# SCALABLE MACHINE INTELLIGENCE SYSTEMS



# GRAPHCORE

# SCALABLE IPU SYSTEMS

## SILICON

A MULTI-GENERATIONAL SILICON ARCHITECTURE PROVIDING
OPTIMIZED SUPPORT FOR HIGH PERFORMANCE MACHINE
INTELLIGENCE APPLICATIONS AND WORKLOADS AT SCALE

## PLATFORMS

HARDWARE PLATFORMS DESIGNED TO DEPLOY IPU DEVICES WHICH
ENABLE OPTIMIZED APPLICATIONS TO EXECUTE EFFICIENTLY IN
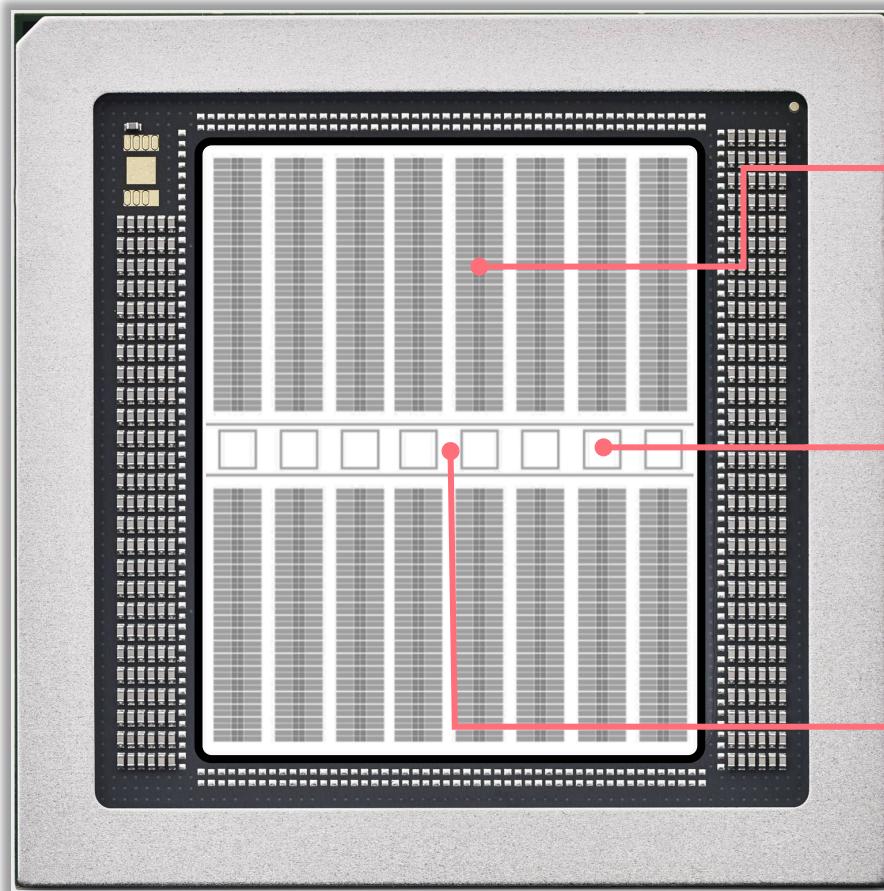SUPPORT OF INDUSTRY STANDARD DEPLOYMENT USE CASES

## SOFTWARE

SUPPORT FOR COMMON MACHINE INTELLIGENCE DEVELOPMENT
FRAMEWORKS AND DIRECT IPU PROGRAMMING THAT ENABLES
DEVELOPERS TO SEAMLESSLY INTEGRATE IPUS INTO APPLICATIONS

GRAPHCORE

# SCALING ON DEVICE



**COMPUTE**
- 1216 FULLY PROGRAMMABLE IPU COMPUTE TILES
- EACH WITH 256KB IN-PROCESSOR MEMORY™
- ACCELERATED FLOATING POINT HARDWARE ENGINES
- 6 INDEPENDENT HARDWARE WORKERS PER TILE

