GRAPHCORE



AI MEETS HPC IN SCIENTIFIC COMPUTING USING GRAPHCORE IPUS

Paweł Gepner Senior System Field Engineer





AGENDA

- WELCOME AND GRAPHCORE INTRODUCTION
- IPU ARCHITECTURE OVERVIEW
- INSIGHTS INTO POD SYSTEM
- SHORT INTRODUCTION TO THE POPLAR SDK
- ACCELERATION OF SCIENTIFIC APPLICATIONS AND ENABLING HETEROGENEITY OF AI AND HPC
- HYBRIDIZATION OF HPC AND AI SYSTEMS
- Q&A SESSION

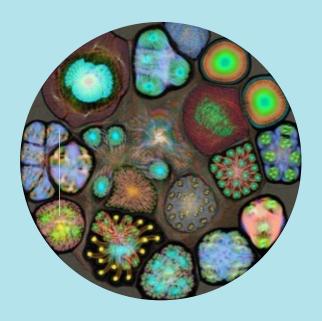
ABOUT US...

Hardware



IPU processors designed for AI

Software



Poplar® software stack & development tools

Platforms

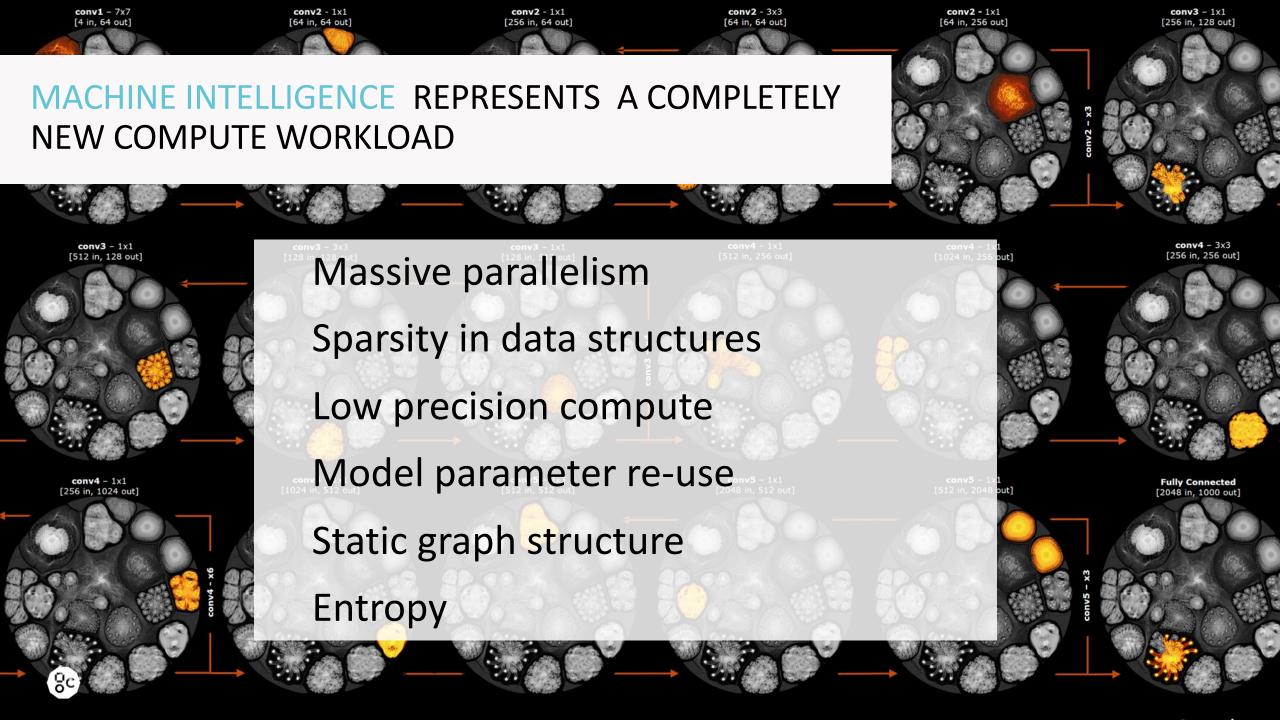


M2000 and Server IPU-POD₆₄ scale-out



IPU ARCHITECTURE OVERVIEW

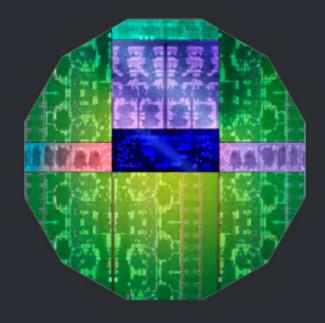




LEGACY PROCESSOR ARCHITECTURES HAVE BEEN REPURPOSED FOR ML



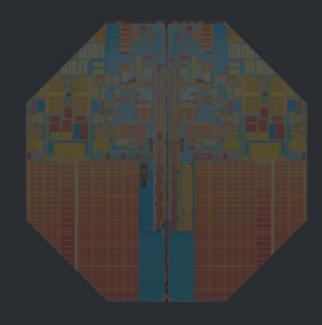
CPUApps and Web/
Scalar



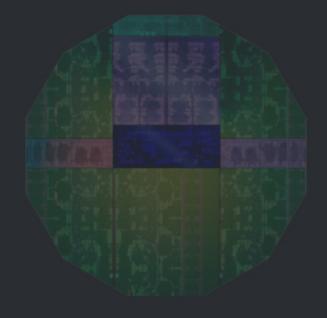
GPUGraphics and HPC/
Vector



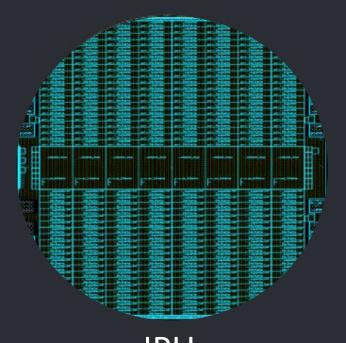
A NEW PROCESSOR IS REQUIRED FOR THE FUTURE



