
- K80 and P100 (SXM2); Base clocks; cube mesh topolgy (DGX-1)
- CUDA 8 GA (8.0.44) with r361.79 (K80) and r361.96 (P100)
- CPU System: Intel Xeon Broadwell dual socket 22-core E5-2699 v4@2.2GHz 3.6GHz Turbo with Ubuntu 14.04.3 x86-64 and 256GB memory
- Full system configurations including benchmark versions and data sets used available in the Appendix

⊚ nVIDIA.

# OpenACC: >2X—12X FASTER WITH P100

## 2.7x Speedup on SPEC ACCEL

Legend: ■ K80  ■ P100-PCIe

Categories (top to bottom):
- 370.bt
- 363.swim
- 360.ilbdc
- 359.miniGhost
- 357.csp
- 356.sp
- 355.seismic
- 354.cg
- 353.clvrleaf
- 352.ep
- 351.palm
- 350.md
- 314.omriq
- 304.olbm
- 303.ostencil

X-axis: Ratio (Score provided by SPEC ACCEL) — 0 to 20

## 12.5x Faster than CPU on CloverLeaf

Y-axis: Time (seconds) — 0 to 45

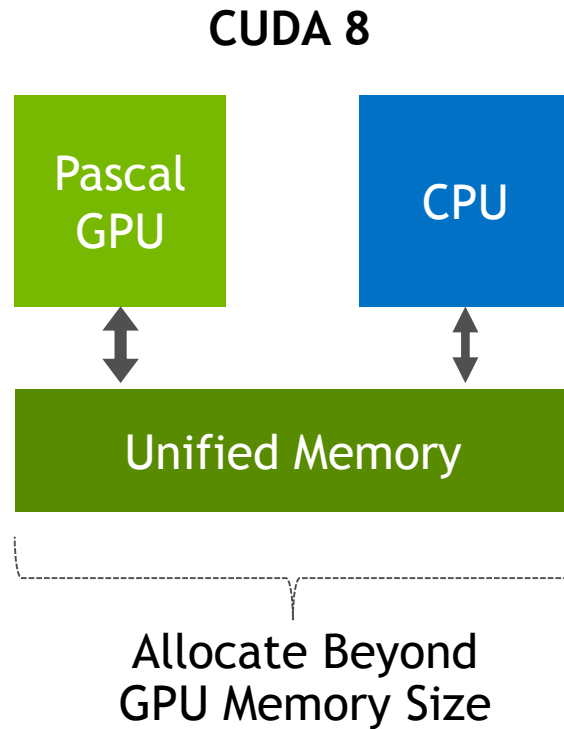Bars: CPU (~40), K80 (~7.5), P100-PCIe (~3), annotated 12.5x

- SPEC ACCEL Base runs on 1x K80 (1-GPU), 1x P100 (PCIe) ; Base clocks and ECC ON
- CloverLeaf runs on 1x K80 (2-GPUs), 2x P100 (PCIe) ; Base clocks and ECC ON
- PGI 16.1 with CUDA 7.5 on K80 (r352) and PGI 16.7 with CUDA 8 on P100 (r361)
- CPU measurements for CloverLeaf on Intel Xeon Broadwell single-socket 20-core E5-2699 v4@2.2GHz 3.6GHz Turbo with Intel MKL 2017
- MPI versions 1.8.8 (K80), 1.10.2 (P100) and 2017.0.098 (CPU)
- Host System: Intel Xeon Haswell dual socket 16-core/socket E5-2698 v3@2.3GHz 3.6GHz Turbo with Ubuntu 14.04.3 x86-64 and 256GB memory

Performance may vary based on OS and software versions, and motherboard configuration

5 NVIDIA.

# UNIFIED MEMORY ON PASCAL

## Large datasets, Simple programming, High performance

**CUDA 8**

Pascal GPU

CPU

Unified Memory

Allocate Beyond
GPU Memory Size

**Enable Large Data Models**
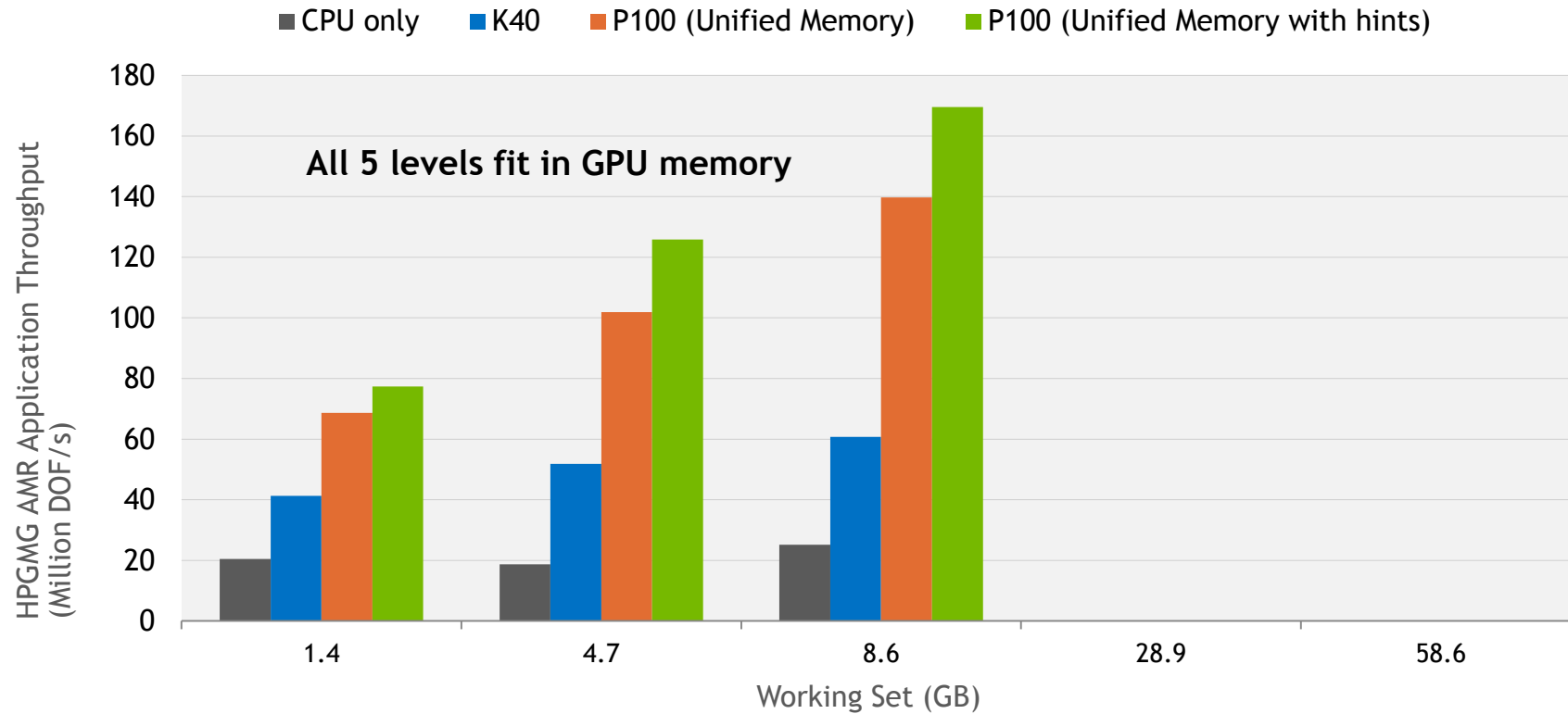- Oversubscribe GPU memory
- Allocate up to system memory size

**Higher Application Performance**
- Demand paging
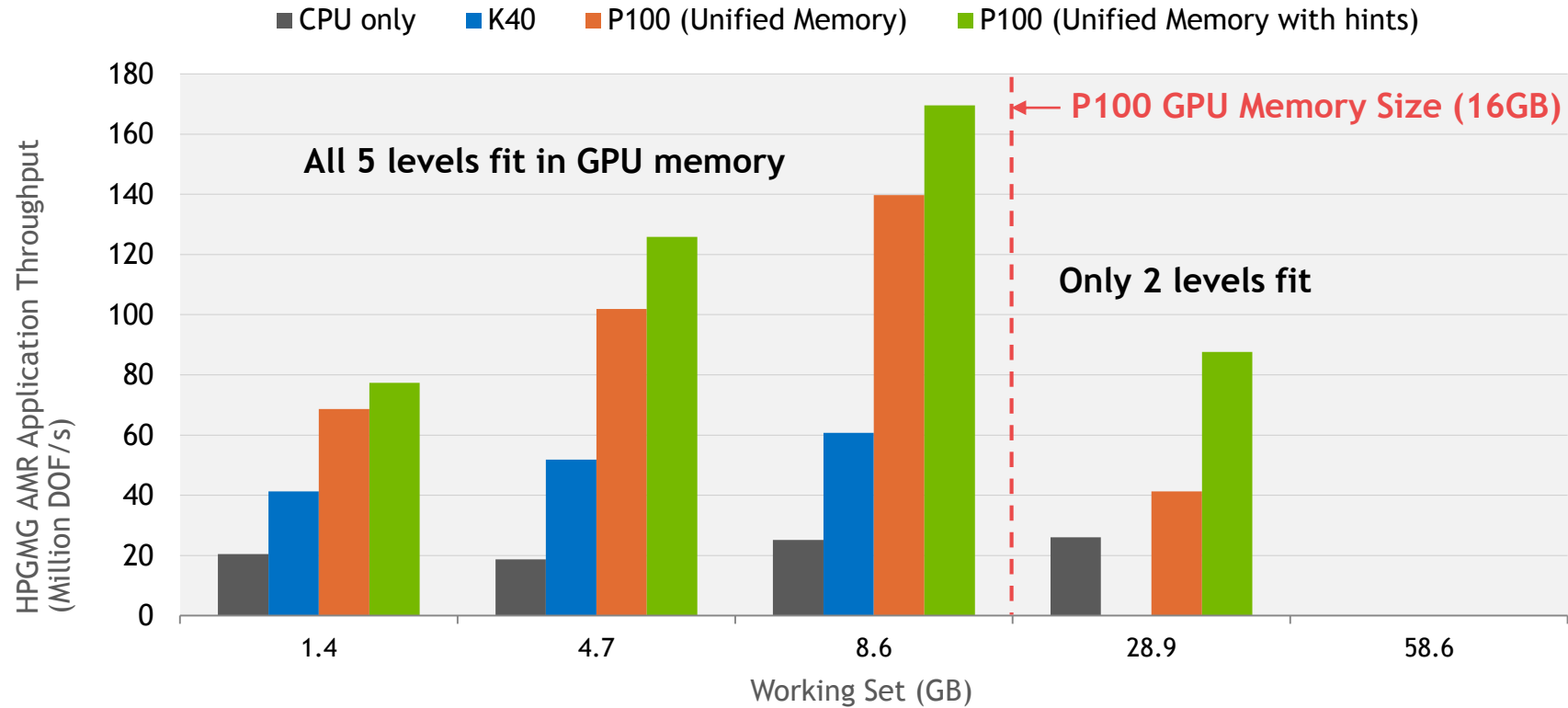- User APIs for prefetching & migration hints

**Simpler Data Access**
- CPU/GPU Data coherence
- Unified memory atomic operations

# SOLVE LARGER PROBLEMS WITH HIGHER THROUGHPUT



■ CPU only  ■ K40  ■ P100 (Unified Memory)  ■ P100 (Unified Memory with hints)