**EXCHANGE**
- 8 TB/S TILE TO TILE COMMUNICATION
- SUPPORTS ANY COMMUNICATION PATTERN
- FULLY TIMING DETERMINISTIC EXECUTION
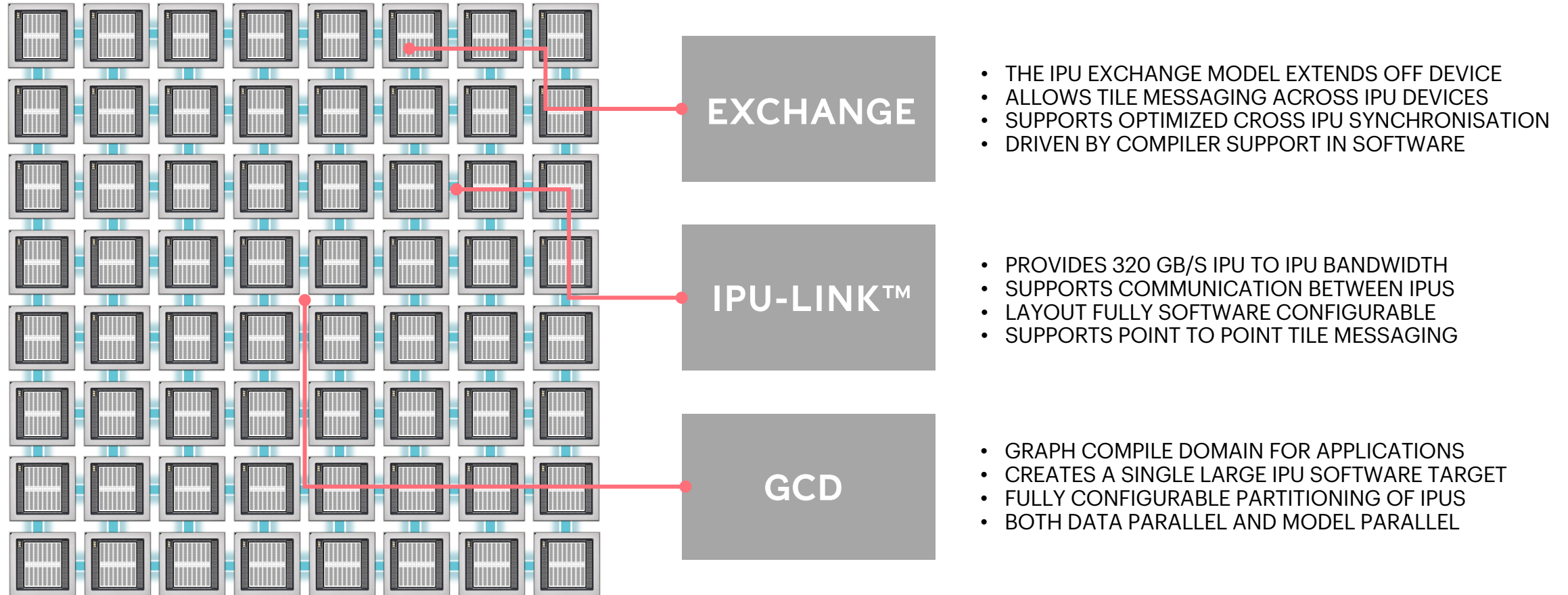- DRIVEN BY COMPILER SUPPORT IN SOFTWARE

**SYNC**
- HARDWARE ACCELERATED TILE SYNCHRONISATION
- LOW LATENCY MESSAGING OF WORK COMPLETION
- USED TO MOVE BETWEEN COMPUTE AND EXCHANGE
- EXTENDED INSTRUCTION SET FOR SYNCHRONISATION