CPUApps and Web/



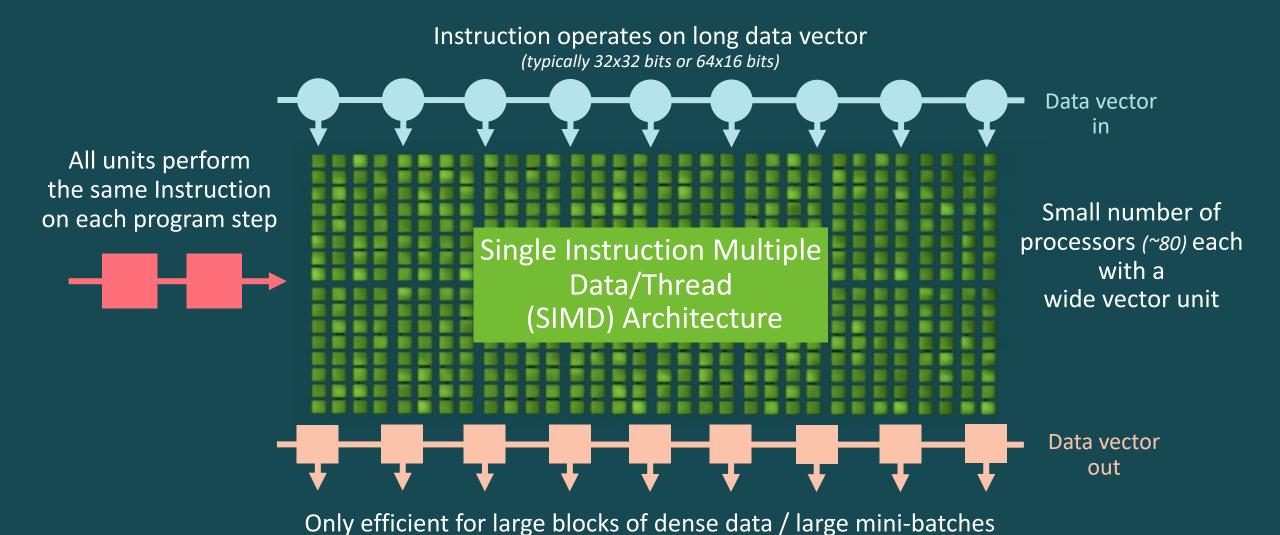
GPUGraphics and HPC/
Vector



IPU
Artificial Intelligence/
Graph

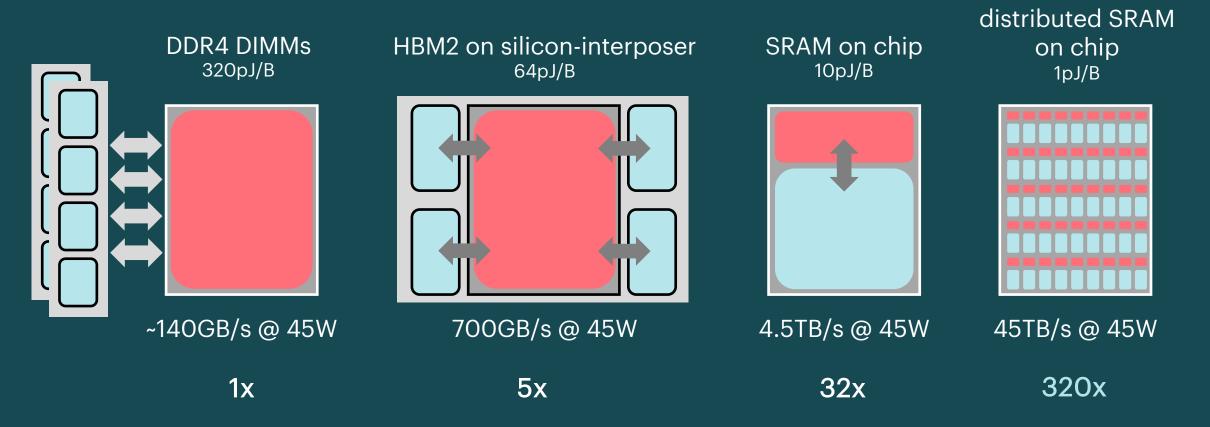


TODAY'S PARALLEL MACHINES ARE INFLEXIBLE





MEMORY TRADE-OFF





MASSIVE PARALLELISM WITH ULTRAFAST MEMORY ACCESS

CPU GPU IPU

Parallelism

Designed for scalar processes

SIMD/SIMT architecture.

Designed for large blocks of dense contiguous data

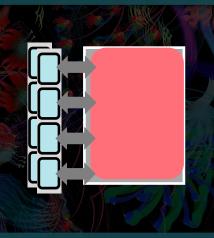
Massively parallel MIMD.

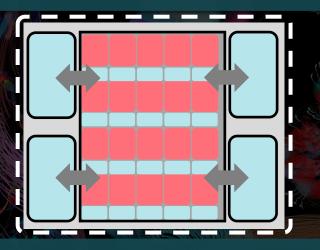
High performance/efficiency for future ML

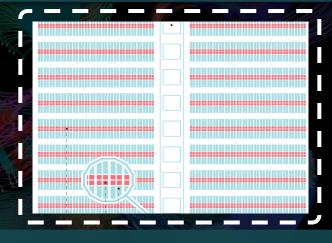
trends

Processor

Memory







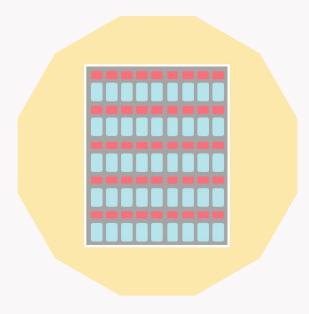
Memory Access

Off-chip memory Model and Data spread across off-chip and small on-chip cache and shared mem.

Model & Data in tightly coupled large locally distributed SRAM

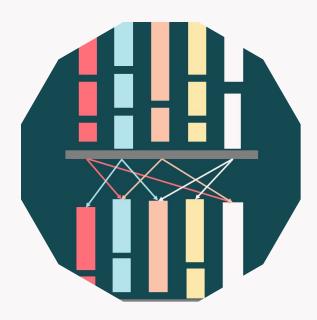
IPU PROVIDES UNIQUE SPEED AND FLEXIBILITY