**All 5 levels fit in GPU memory**

HPGMG AMR Application Throughput (Million DOF/s)
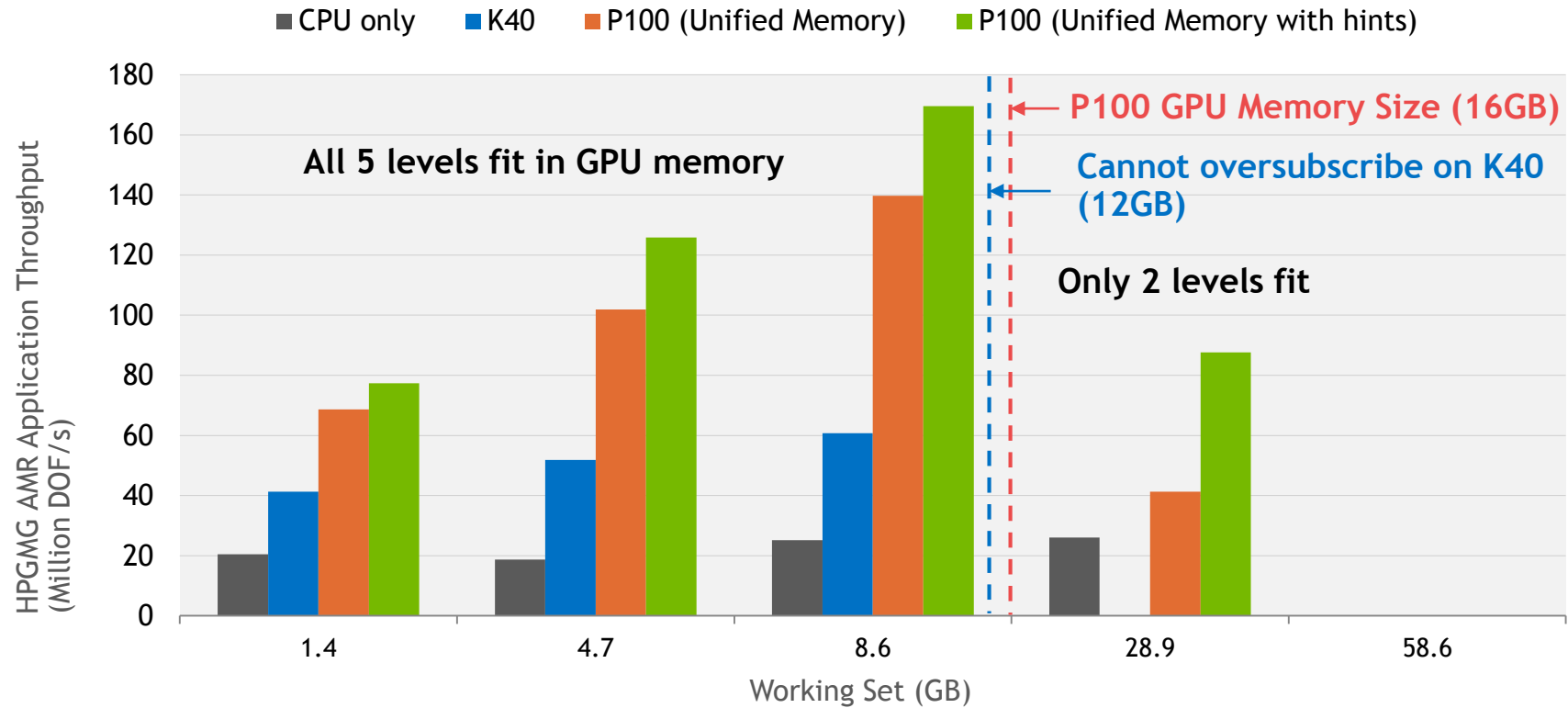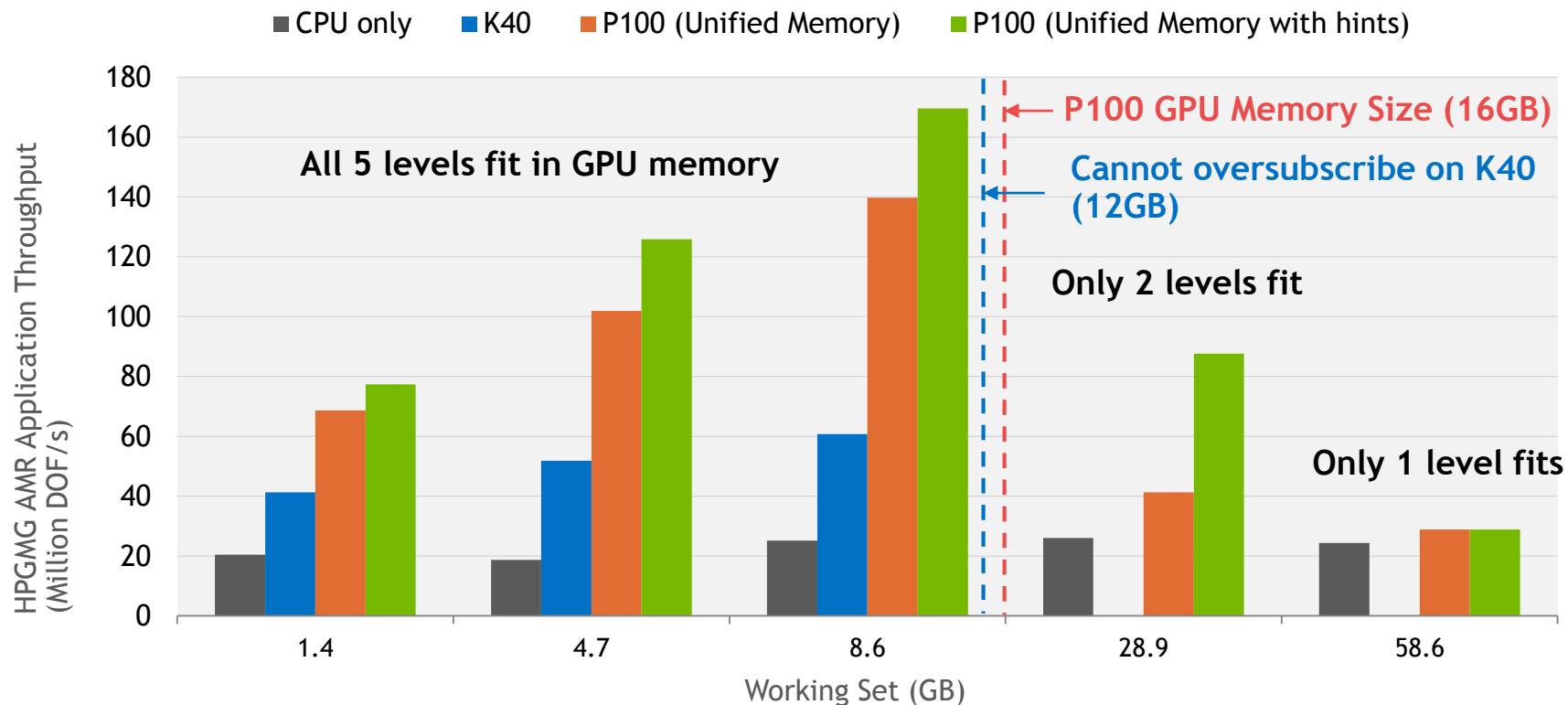
Working Set (GB)

◎ NVIDIA.

# SOLVE LARGER PROBLEMS WITH HIGHER THROUGHPUT



- HPGMG AMR on 1xK40, 1xP100 (PCIe) with CUDA 8 (r361)
- CPU measurements with Intel Xeon Haswell dual socket 10-core E5-2650 v3@2.3 GHz 3.0 GHz Turbo, HT on
- Host System: Intel Xeon Haswell dual socket 16-cores E5-2630 v3@2.4GHz 3.2GHz Turbo

⬡ nVIDIA.

# SOLVE LARGER PROBLEMS WITH HIGHER THROUGHPUT

**Legend:** ■ CPU only ■ K40 ■ P100 (Unified Memory) ■ P100 (Unified Memory with hints)

**All 5 levels fit in GPU memory**

P100 GPU Memory Size (16GB)

Cannot oversubscribe on K40 (12GB)

**Only 2 levels fit**

Y-axis: HPGMG AMR Application Throughput (Million DOF/s) — 0, 20, 40, 60, 80, 100, 120, 140, 160, 180

X-axis: Working Set (GB) — 1.4, 4.7, 8.6, 28.9, 58.6

NVIDIA.

# SOLVE LARGER PROBLEMS WITH HIGHER THROUGHPUT



Legend: ■ CPU only ■ K40 ■ P100 (Unified Memory) ■ P100 (Unified Memory with hints)

Y-axis: HPGMG AMR Application Throughput (Million DOF/s), scale 0–180

X-axis: Working Set (GB) — 1.4, 4.7, 8.6, 28.9, 58.6

Annotations:
- All 5 levels fit in GPU memory
- P100 GPU Memory Size (16GB)
- Cannot oversubscribe on K40 (12GB)
- Only 2 levels fit
- Only 1 level fits

⊗ nVIDIA.

# CUDA 8 NVCC: > 2X FASTER COMPILE TIMES

## Improved Developer Productivity



- Average total compile times (per translation unit)
- Host system: Intel Core i7-3930K 6-cores @ 3.2GHz
- CentOS x86_64 Linux release 7.1.1503 (Core) with GCC 4.8.3 20140911
- GPU target architecture sm_52

Performance may vary based on OS and software versions, and motherboard configuration

NVIDIA.

# GPU ACCELERATED LIBRARIES:
## GRAPH ANALYTICS

# nvGRAPH

## GPU Accelerated Graph Analytics



**Parallel Library for Interactive and High Throughput Graph Analytics**

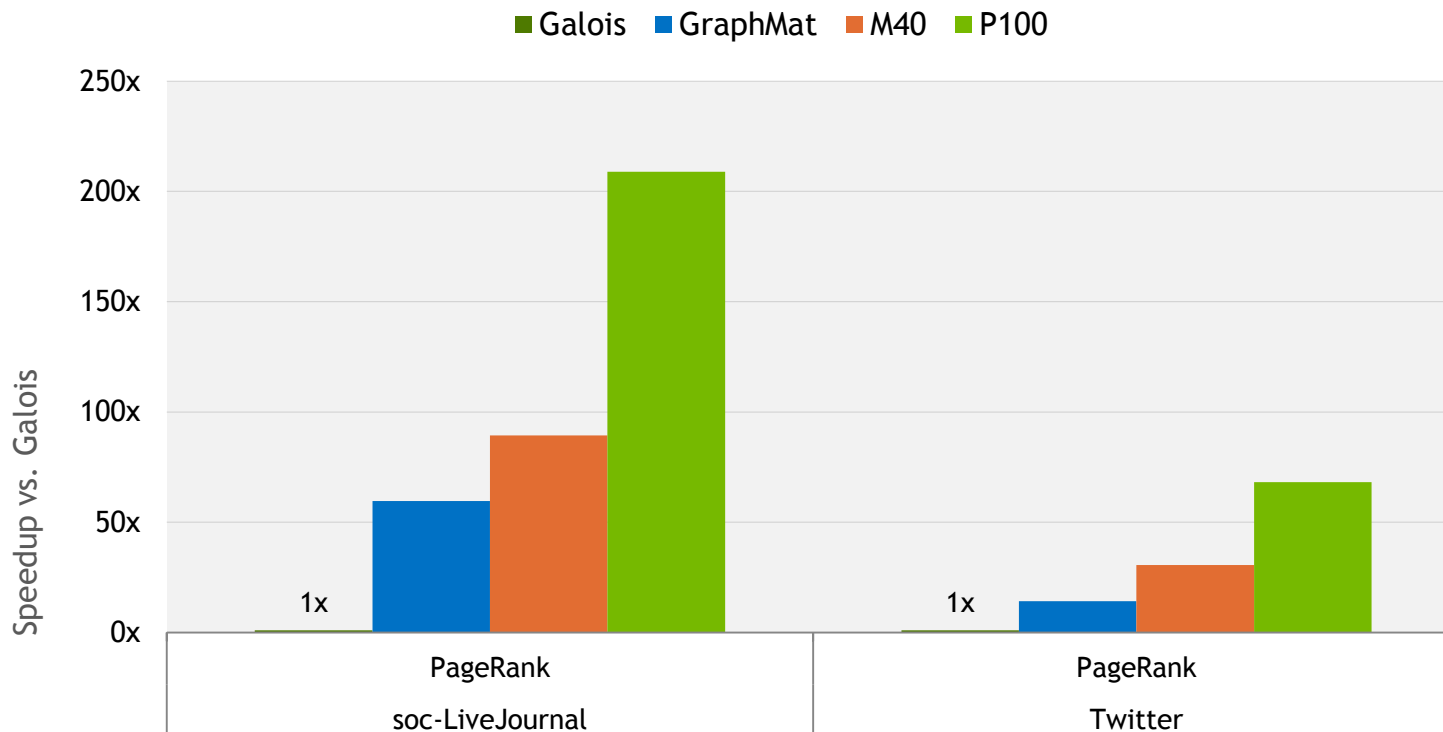Solve graphs with up to 2.5 Billion edges on a single GPU (Tesla M40)