**FIRST GENERATION IPU – COLOSSUS MK1**

23.6 BILLION TRANSISTORS, 7296 FULLY INDEPENDENT WORKERS, 45 TB/S MEMORY BANDWIDTH

300MB IN-PROCESSOR MEMORY™, PCIe GEN4 INTERFACE, IPU-LINK™ INTERFACE
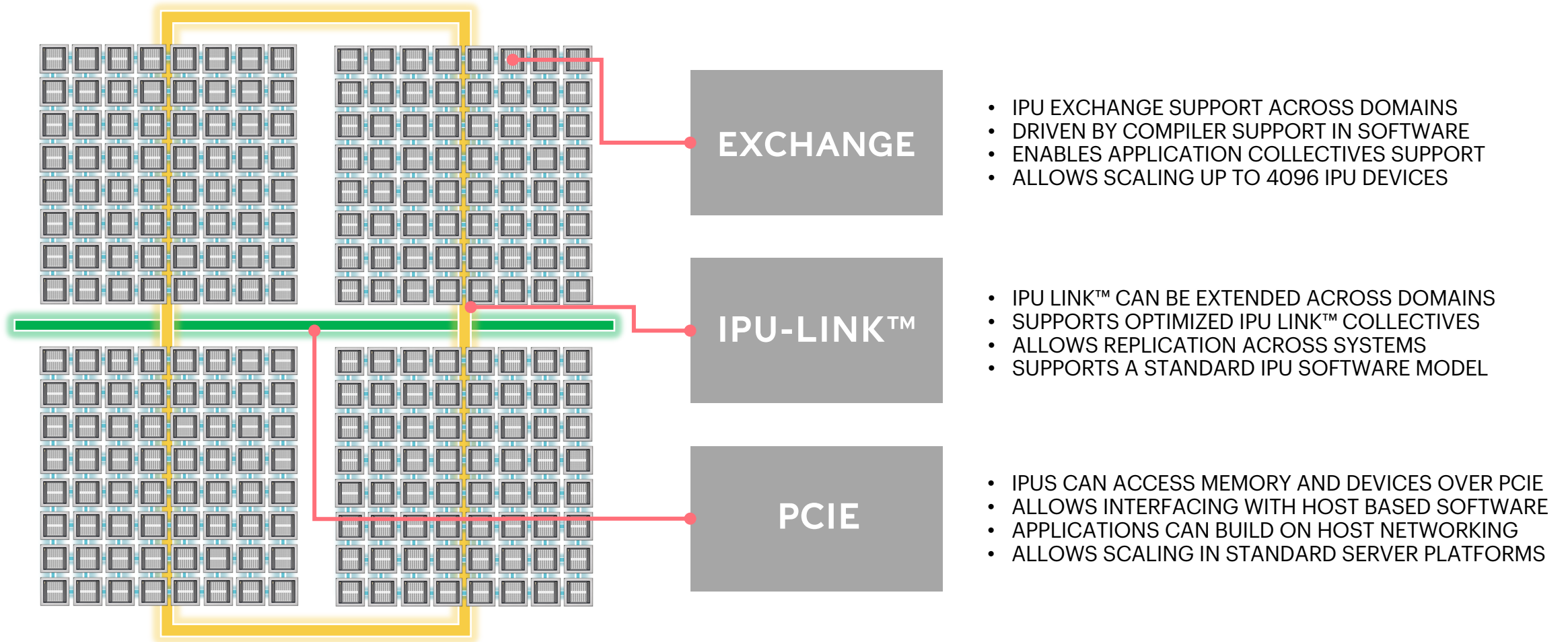
**GRAPHCORE**

# SCALING ACROSS DEVICES

**EXCHANGE**

- THE IPU EXCHANGE MODEL EXTENDS OFF DEVICE
- ALLOWS TILE MESSAGING ACROSS IPU DEVICES
- SUPPORTS OPTIMIZED CROSS IPU SYNCHRONISATION
- DRIVEN BY COMPILER SUPPORT IN SOFTWARE

**IPU-LINK™**

- PROVIDES 320 GB/S IPU TO IPU BANDWIDTH
- SUPPORTS COMMUNICATION BETWEEN IPUS
- LAYOUT FULLY SOFTWARE CONFIGURABLE
- SUPPORTS POINT TO POINT TILE MESSAGING

**GCD**

- GRAPH COMPILE DOMAIN FOR APPLICATIONS
- CREATES A SINGLE LARGE IPU SOFTWARE TARGET
- FULLY CONFIGURABLE PARTITIONING OF IPUS
- BOTH DATA PARALLEL AND MODEL PARALLEL

**UP TO 64 IPU DEVICES USABLE AS A SINGLE LARGE IPU FROM APPLICATIONS**

466944 FULLY INDEPENDENT WORKERS, 19.2GB IN-PROCESSOR MEMORY™, LEVERAGING OVER 1.5 TRILLION TRANSISTORS

GRAPHCORE

# SCALING ACROSS SYSTEMS



**EXCHANGE**

- IPU EXCHANGE SUPPORT ACROSS DOMAINS
- DRIVEN BY COMPILER SUPPORT IN SOFTWARE
- ENABLES APPLICATION COLLECTIVES SUPPORT
- ALLOWS SCALING UP TO 4096 IPU DEVICES