Code and data are always in SRAM



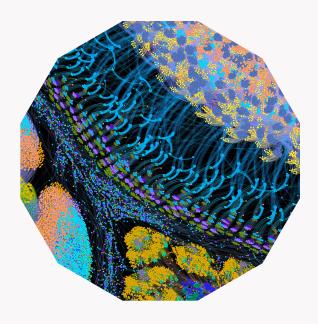
45 TB/s memory bandwidth

Parallel programs are truly independent



7,296 independent threads

Native support for computational graphs

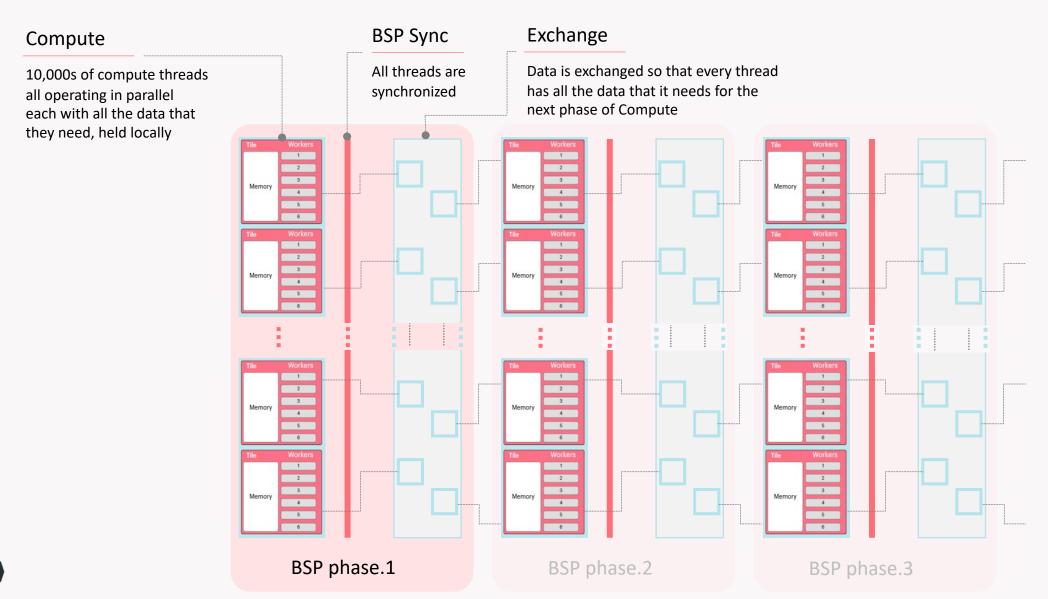


Poplar[®] SDK



BULK SYNCHRONOUS PARALLEL (BSP)