Includes — PageRank, Single Source Shortest Path and Single Source Widest Path algorithms

Semi-ring SPMV operations provides building blocks for graph traversal algorithms

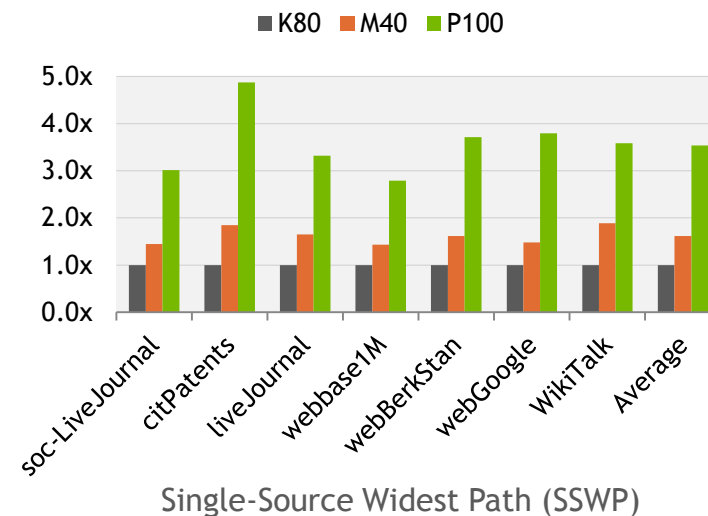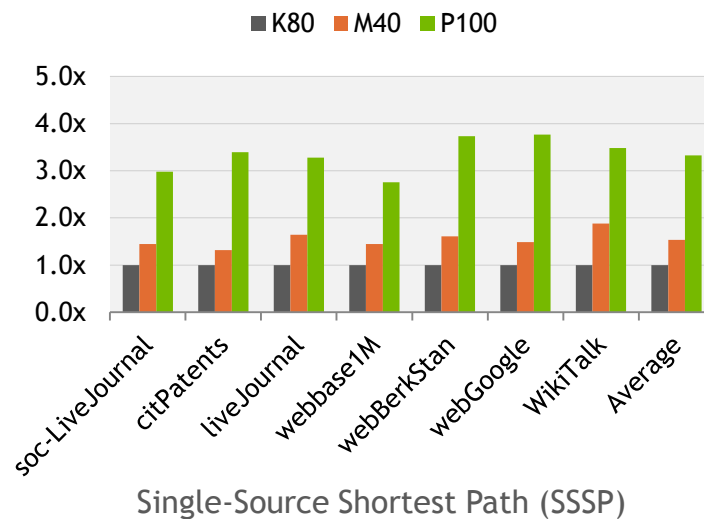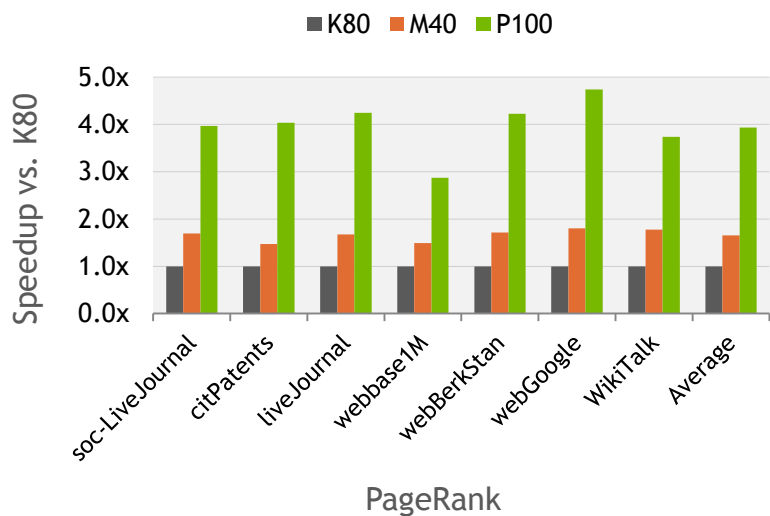| PageRank | Single Source Shortest Path | Single Source Widest Path |
|---|---|---|
| Search | Robotic Path Planning | IP Routing |
| Recommendation Engines | Power Network Planning | Chip Design / EDA |
| Social Ad Placement | Logistics & Supply Chain Planning | Traffic sensitive routing |

13 ⬡ nVIDIA.

# > 200X SPEEDUP ON PAGERANK VS GALOIS



Legend: ■ Galois ■ GraphMat ■ M40 ■ P100

- nvGRAPH on M40 (ECC ON, r352), P100 (r361), Base clocks, input and output data on device
- GraphMat, Galois (v2.3) on Intel Xeon Broadwell dual-socket 22-core/socket E5-2699 v4 @ 2.22GHz, 3.6GHz Turbo
- Comparing Average Time per Iteration (ms) for PageRank
- Host System: Intel Xeon Haswell single-socket 16-core E5-2698 v3 @ 2.3GHz, 3.6GHz Turbo
- CentOS 7.2 x86-64 with 128GB System Memory

Performance may vary based on OS and software versions, and motherboard configuration

14 ⬢ nVIDIA.

# > 4X SPEEDUPS WITH P100

## Using Different Algorithms in nvGRAPH



PageRank

Single-Source Shortest Path (SSSP)

Single-Source Widest Path (SSWP)

- nvGRAPH on K80, M40, P100, ECC ON, Base clocks, input and output data on device
- GraphMat, Galois (v2.3) on Intel Xeon Broadwell dual-socket 44-core E5-2699 v4 @ 2.22GHz, 3.6GHz Turbo
- Comparing Average Time per Iteration (ms) for PageRank and Total Solver Time (ms) for SSSP and SSWP
- Host System: Intel Xeon Haswell single-socket 16-core E5-2698 v3 @ 2.3GHz, 3.6GHz Turbo
- CentOS 7.2 x86-64 with 128GB System Memory

Performance may vary based on OS and software versions, and motherboard configuration

15 ⬡ nvIDIA.

# GPU ACCELERATED LIBRARIES:
## FAST FOURIER TRANSFORMS

# cuFFT

## Complete Fast Fourier Transforms Library

## Complete Multi-Dimensional FFT Library

Simple "drop-in" replacement of a CPU FFTW library

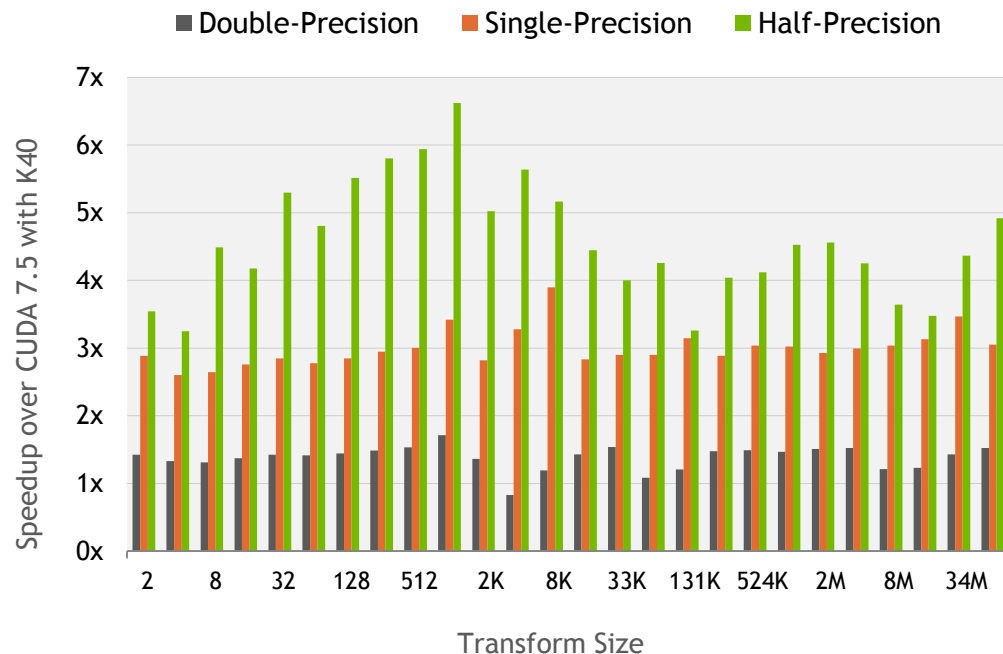Real and complex, single- and double-precision data types

Includes 1D, 2D and 3D batched transforms

Support for half-precision (FP16) data types

Supports flexible input and output data layouts

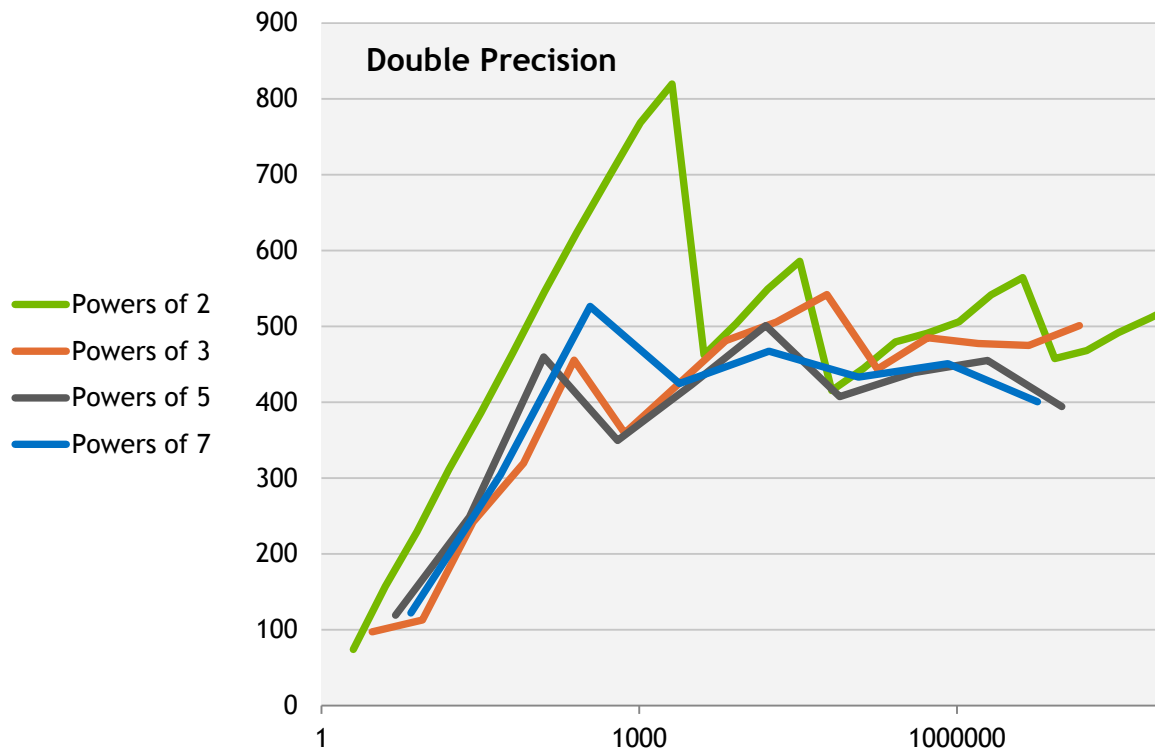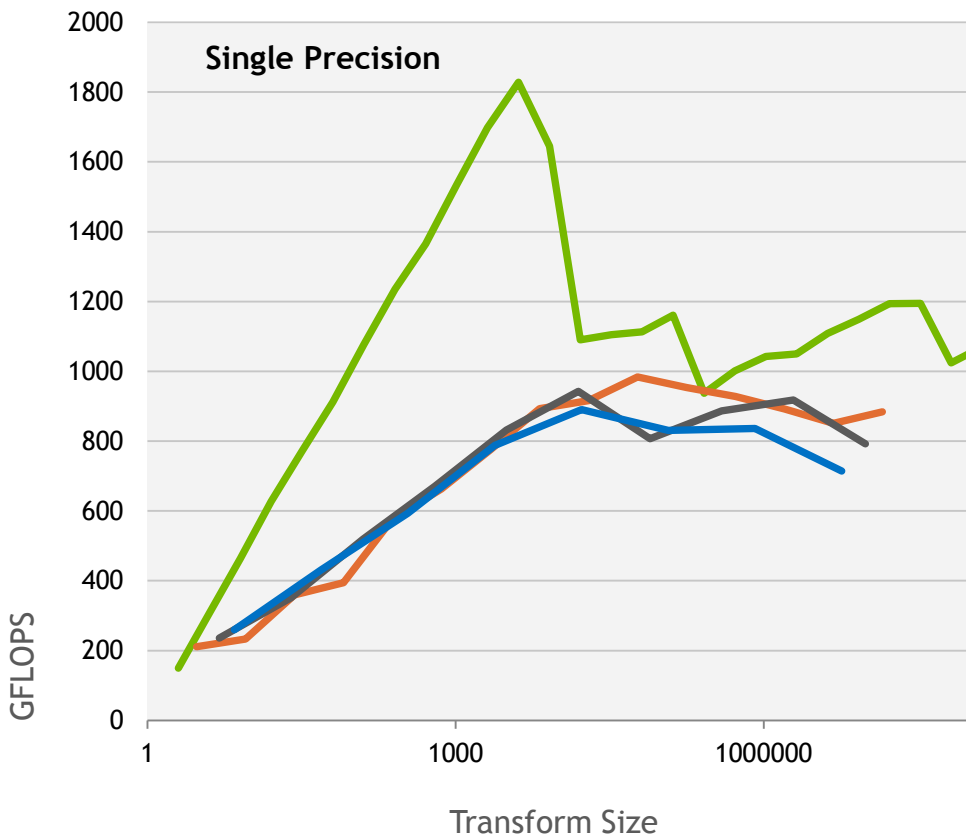XT interface now supports up to 8 GPUs