**IPU-LINK™**

- IPU LINK™ CAN BE EXTENDED ACROSS DOMAINS
- SUPPORTS OPTIMIZED IPU LINK™ COLLECTIVES
- ALLOWS REPLICATION ACROSS SYSTEMS
- SUPPORTS A STANDARD IPU SOFTWARE MODEL

**PCIE**

- IPUS CAN ACCESS MEMORY AND DEVICES OVER PCIE
- ALLOWS INTERFACING WITH HOST BASED SOFTWARE
- APPLICATIONS CAN BUILD ON HOST NETWORKING
- ALLOWS SCALING IN STANDARD SERVER PLATFORMS

**256 IPU APPLICATION TARGET BUILT FROM INTERCONNECTED 64 IPU DOMAINS**

SCALE OUT SUPPORT UP TO A MAXIMUM OF 4096 IPUS WITH FIRST GENERATION COLOSSUS MK1

**GRAPHCORE**

# EXECUTION MODEL

## COMPUTATIONAL GRAPH

**data**

data [0]  data [1]  data [2]  data [3]

AdderVertex [v0]  AdderVertex [v1]

Compute Set 0

**output**

output [0]  output [1]

AdderVertex [v2]

Compute Set 1

**result**

result [0]

## BSP SCHEDULE

EXCHANGE

SYNC

COMPUTE

SYNC

EXCHANGE

SYNC

COMPUTE

SYNC

EXCHANGE

## OPTIMIZED IPU EXECUTION

BSP EXECUTION TRACE - IPU TILES 0 - 1215



OUTPUT FROM POPVISION GRAPH ANALYSER

# GRAPHCORE

# POPLAR® SDK



MACHINE LEARNING FRAMEWORKS

PyTorch

PaddlePaddle

ONNX

TensorFlow

POPLAR ADVANCED RUNTIME

XLA

POPLAR®

POPNN

POPLIN

POPOPS

POPRAND

POPUTIL

GCL

POPLAR STANDARD LIBRARIES

POPLAR GRAPH LIBRARY

GRAPH COMPILER

GRAPH ENGINE

POPLAR DEVICE INTERFACE

DRIVERS

GRAPHCORE DEVICE ACCESS

IPU HARDWARE ABSTRATION LAYER

USER SPACE API

PCIe DRIVER

IPUoF DRIVER

GRAPHCORE

# MULTI-IPU CONSTRUCTS

## MODEL SHARDING



SUPPORT THE SPLITTING OF MODELS
ACROSS MULTIPLE IPU DEVICES

ALLOW USER DRIVEN SOFTWARE
CONTROL OF MODEL PARALLELISM

## MODEL REPLICATION



SUPPORT THE REPLICATION OF MODELS
ACROSS AN ENTIRE IPU SYSTEM

ENABLE DATA PARALLEL TRAINING AND
AUTOMATIC REPLICATION OF MODELS

## MODEL PIPELINING



SUPPORT THE PIPELINING OF MODELS
ACROSS MULTIPLE IPU DEVICES

EXTRACT MAXIMUM PERFORMANCE FOR
MODEL PARALLEL EXECUTION

FULLY SUPPORTED IN PYTORCH, TENSORFLOW, POPART™ AND POPLAR®

GRAPHCORE

# POPLAR®

## GRAPH LIBRARY

- SIMPLE C++ GRAPH BUILDING API
- IMPLEMENT ANY APPLICATION
- FULL CONTROL FLOW SUPPORT