Software bridging model for parallel computing





COLOSSUS MK2

the worlds most complex processor

59.4Bn transistors, TSMC 7nm @ 823mm²

250TFlops AI-Float | 900MB In-Processor-Memory™

1472 independent processor cores

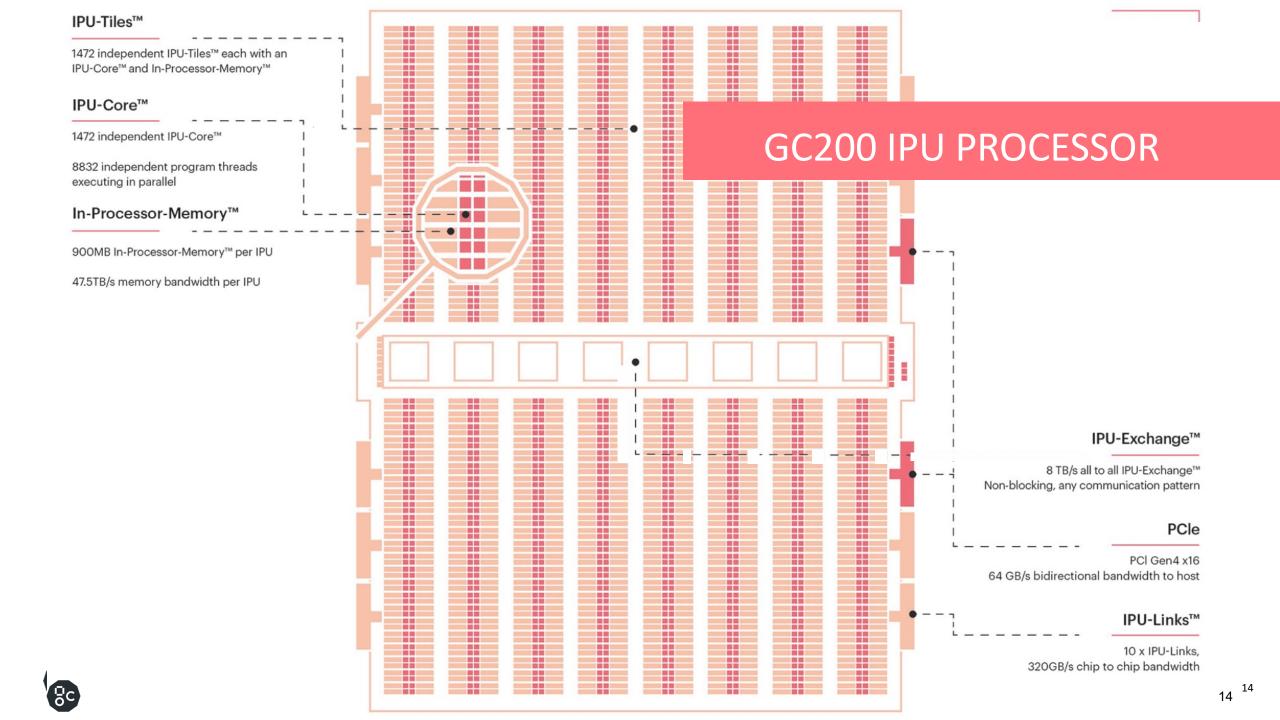
8832 separate parallel threads

>8x step-up in system performance vs Mk1



GC200 IPU





IPU POD SYSTEMS

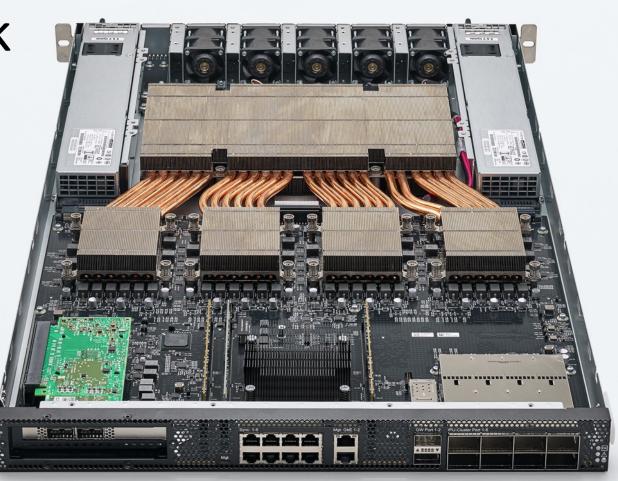


IPU-M2000 PODs Building Block

1U Machine Intelligence Compute Blade

1 PetaFlop IPU compute

2.8Tbps IPU-Fabric™



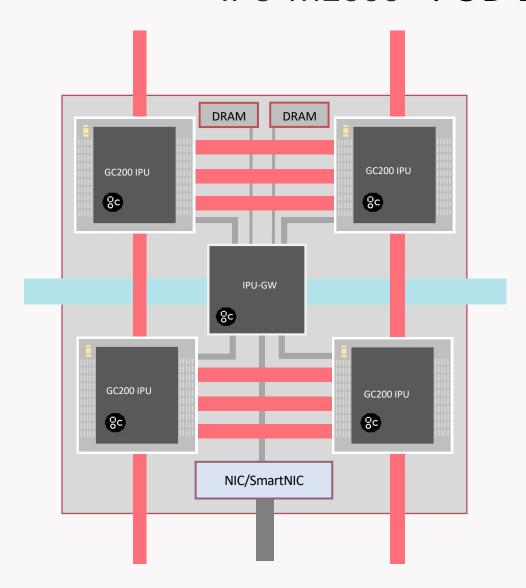
DATA

COMMUNICATIONS

COMPUTE



IPU-M2000 - POD BUILDING BLOCK





4x GC200 IPUs

- 1 PFLOP_{16.16} compute
- 5,888 processor cores
- > 35,000 independent parallel threads