## > 6x Speedup with Half-Precision on P100

■ Double-Precision  ■ Single-Precision  ■ Half-Precision

_Speedup over CUDA 7.5 with K40_ (y-axis)

y-axis: 0x, 1x, 2x, 3x, 4x, 5x, 6x, 7x

x-axis (Transform Size): 2, 8, 32, 128, 512, 2K, 8K, 33K, 131K, 524K, 2M, 8M, 34M

- Speedup of P100 with CUDA 8 vs. K40m with CUDA 7.5
- cuFFT 7.5 on K40m, Base clocks, ECC on (r352)
- cuFFT 8.0 on P100, Base clocks, ECC on (r361)
- 1D Complex, Batched transforms on 28-33M elements
- Input and output data on device
- Host system: Intel Xeon Haswell single-socket 16-core E5-2698 v3 @ 2.3GHz, 3.6GHz Turbo
- CentOS 7.2 x86-64  with 128GB System Memory

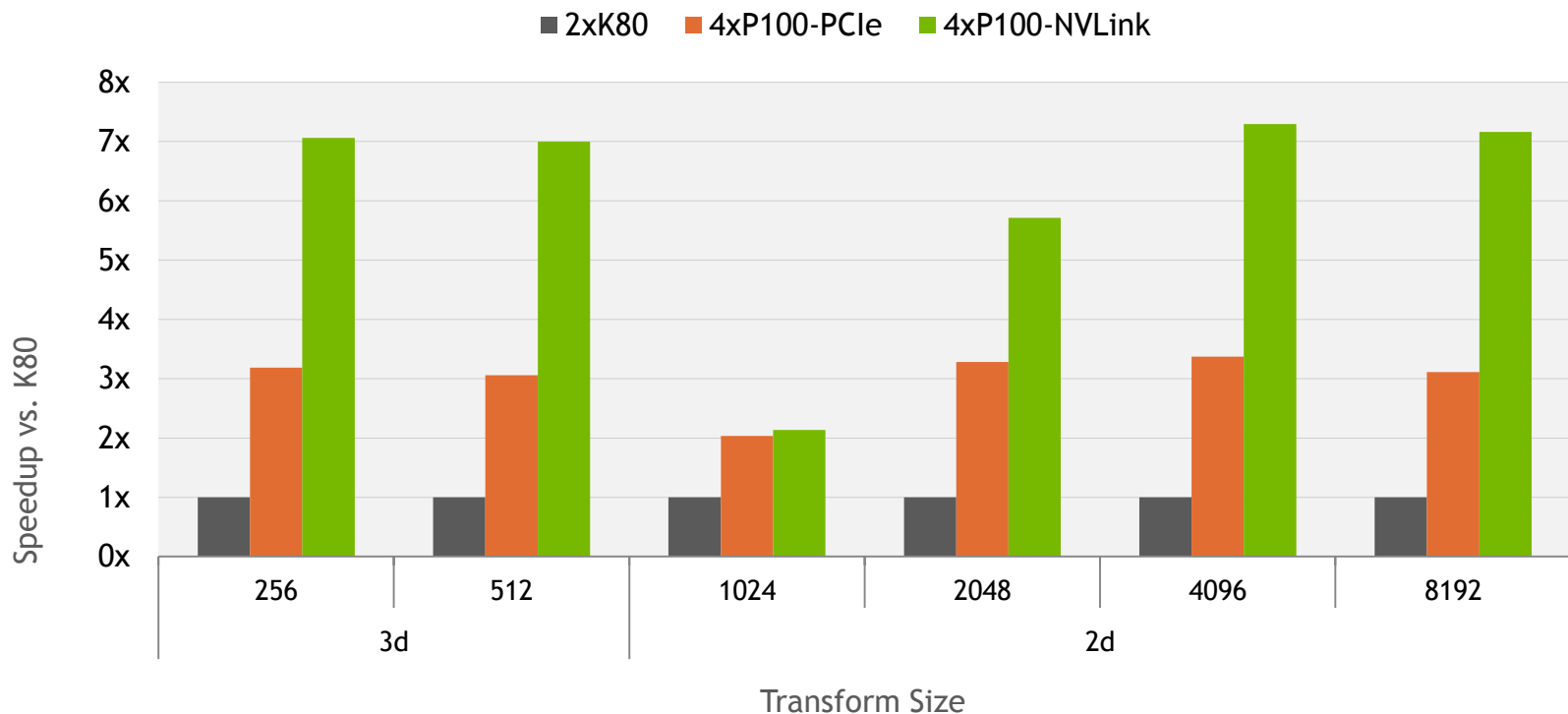# cuFFT: > 1800 GFLOPS SINGLE PRECISION

## 1D Complex, Batched FFTs



- cuFFT 8 on P100, Base clocks (r361)
- Batched transforms on 28-33M elements
- Input and output data on device
- Excludes time to create cuFFT "plans"
- Host system: Intel Xeon Haswell single-socket 16-core E5-2698 v3@ 2.3GHz, 3.6GHz Turbo
- CentOS 7.2 x86-64  with 128GB System Memory

Performance may vary based on OS and software versions, and motherboard configuration

18  NVIDIA.

# cuFFT-XT: > 7X IMPROVEMENTS WITH NVLINK
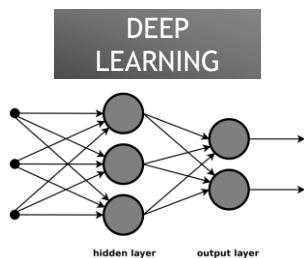
## 2D and 3D Complex FFTs



- cuFFT 7.5 on 2xK80m, ECC ON, Base clocks (r352)
- cuFFT 8 on 4xP100 with PCIe and NVLink (DGX-1), Base clocks (r361)
- Input and output data on device
- Excludes time to create cuFFT "plans"
- Host system: Intel Xeon Haswell single-socket 16-core E5-2698 v3@ 2.3GHz, 3.6GHz Turbo
- CentOS 7.2 x86-64 with 128GB System Memory
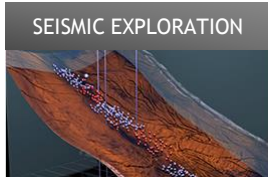
NVIDIA.

# DENSE AND SPARSE LINEAR ALGEBRA
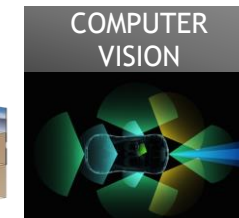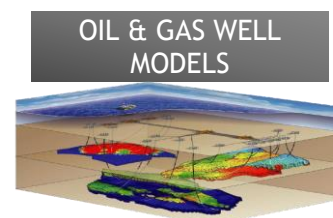
## Applications

### cuBLAS

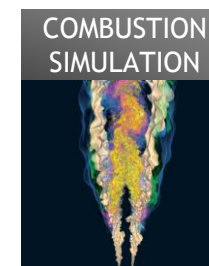**BASIC LINEAR ALGEBRA SUBPROGRAMS**


DEEP LEARNING

### cuSPARSE


COMPUTATIONAL FLUID DYNAMICS


CAD/CAM/CAE


SEISMIC EXPLORATION


RECOMMENDATION ENGINES


NLP

### cuSOLVER


COMBUSTION SIMULATION


OIL & GAS WELL MODELS


COMPUTER VISION

# cuBLAS

## Dense Linear Algebra on GPUs

## Complete BLAS Library Plus Extensions

Supports all 152 standard routines for single, double, complex, and double complex
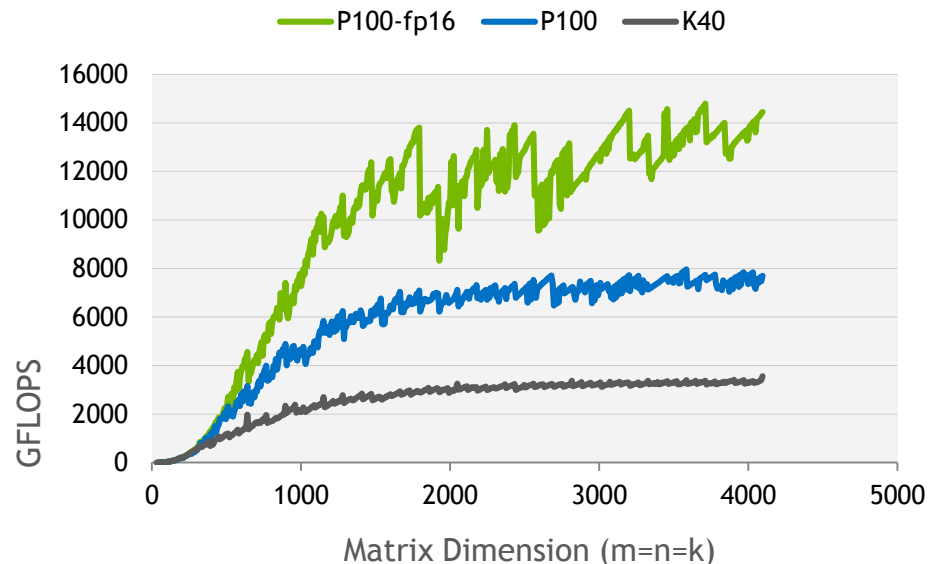
Supports half-precision (FP16) and integer (INT8) matrix multiplication operations

Batched routines for higher performance on small problem sizes

Host and device-callable interface

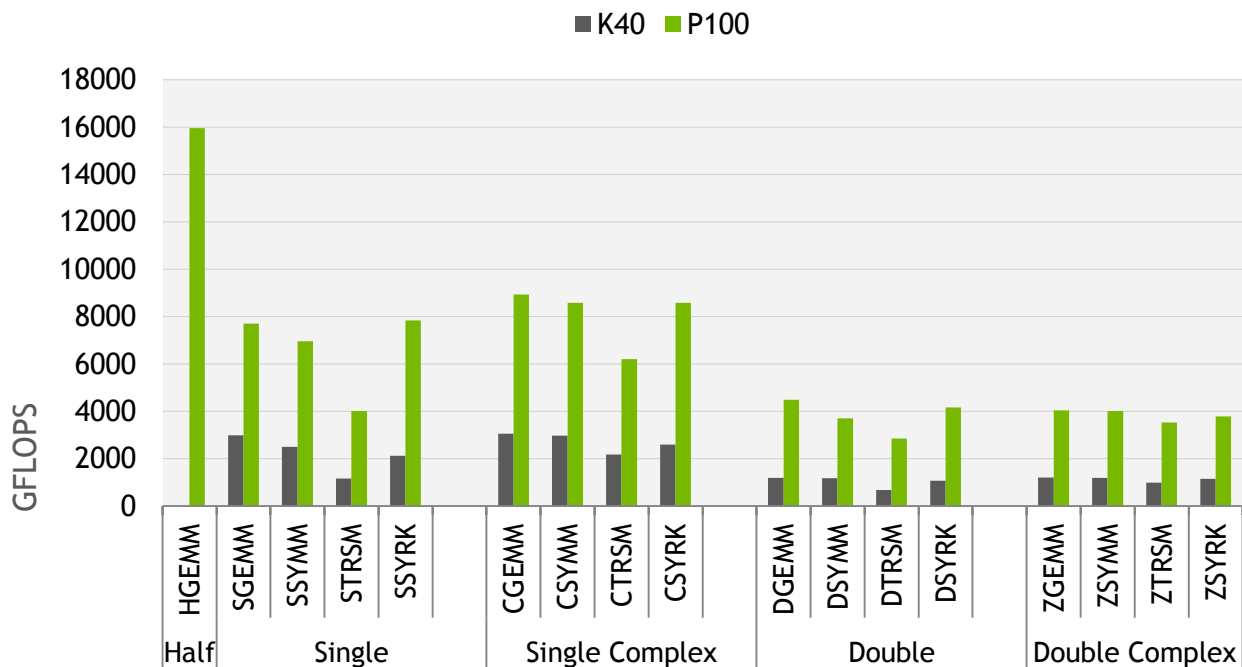XT interface supports distributed computations across multiple GPUs