USER DEFINED COMPUTATIONAL GRAPH

## GRAPH COMPILER

- OPTIMIZING IPU GRAPH COMPILER
- IMPLEMENTS IPU EXECUTION MODEL
- CODE GENERATION USING LLVM



GRAPH COMPILER INTERMEDIATE REPRESENTATION

## GRAPH ENGINE

- HIGH PERFORMANCE GRAPH RUNTIME
- INTERFACES TO HOST MEMORY SYSTEM
- HIGHLY OPTIMIZED IPU DATA TRANSFER



OPTIMIZED GRAPH EXECUTION WITH HOST SYSTEM

---

**WWW.GRAPHCORE.AI/DEVELOPER**



COMPREHENSIVE POPLAR USER GUIDE



UISING MULTIPLE IPUS IN APPLICATIONS



TUTORIALS AND CODE EXAMPLES



DETAILED DESCRIPTIONS OF UNDERLYING CONCEPTS

**GRAPHCORE**

# POPLAR® LIBRARIES

- OVER 50 OPTIMISED FUNCTIONS FOR COMMON ML MODELS
- MORE THAN 750 HIGH PERFORMANCE COMPUTE ELEMENTS

**POPNN** — FUNCTIONS USED IN NEURAL NETWORKS (NON-LINEARITIES, POOLING, LOSS FUNCTIONS)

**POPLIN** — OPTIMIZED LINEAR ALGEBRA FUNCTIONS (MATRIX MULTIPLICATION, CONVOLUTIONS)

**POPOPS** — FUNCTIONS FOR PERFORMING ELEMENTWISE OPERATIONS ON TENSOR DATA

**POPRAND** — HIGH PERFORMANCE FUNCTIONS FOR POPULATING TENSORS WITH RANDOM NUMBERS

**POPUTIL** — GENERAL UTILITY FUNCTIONS FOR BUILDING GRAPHS FOR IPU DEVICES

**GCL** — OPTIMIZED COLLECTIVES LIBRARY SUPPORTING MODEL AND DATA PARALLEL

OPTIMIZED WORK PLANNING OF FUNCTIONS ACROSS IPU DEVICES



POPLAR GRAPH COMPILER INTERMEDIATE REPRESENTATION FOR MATRIX MULTIPLICATION OPERATIONS, CUSTOM GRAPH LAYOUT SPECIALISED BASED ON SHAPES OF INPUTS AND OUTPUTS

GRAPHCORE



WWW.GRAPHCORE.AI/DEVELOPER

# MACHINE LEARNING FRAMEWORKS

## TensorFlow

- TENSORFLOW SUPPORT FOR IPU AS FAMILIAR TARGET FOR MODELS
- FULL PERFORMANT INTEGRATION WITH TENSORFLOW XLA BACKEND
- SUPPORT FOR VERSION 1 & 2 WITH EXAMPLES AND DOCUMENTATION

## PyTorch

- PYTORCH SUPPORT FOR TARGETING IPU WITH SIMPLE EXTENSIONS
- TAKE NATIVE PYTORCH MODELS, DEPLOY AND TRAIN ON IPU DEVICES
- SUPPORT FOR MULTI-IPU PRIMATIVES FROM PYTORCH MODELS

## POPART™

- THE POPLAR ADVANCED RUNTIME FOR INFERENCE AND TRAINING
- SUPPORTS ONNX MODEL INPUT AND PYTHON / C++ MODEL BUILDING
- AN OPTIMIZED LIGHTWEIGHT APPLICATION RUNTIME FOR DEPLOYMENT