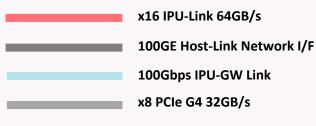
Up to ~450GB Exchange Memory™

- Up to 448GB Streaming Memory[™] DRAM
- 3.6GB In-Processor Memory TM

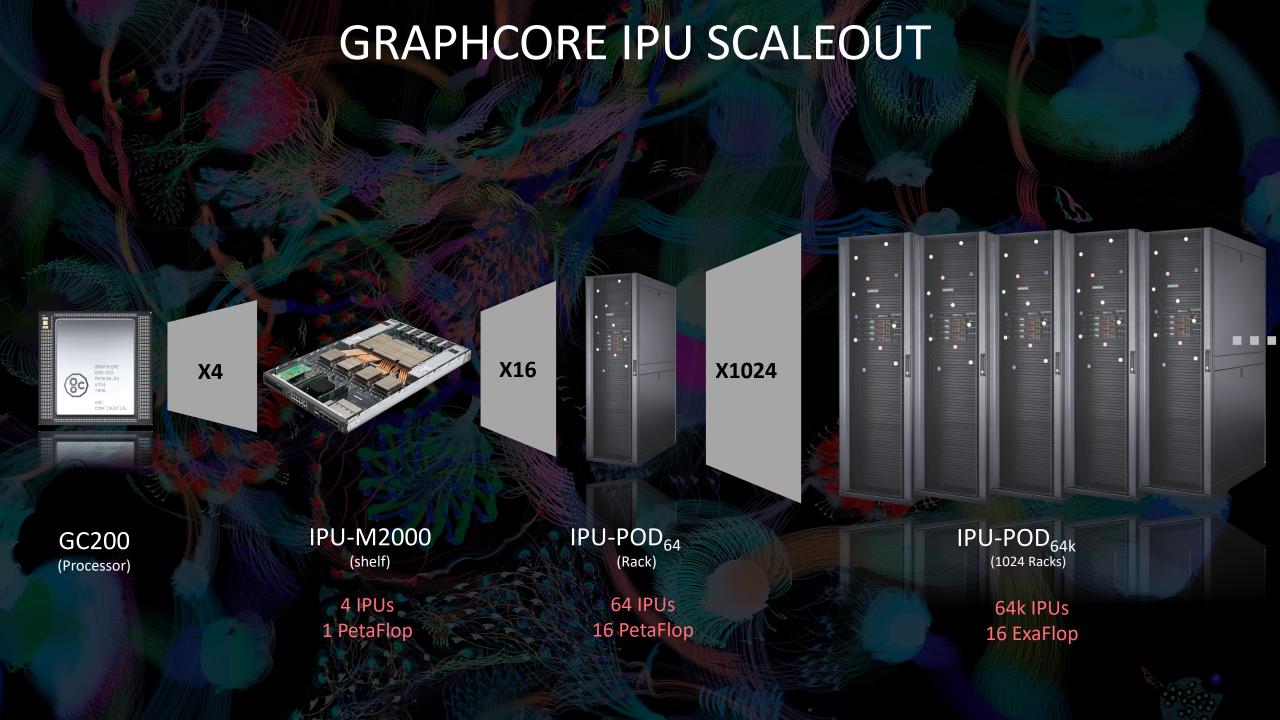


IPU-Fabric for Compiled-In Networking

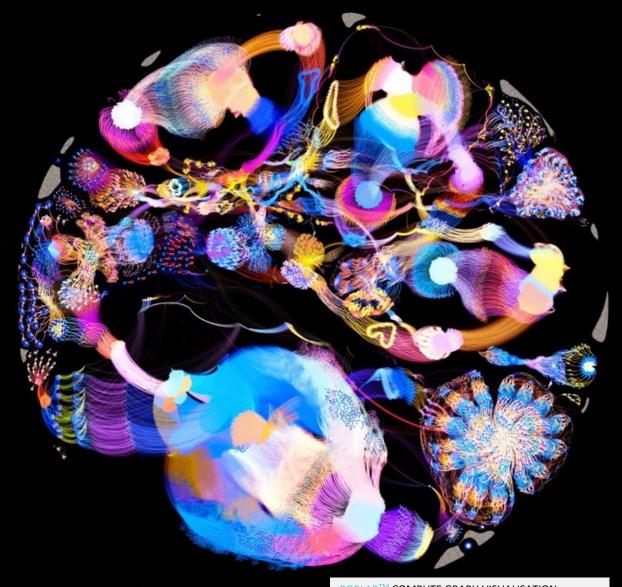
- Host-Link 100GE link to Poplar Server
- IPU-Link 4x 512Gbps for intra IPU-POD64 communication
- GW-Link 2x 100bps Gateway-Links for inter IPU-POD64 communication (0.4 Tbps bidir.)
- Sync-Link IPU-POD hardware sync signal







SOFTWARE



POPLAR™ COMPUTE GRAPH VISUALISATION

POPLAR EASE OF USE

















Open & Extensible Poplar Libraries

Get access to 50+ optimised functions for common ML models and 750 high performance compute elements. Modify and write custom libraries.

ML Frameworks Support

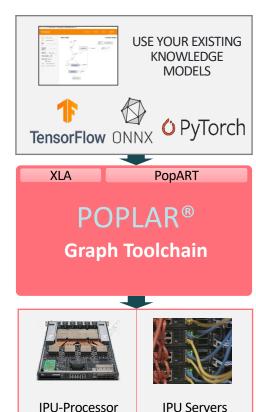
Support for standard ML frameworks: TensorFlow 1 and 2, ONNX and PyTorch with PaddlePaddle coming soon.

Straightforward Deployment

Pre-built Docker containers with Poplar SDK, Tools and frameworks images to get up and running fast.