## > 4x Faster GEMM Performance with FP16 on P100

**Legend:** P100-fp16 ── P100 ── K40

**Y-axis:** GFLOPS (0, 2000, 4000, 6000, 8000, 10000, 12000, 14000, 16000)

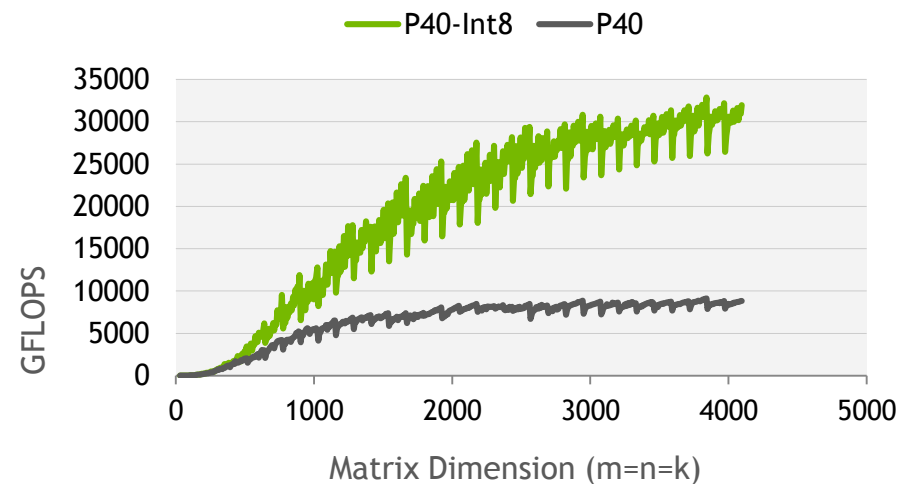**X-axis:** Matrix Dimension (m=n=k) (0, 1000, 2000, 3000, 4000, 5000)

- Comparing GEMM performance on K40m (FP32) and P100 (FP32 and FP16)
- cuBLAS 8 on P100, Base clocks (r361)
- cuBLAS 8 on P40, Base clocks (r367)
- cuBLAS 7.5 on K40m, Base clocks, ECC ON (r352)
- Input and output data on device
- Host system: Intel Xeon Haswell single-socket 16-core E5-2698 v3@ 2.3GHz, 3.6GHz Turbo
- CentOS 7.2 x86-64 with 128GB System Memory
- m=n=k=4096

NVIDIA.

# cuBLAS: > 8 TFLOPS SINGLE PRECISION
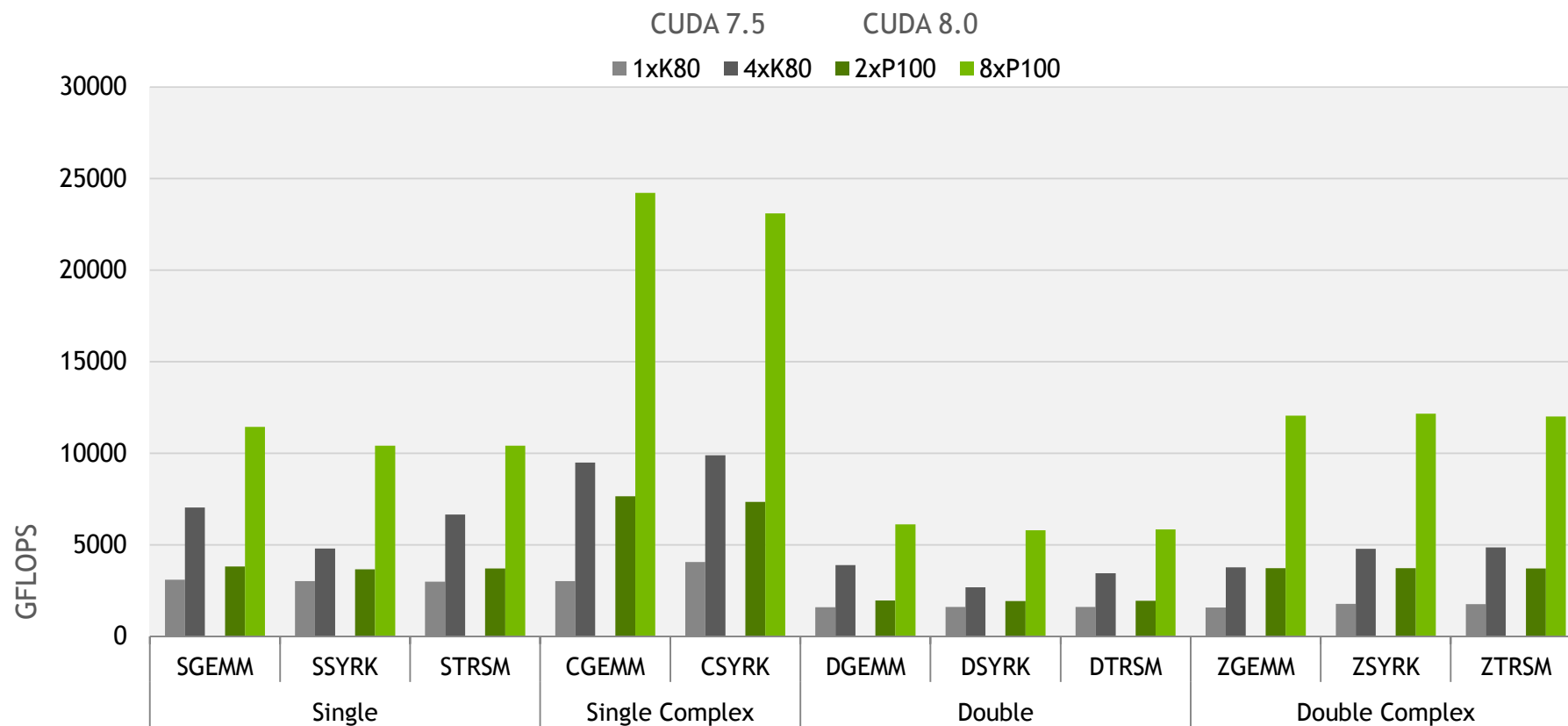
## 16 TFLOPS FP16 GEMM Performance

Legend: ■ K40 ■ P100



## 32 TFLOPS INT8 GEMM Performance

Legend: — P40-Int8 — P40



- cuBLAS 8 on P100 (r361) and P40 (r367) ; Base clocks
- cuBLAS 7.5 on K40m ; Base clocks, ECC ON (r352)
- Input and output data on device
- m=n=k=4096, transpose=no, side=right, fill=lower
- Host system: Intel Xeon Haswell single-socket 16-core E5-2698 v3@ 2.3GHz, 3.6GHz Turbo
- CentOS 7.2 x86-64  with 128GB System Memory

# cuBLAS-XT: > 24 TFLOPS ON A SINGLE NODE



Performance may vary based on OS and software versions, and motherboard configuration

- cuBLAS 8 on P100 (r361); Base clocks
- cuBLAS 7.5 on K80 ; Base clocks, ECC ON (r352)
- 1xK80 indicates 2-GPUs (or one K80 board)
- Input and output data on device
- m=n=k=4096, transpose=no, side=right, fill=lower
- Host system: Intel Xeon Haswell dual-socket 22-core E5-2699 v4@ 2.2GHz, 3.6GHz Turbo
- CentOS 7.2 x86-64  with 256GB System Memory

23  NVIDIA.

# cuSPARSE

## Sparse Linear Algebra on GPUs

## Optimized Sparse Matrix Library

Optimized sparse linear algebra BLAS routines for matrix-vector, matrix-matrix, triangular solve

Support for variety of formats (CSR, COO, block variants)
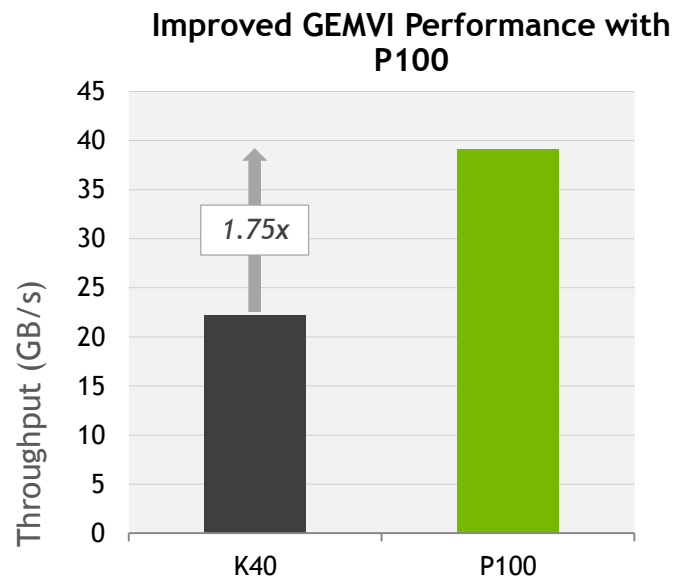
Incomplete-LU and Cholesky preconditioners