GRAPHCORE

# SOFTWARE ECOSYSTEM

DEPLOYMENT

ORCHESTRATION

CONTAINERS

VIRTUALIZATION
& SECURITY

STANDARD
PCI-E DRIVER

DELL EMC DSS8440

GRAPHCORE

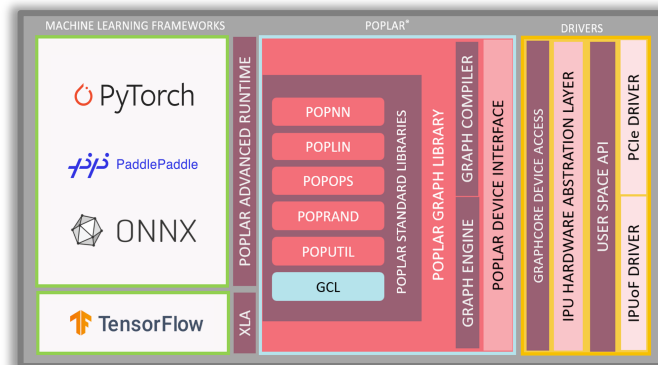WWW.GRAPHCORE.AI/DEVELOPER

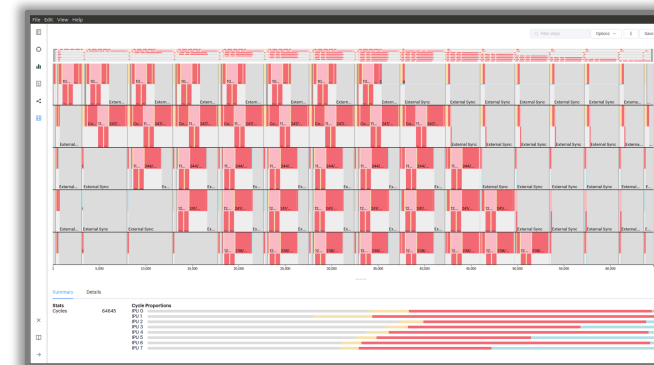# DEVELOPERS

CLOSE TO 1000 USERS SIGNED UP TO WORK WITH GRAPHCORE TECHNOLOGY

## POPLAR SDK



## POPVISION TOOLS



## DOCUMENTATION



## APPLICATION EXAMPLES



GRAPHCORE

# SUMMARY



Machine Learning, Machine Intelligence, AI, Machine learning hardware, Baidu

**Graphcore joins Baidu PaddlePaddle Hardware Ecosystem**

Poplar, Software, Developer

**Major new Graphcore Poplar® SDK 1.1. Released**

Machine Intelligence, AI, Memory, Poplar, IPU, Software, Developer

**Graphcore launches new Poplar® Analysis Tool**

Machine Intelligence, AI, IPU, Computer Vision, Robotics

**Imperial College London accelerate Classical Computer Vision Problem on IPU**

Machine Learning, Machine Intelligence, AI, Graphcore, Poplar

**Graphcore makes Poplar SDK Docs publicly available**

Machine Intelligence, AI, Machine learning hardware, microsoft, ResNeXt, IPU, Healthcare

**Microsoft accelerates ResNeXt-50 Medical Imaging Inference on the IPU**

Machine Intelligence, AI, Machine learning hardware, ResNeXt, C2 Card, IPU, Microsoft Azure

**Qwant publishes New Paper evaluating IPU Performance for Image-Based Deep Learning**
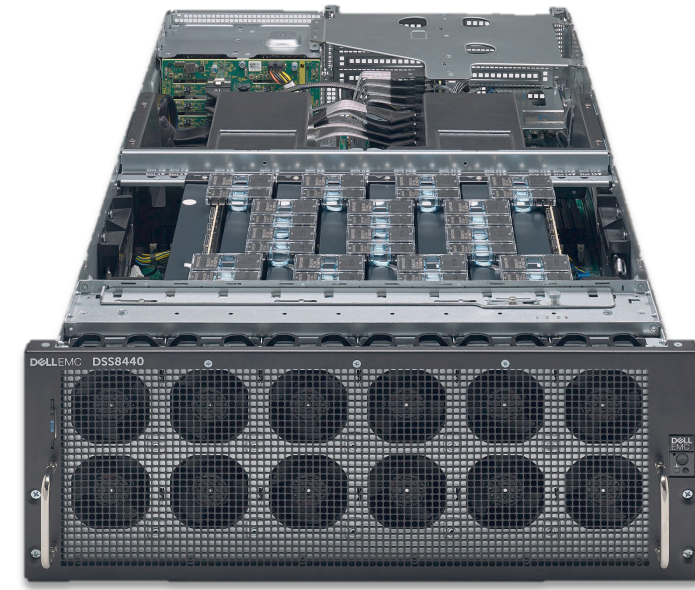
- RAPID PROGRESS WITH PLATFORM CAPABILITY

- CLOSE TO 1000 USERS AND GROWING FAST

- POPLAR SDK CONTINUES TO EVOLVE AT PACE

- DEVELOPER PORTAL MAKES USING IPU SIMPLE

- TRY IPU TODAY IN MICROSOFT AZURE CLOUD

GRAPHCORE

GRAPHCORE

WWW.GRAPHCORE.AI

BERT BASE POPART MODEL – POPLAR GRAPH COMPILER IR VISUALISATION