Standard Ecosystem Support

Ready for production with Microsoft Azure deployment, Kubernetes orchestration, Docker containers and Hyper-V virtualization & security.



and System

Platforms



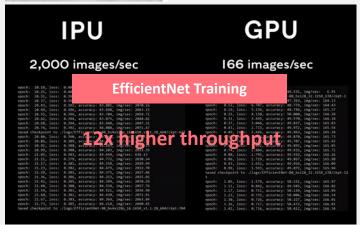
IPU ADVANTAGE



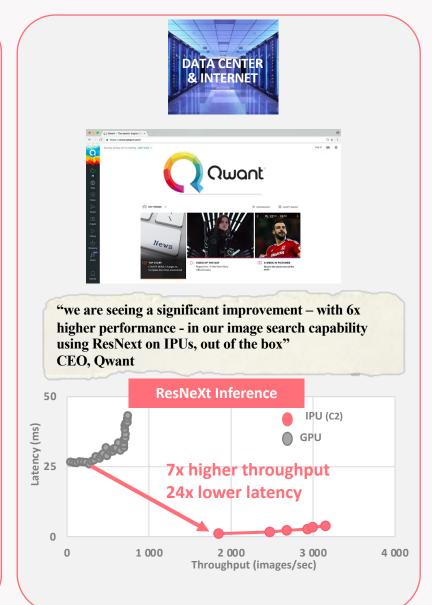




Analysis of Chest X-Ray images using variant of EfficientNet-B0 Training for COVID-19/pneumonia classification model development



Images & Work presented by
Microsoft AI & Advanced Architectures
Research Lead, Sujeeth Bharadwaj, PhD at Intelligent
Health Inspired conference May 2020

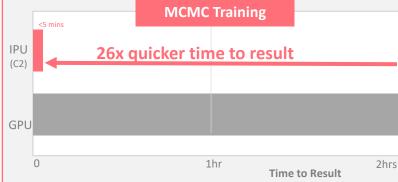








"We were able to train one of our proprietary models in 4.5minutes instead of 2hours. That's 26x faster time to train than other leading platforms" George Sokoloff, Founder & CIO Carmot Capital





IPU FOR RESEARCH



IMPERIAL COLLEGE LONDON ACCELERATE CLASSICAL COMPUTER VISION PROBLEM ON IPU



UNIVERSITY OF BRISTOL SOLVES SCIENTIFIC PROBLEMS WITH NEW IPU-BASED AI SYSTEMS

Using the Graphcore IPU for traditional HPC applications

Dorben Louw, Simon McIntooh-Smith Days: of Computer Science Debursity of British Brisisl, United Kingdom

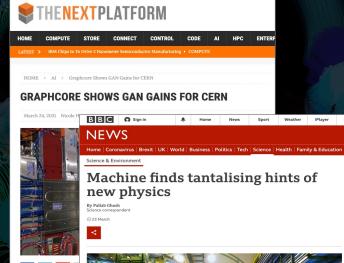
Advance—In increase in marsine harming workflow means and the contraction of the contract