Support for half-precision (fp16) sparse matrix-vector operations

## GEMVI – Dense Matrix X Sparse Vector

$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} = \alpha \begin{bmatrix} A_{11} & A_{12} & A_{13} & A_{14} \\ A_{21} & A_{22} & A_{23} & A_{24} \\ A_{31} & A_{32} & A_{33} & A_{34} \\ A_{41} & A_{42} & A_{43} & A_{44} \end{bmatrix} \begin{bmatrix} - \\ 2.0 \\ - \\ 4.0 \end{bmatrix} + \beta \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix}$$
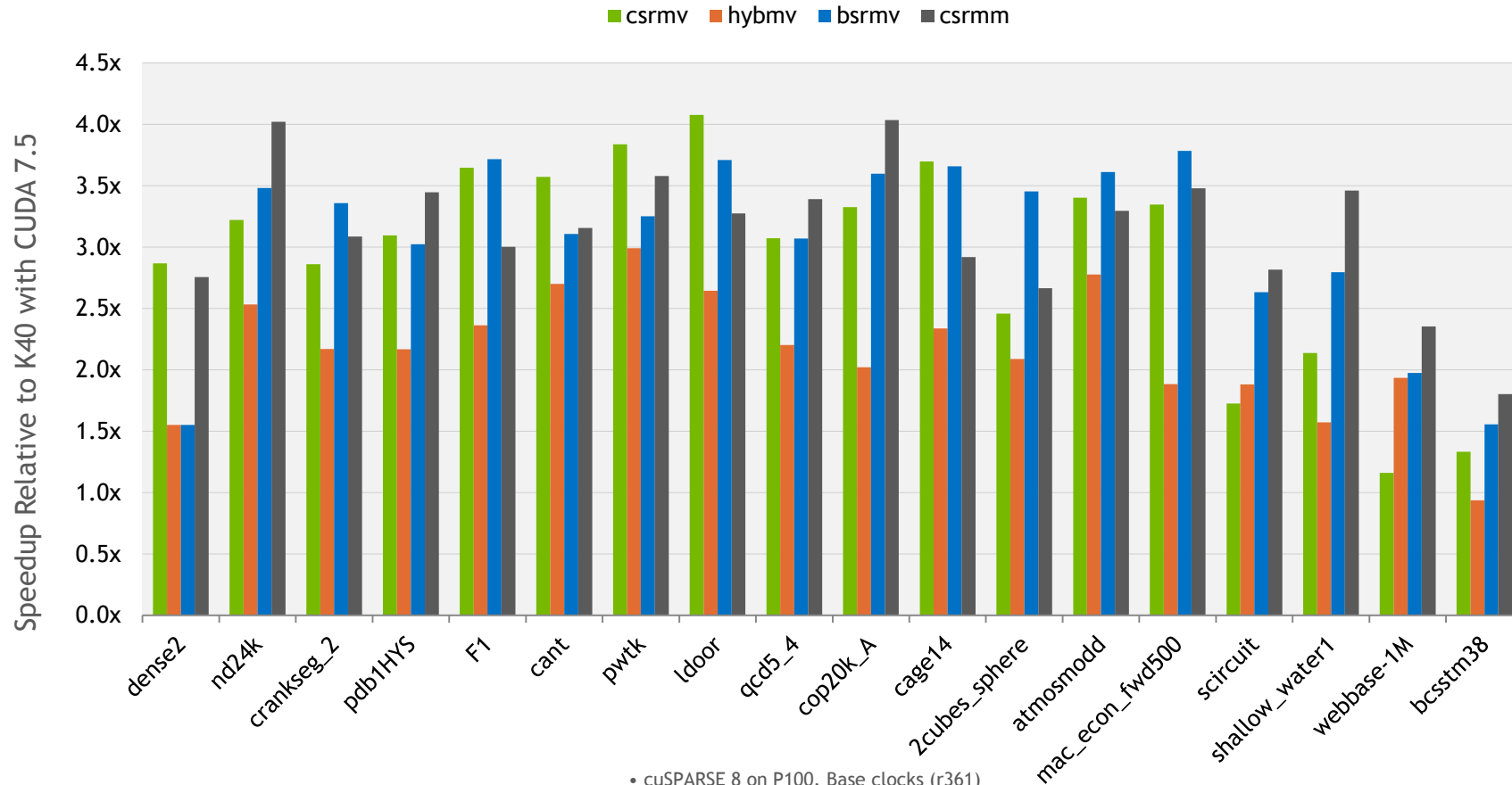
Used in language modeling and auto-encoders for recommender systems

**Improved GEMVI Performance with P100**



Throughput (GB/s) — K40: ~22, P100: ~39 — 1.75x

- cuSPARSE 8 on P100, Base clocks (r361)
- cuSPARSE 7.5 on K40m, Base clocks, ECC ON (r352)
- Input and output data on device
- Dense matrices with 1e6 columns and 1e3 rows; Sparse vectors with less than 100 non-zeros out of 1e6 locations
- Host system: Intel Xeon Haswell single-socket 16-core E5-2698 v3 @ 2.3GHz, 3.6GHz Turbo
- CentOS 7.2 x86-64 with 128GB System Memory

24 NVIDIA

# cuSPARSE: > 4X FASTER WITH P100

Legend: ■ csrmv ■ hybmv ■ bsrmv ■ csrmm



Speedup Relative to K40 with CUDA 7.5

Categories: dense2, nd24k, crankseg_2, pdb1HYS, F1, cant, pwtk, ldoor, qcd5_4, cop20k_A, cage14, 2cubes_sphere, atmosmodd, mac_econ_fwd500, scircuit, shallow_water1, webbase-1M, bcsstm38
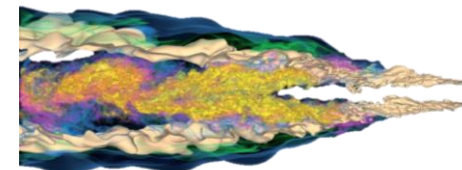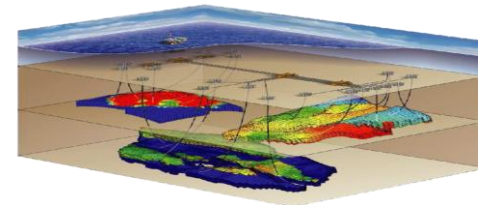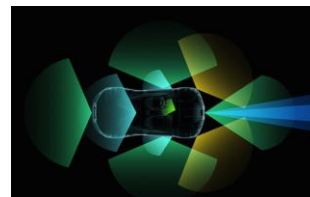
NVIDIA.

# cuSOLVER
## Linear Solver Library

## Library for Dense and Sparse Direct Solvers

Supports Dense Cholesky, LU, (batched) QR, SVD and Eigenvalue solvers (new in CUDA 8)

Sparse direct solvers & Eigensolvers

Includes a sparse refactorization solver for solving sequences of matrices with a shared sparsity pattern

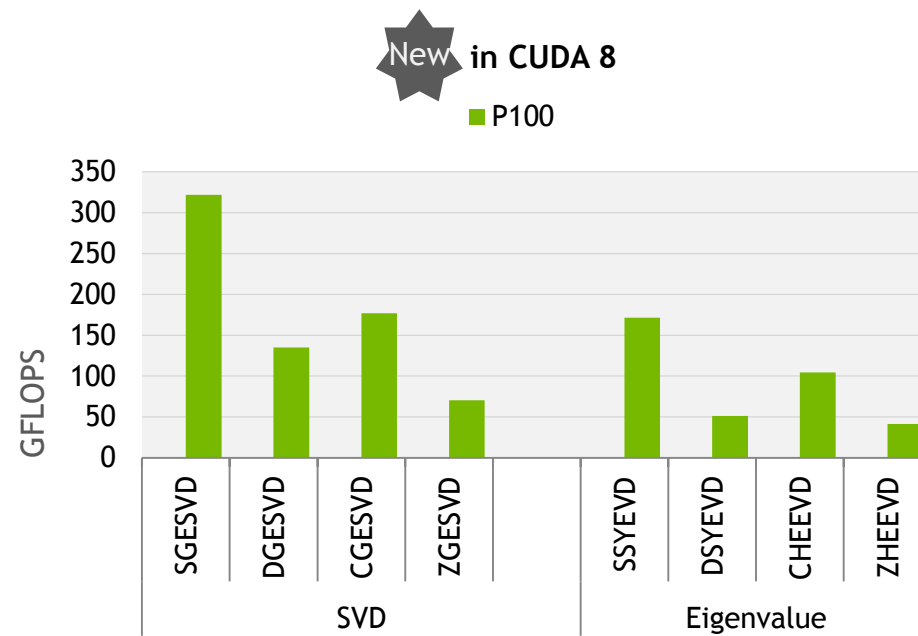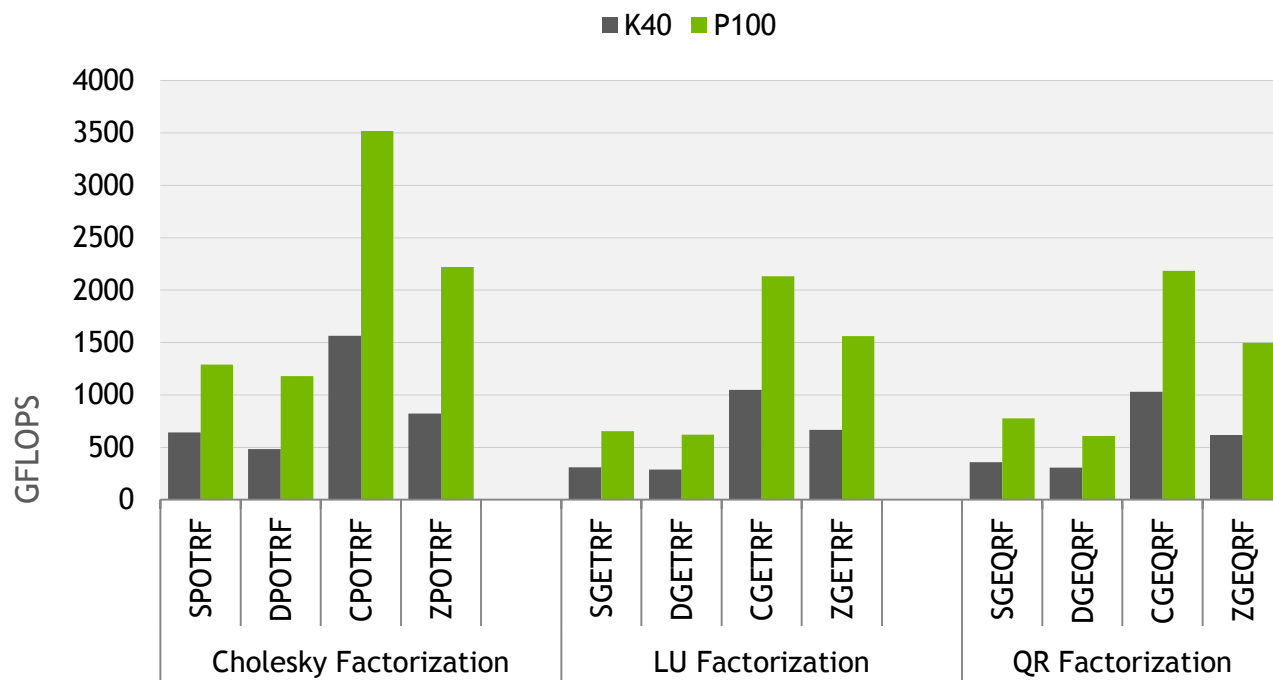Used in a variety of applications such as circuit simulation and computational fluid dynamics



**Sample Applications**
- **Computer Vision**
- **CFD**
- **Newton's method**
- **Chemical Kinetics**
- **Chemistry**
- **ODEs**
- **Circuit Simulation**

NVIDIA.

# DENSE PERFORMANCE: > 2X FASTER

Performance may vary based on OS and software versions, and motherboard configuration

- cuSOLVER 8 on P100, Base clocks (r361)
- cuSOLVER 7.5 on K40m, Base clocks, ECC ON (r352)
- M=N=4096
- Host system: Intel Xeon Haswell single-socket 16-core E5-2698 v3 @ 2.3GHz, 3.6GHz Turbo
- CentOS 7.2 x86-64 with 128GB System Memory
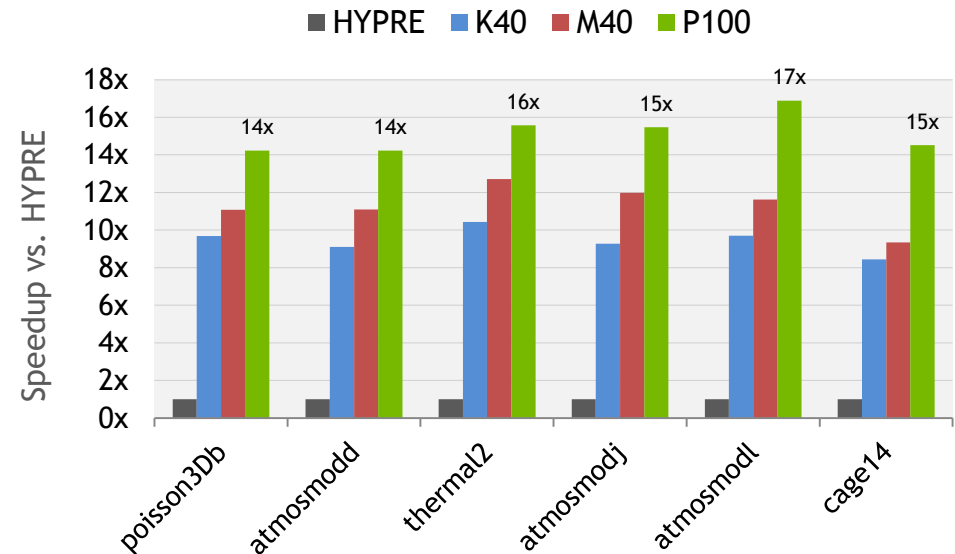
# AmgX
## Algebraic Multi-Grid Solvers

## Flexible Solver Composition System

Easy construction of complex nested solvers and pre-conditioners

Flexible and simple high level C API that abstracts parallelism and GPU implementation

Includes Ruge-Steuben, un-smoothed aggregation, Krylov methods and different smoother algorithms
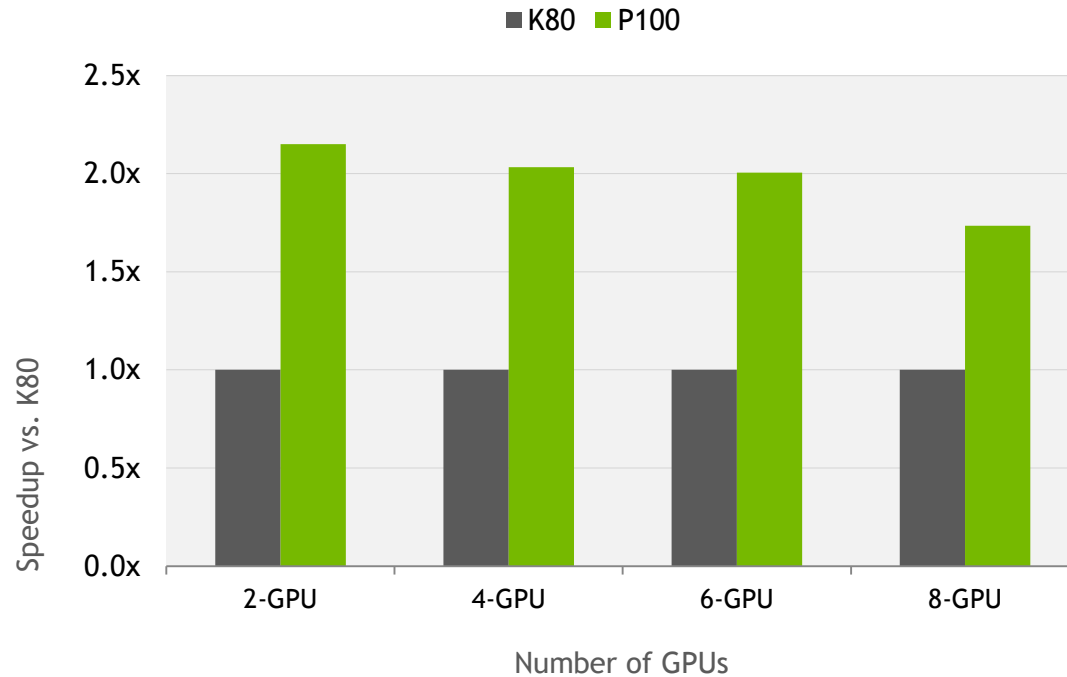
### > 15x Speedup vs HYPRE

Legend: ■ HYPRE ■ K40 ■ M40 ■ P100

Bar chart — Speedup vs. HYPRE (y-axis 0x to 18x):
- poisson3Db: K40 ~9.7x, M40 ~11x, P100 14x
- atmosmodd: K40 ~9.1x, M40 ~11x, P100 14x
- thermal2: K40 ~10.5x, M40 ~12.7x, P100 16x
- atmosmodj: K40 ~9.3x, M40 ~12x, P100 15x
- atmosmodl: K40 ~9.7x, M40 ~11.6x, P100 17x
- cage14: K40 ~8.4x, M40 ~9.4x, P100 15x

- Florida Matrix Collection; Total Time to Solution
- HYPRE AMG Package (http://acts.nersc.gov/hypre) on Intel Xeon E5-2697 v4@2.3GHz, 3.6GHz Turbo, Hyperthreading off
- AmgX on K40, M40, P100 (SXM2); Base clocks
- Host system: Intel Xeon Haswell single-socket 16-core E5-2698 v3 @ 2.3GHz, 3.6GHz Turbo
- CentOS 7.2 x86-64 with 128GB System Memory

# > 2X FASTER SOLUTION TIMES WITH P100
## SPE10 Using Multi-GPUs

■K80 ■P100



- AmgX on K80 and M40 with CUDA 7.5, and P100 (PCIe) with CUDA 8; Base clocks, ECC ON
- Society of Petroleum Engineers 10th comparative simulation model (SPE10)
- Matrix size varied with 100 x Number of GPUs x 200 x 50 for SPE10
- Time to solution includes both setup and solve times

⬨ nVIDIA.

# GPU ACCELERATED LIBRARIES:
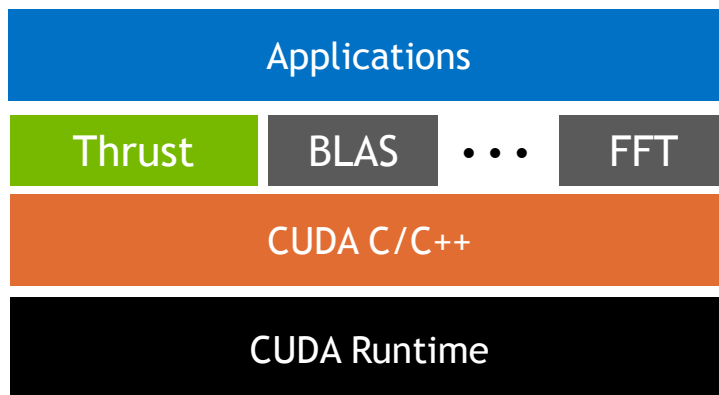## PARALLEL PRIMITIVES, IMAGE PROCESSING AND RNGs

# THRUST

## Parallel Primitives and Data Structures

## Templated Library of Algorithms and Data Structures

GPU-accelerted scan, sort, transform and reduce

Improved developer productivity via C++ template library, allowing developer to focus on high-level tasks

Library simplifies memory management and data transfer

| Applications | | |
|:---:|:---:|:---:|
| Thrust | BLAS ••• | FFT |
| CUDA C/C++ | | |
| CUDA Runtime | | |

**Thrust is a C++ template library of parallel algorithms inspired by C++ STL**

NVIDIA.

# HIGH PERFORMANCE PARALLEL PRIMITIVES

## > 2.5x Primitives Performance



## > 3x Sorting Performance



- Thrust 8 on P100, Base clocks (r361)
- Thrust 7.5 on K40m, M40, Base clocks, ECC ON (r352)
- Input and output data on device
- Host system: Intel Xeon Haswell single-socket 16-core E5-2698 v3 @ 2.3GHz, 3.6GHz Turbo
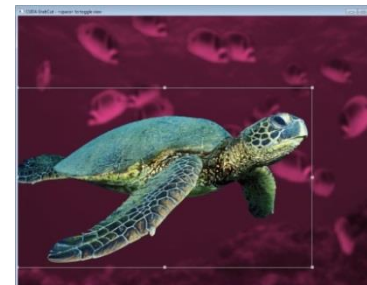- CentOS 7.2 x86-64  with 128GB System Memory

NVIDIA.

# NPP

## NVIDIA Performance Primitives Library

GPU-accelerated Building Blocks for Image, Video Processing & Computer Vision

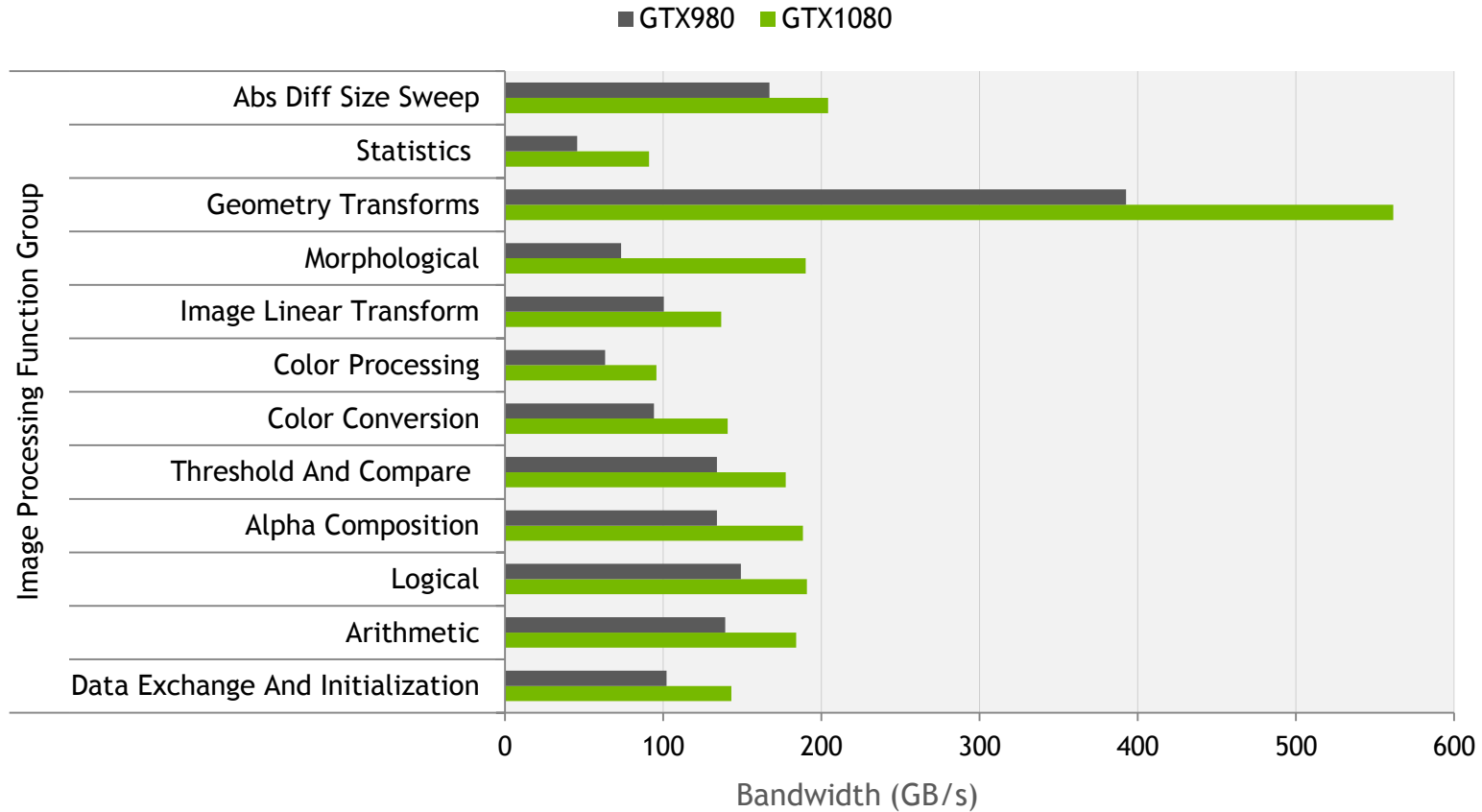Over 2500 image, signal processing and computer vision routines

Color transforms, geometric transforms, move operations, linear filters, image & signal statistics, image & signal arithmetic, building blocks, image segmentation, median filter, BGR/YUV conversion, 3D LUT color conversion

Eliminate unnecessary copying of data to/from CPU memory



"Grabcut" example using NPP graphcut primitive

# > 2X FASTER IMAGE PROCESSING WITH PASCAL

■ GTX980 ■ GTX1080

Chart: Bandwidth (GB/s) by Image Processing Function Group

- Abs Diff Size Sweep
- Statistics
- Geometry Transforms
- Morphological
- Image Linear Transform
- Color Processing
- Color Conversion
- Threshold And Compare
- Alpha Composition
- Logical
- Arithmetic
- Data Exchange And Initialization

X-axis: Bandwidth (GB/s) — 0, 100, 200, 300, 400, 500, 600

Y-axis: Image Processing Function Group

- NPP 7.5 on GTX980 (GM204) and NPP 8 on GTX1080, Base clocks
- Input and output data on device
- Each bar represents the average bandwidth over all routines in the function group
- Host system: Intel Xeon Haswell single-socket 16-core E5-2698 v3 @ 2.3GHz, 3.6GHz Turbo
- CentOS 7.2 x86-64 with 128GB System Memory

⬢ nVIDIA.