The state of the s

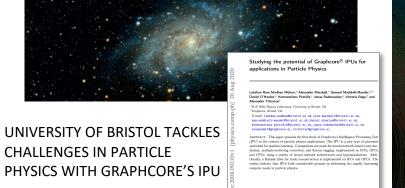
And the second s

UNIVERSITÉ DE PARIS ACCELERATES COSMOLOGY WITH IPUS

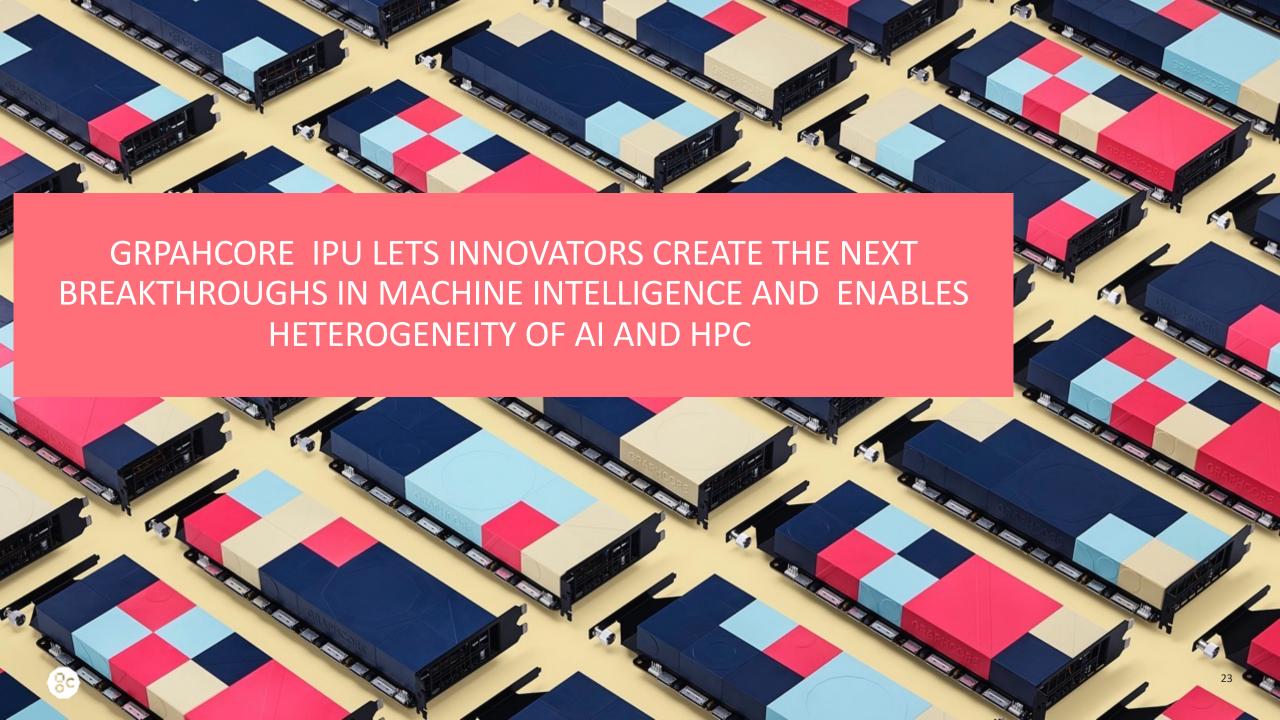




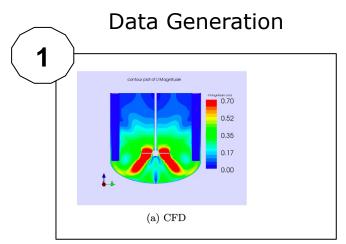


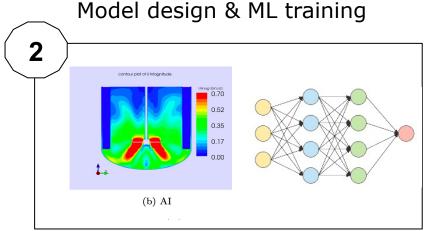


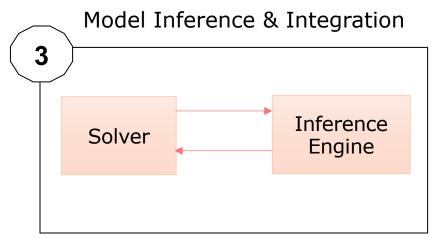


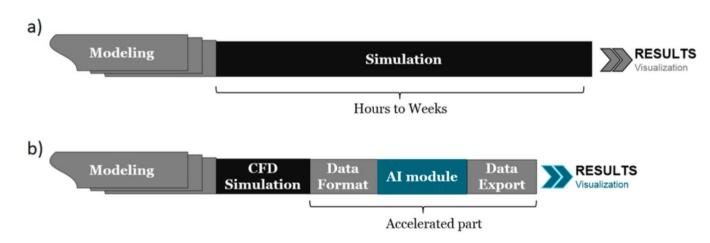


Acceleration of Scientific Applications and Enabling Heterogeneity of Al and HPC to Improve Simulation Time









Areas where this approach makes sense and where it has already been proven

- HEP -High Energy Physics
- CFD -Computational Fluid Dynamics
- PDE -Partial Differential Equation
- PF Protein folding
- MC Monte Carlo simulation
- WRF Weather Research and Forecasting
- O&G Oil and Gas Exploration Simulation



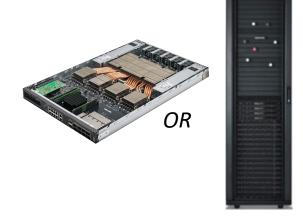
Hybridization of HPC and AI systems

For typical HPC workloads



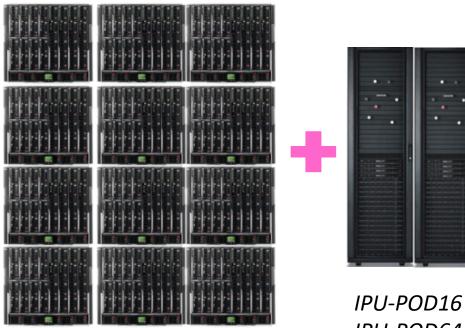
CPUs and GPUs type of configurations

For AI workloads



IPU-POD16 IPU-POD64 IPU-POD512

For new heterogeneous types of workloads



CPUs and GPUs type of configurations

IPU-POD16
IPU-POD64
IPU-POD512

The New Way of Performance Acceleration for Exascale Systems