# GPU ACCELERATED LIBRARIES:
## RNGs
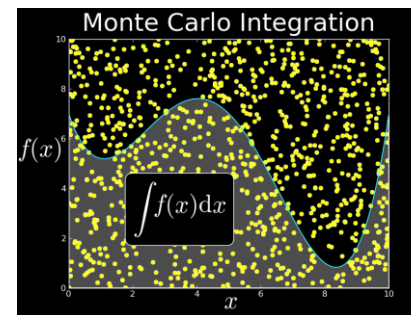
# cuRAND

## Random Number Generation (RNG) Library

**High Performance Random Number Generation**

Flexible interfaces for RNG on the host or within GPU kernels

Pseudo- and Quasi-RNGs — MRG32k3a, MTGP Mersenne Twister, XORWOW, Sobol
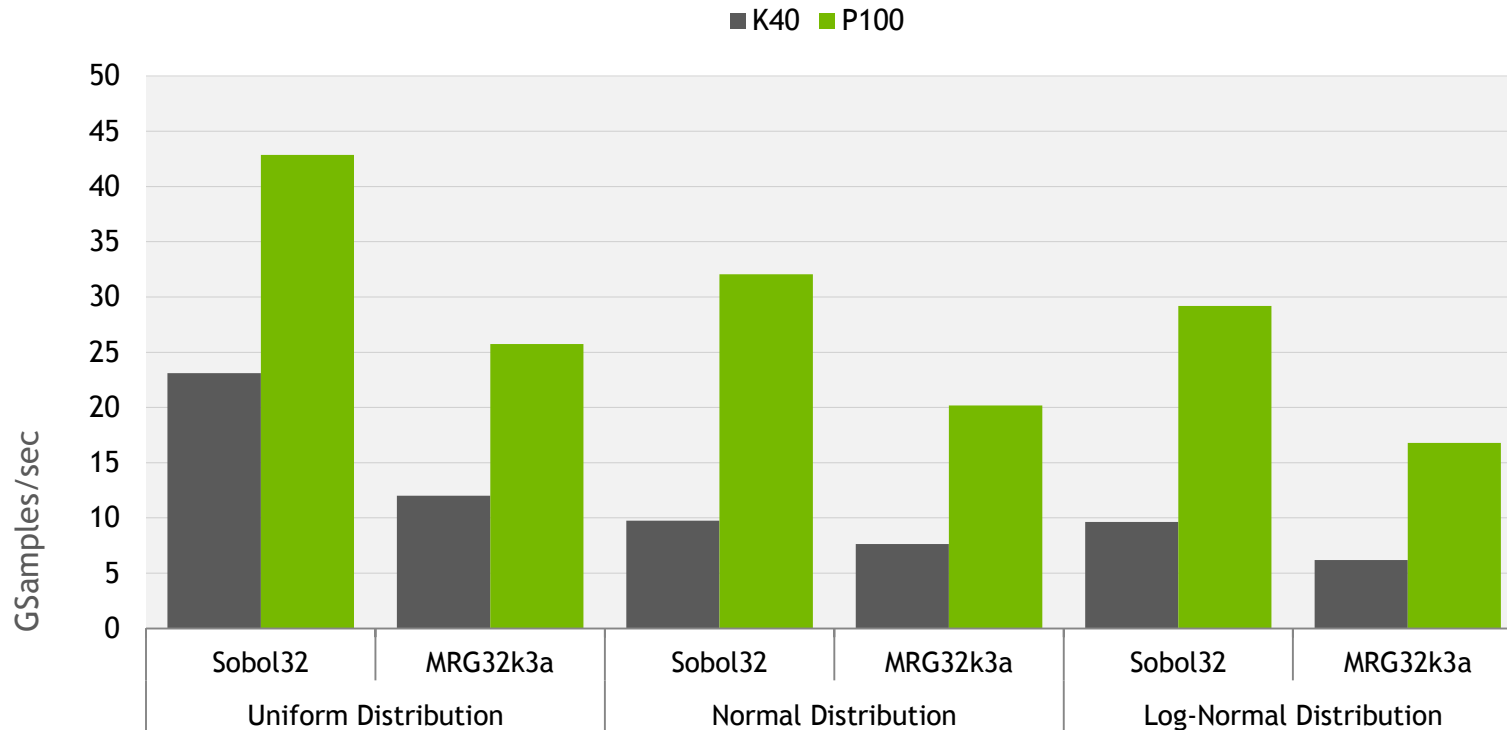
Supports several output distributions

Tested against well-known statistical test batteries (test results available in documentation)
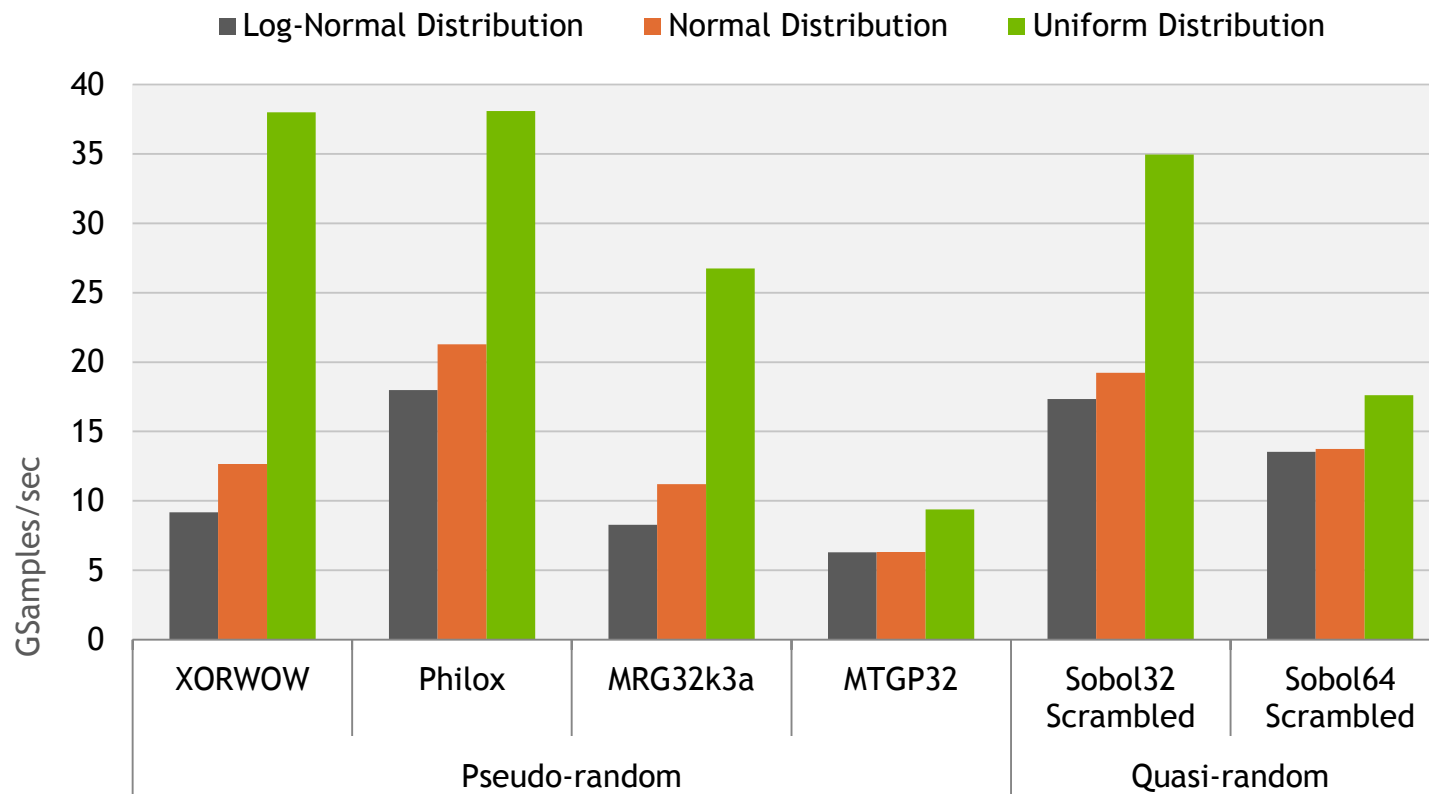


**Sample Applications of cuRAND**

# cuRAND: > 3X SPEEDUPS WITH P100



- cuRAND 8 on P100, Base clocks (r361)
- cuRAND 7.5 on K40m, Base clocks, ECC ON, (r352)
- Single-precision input and output data on device
- Host system: Intel Xeon Haswell single-socket 16-core E5-2698 v3 @ 2.3GHz, 3.6GHz Turbo
- CentOS 7.2 x86-64 with 128GB System Memory

NVIDIA.

# cuRAND: HIGH PERFORMANCE RNGS



Legend: ■ Log-Normal Distribution ■ Normal Distribution ■ Uniform Distribution

Chart categories (left to right):
XORWOW, Philox, MRG32k3a, MTGP32 (Pseudo-random); Sobol32 Scrambled, Sobol64 Scrambled (Quasi-random)

Y-axis: GSamples/sec

- cuRAND 8 on P100, Base clocks (r361)
- Double-precision input and output data on device
- Host system: Intel Xeon Haswell single-socket 16-core E5-2698 v3 @ 2.3GHz, 3.6GHz Turbo
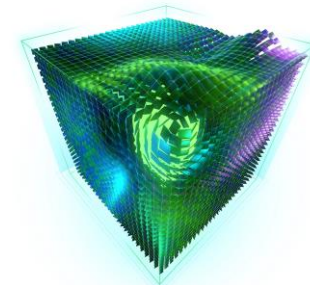- CentOS 7.2 x86-64 with 128GB System Memory

<NVIDIA.>

# CUDA 8 – DOWNLOAD TODAY!

## Everything You Need to Accelerate Applications

Google | cuda 8 download

- CUDA Applications in your Industry: www.nvidia.com/object/gpu-applications-domain.htm

- Additional Webinars:

  - Inside PASCAL
  - Deep Learning and Hyperscale Performance Overview
  - CUDA 8 Tools
  - CUDA 8 Unified Memory

- CUDA 8 Release Notes: www.docs.nvidia.com/cuda/cuda-toolkit-release-notes/index.html#abstract

developer.nvidia.com/cuda-toolkit

NVIDIA.

# APPENDIX

# CONFIGURATIONS USED FOR HPC APPS

| Application | Dataset and CUDA | Metric |
|---|---|---|
| HOOMD-BLUE – Particle dynamics package written grounds up for GPUs<br>http://codeblue.umich.edu/hoomd-blue/index.html | lj_liquid_1m | Avg. Timesteps (seconds) |
| HPCG – Benchmark that exercises computational and data access patterns that closely match a broad set of important HPC applications<br>http://www.hpcg-benchmark.org/index.html | Local Size: 256x256x256<br>CUDA Versions: CUDA 6.5.19 | GFLOPS/second |
| LAMMPS – Classical molecular dynamics package<br>http://lammps.sandia.gov/index.html | EAM Mixed Precision;<br>Problem Size: 4x4x4; 2048000 atoms | Avg. Loop Time (seconds) |
| QUDA  -- A library for Lattice Quantum Chromo Dynamics on GPUs<br>https://lattice.github.io/quda/ | QUDA (GPU):<br>Dslash Wilson-Clover<br>Precision: Double<br>Problem Size: 32x32x32x64<br>QPhiX (CPU):<br>Dslash Wilson-Clover<br>Precision: Double<br>Problem Size: 32x32x32x64 | GFLOPS |
| MiniFE – Finite Element Analysis<br>https://mantevo.org/about/applications | Problem Size: 350x350x350 | Total Time (seconds) |
| MILC – Lattice Quantum Chromodynamics (LQCD) codes simulate how elemental particles are formed and bound by the "strong force" to create larger particles like protons and neutrons<br>http://physics.indiana.edu/~sg/milc.html | Precision: Double | Total Time (seconds) |
| VASP<br>http://www.vasp.at/index.php/news/44-administrative/115-new-release-vasp-5-4-1-with-gpu-support | Silica IFPEN | Elapsed Time (seconds) |

CPU Server:       Intel Xeon dual-socket 22-core E5-2699 v4@2.2GHz 3.6GHz Turbo with CentOS 7.2 x86_64 and 128GB system memory
GPU Host Server: Intel Xeon dual-socket 20-core E5-2698 v4@2.2GHz 3.6GHz Turbo with Tesla P100 SXM2 ; Ubuntu 14.04.4 x86_64 with 128GB system memory
GPU Host Server: Intel Xeon dual-socket 22-core E5-2699 v4@2.2GHz 3.6GHz Turbo with Tesla K80  ; CentOS 7.2 x86_64 and 128GB system memory
CUDA Versions:   CUDA 8.0.44 with r361.96 (P100) and 361.79 (K80) unless otherwise specified

NVIDIA.

# SYSTEM TOPOLOGY USED FOR K80 AND P100



2 CPU – 2x (K80)

2 CPU – 2x P100 SXM2

2 CPU – 4x P100 SXM2

2 CPU – 8x P100 SXM2

NVLink    PCIE Gen 3