Graphcore IPUs: Accelerating Argonne's AI/ML Applications



TOPICS

- **Graphcore Introduction** \bullet
 - Company ullet
 - **Processor Architecture (IPU)** •
 - System Architecture •
 - Software & Enablement ullet
- **General Benchmarks** \bullet
- **Customer Case Studies** \bullet
- **Argonne Applications** \bullet
 - UNO \bullet
 - BraggNN ullet
- "Good" Computer \bullet





GRAPHCORE





GRAPHCORE ENABLING MACHINE INTELLIGENCE



- Founded in 2016
- Technology: Intelligence Processor Unit (IPU)
- Team: 650+ globally
- Offices: UK, US, China, Norway, Poland
- Raised >\$730M



GRAPHCORE



ABOUT US...

Hardware



IPU processors designed for AI

Software



Poplar[®] software stack & development tools

Systems



BOW-2000 and Server IPU-POD₆₄ scale-out



PROCESSOR ARCHITECTURE





INTRODUCING THE BOW IPU WORLD'S FIRST 3D WAFER-ON-WAFER PROCESSOR



3D silicon wafer stacked processor **350 TeraFLOPS AI compute Optimized** silicon power delivery 0.9 GigaByte In-Processor-Memory @ 65TB/s 1,472 independent processor cores 8,832 independent parallel programs 10x IPU-Links[™] delivering 320GB/s



65.4TB/s memory bandwidth per IPU

BULK SYNCHRONOUS PARALLEL (BSP)

"Bridging Model" for parallel computation

Repeat Supercycles composed of 3 phases

- Compute execution from local memory
- Sync global barrier synchronization
- Exchange global communication

Graphcore: Sync, Exchange, Compute

- Globally Exchange immediately follows Sync
- Each Tile Compute immediately follows Exchange





SYSTEM ARCHITECTURE





BOW-2000 IPU MACHINE

4 x Bow 3D Wafer-on-Wafer IPUs

1.4 PetaFLOPS AI Compute

Up to 256 GB IPU Streaming Memory

2.8 Tbps IPU-Fabric™

3.6 GB In-Processor-Memory @ 260TB/s

1U blade form factor



GRAPHCORE

BOW-2000 IPU MACHINE



IPU-FABRIC FOR SCALE-OUT



BOW-2000/IPU-POD RECONFIGURABILITY

Reusable, simple migration between platforms



BOW: 3RD GENERATION IPU SYSTEMS SHIPPING TO CUSTOMERS TODAY







4x Bow-2000 5.6 PetaFLOPS 1 CPU server

BOW POD₃₂

8x Bow-2000 11.2 PetaFLOPS 1 CPU server

BOW POD₆₄

16x Bow-2000 22.4 PetaFLOPS 1-4 CPU server(s)

BOW POD₂₅₆

64x Bow-2000 **89.6 PetaFLOPS** 4-16 CPU server(s)



BOW POD₁₀₂₄

256x Bow-2000 **358.4 PetaFLOPS** 16 - 64 CPU server(s) Early access



SOFTWARE & ENABLEMENT





GRAPHCORE SOFTWARE MATURITY



GRAPHCORE

USEFUL SW CAPABILITIES



of the Darmon 1 Bentles & Row Hillord I I least have 10.061.641 101,449,313 ALC: UNK 101, 101, 101 1.04.00 Million in Cost, Name 194,186,88 W. State, State Let Use PL-10, IN 1.616.000 14, 20, 10 6,416,416 81,81,00 1.04,014 10.0 N.ST.A Non-Advanced and 344,424 Transmission (Collins) (412,254 16,458,520 1.44.368 10,00,00 1.44,546 78,811.mm 1.62,200 10,011.0 4.1.1 1014 1114 41. STATE NAME AND ADDRESS OF TAXABLE PARTY. 41.5% 141.66 14.6 No. 84 14 50

-

Debug Tools





✿ gpu_cnn_keras.py ↔ ipu_cnn_keras.py tf_keras					
and the second se	import tensorflow as tf				
Crac ras.layers import *	from tensorflow.keras.layers import *				
ICIAS ////////////////////////////////////	+ from tensorflow.python import ipu				
	+ cfg = ipu.config.IPUConfig()				
	+ cfg.auto_select_ipus = 1				
	+ cfg.configure_ipu_system()				
(x train x train) (x test x test) = tf keres detects sifer10 lead dets()	+ with ipu.ipu_strategy.iPustrategy().scope():				
$(x_train, y_train), (x_test, y_test) = (1.keras.datasets.cirario.toad_data())$	$(x_train, y_train), (x_test, y_test) = (1.keras.uatasets.cirario.toau_uata())$				
$x_{train} = x_{train.astype}(1000152) / 255.0$	$x_{train} = x_{train.astype}(1000000000000000000000000000000000000$				
ds train = tf data Dataset from tensor slices((x train, y train)) batch(64, dron remainds)	ds train = tf data Dataset from tensor slices((v train, v train)) batch(64, drop rema				
us_train = tridata.bataset.from_tensor_stites((x_train, y_train)).batch(04, urop_remainde					
<pre>model = tf.keras.Sequential([</pre>	<pre>model = tf.keras.Sequential([</pre>				
<pre>Conv2D(32, (3, 3), padding='same', input_shape=x_train.shape[1:]),</pre>	<pre>Conv2D(32, (3, 3), padding='same', input shape=x train.shape[1:]),</pre>				
Activation('relu'),	Activation('relu'),				
Conv2D(32, (3, 3)),	Conv2D(32, (3, 3)),				
Activation('relu'),	Activation('relu'),				
<pre>MaxPooling2D(pool_size=(2, 2)),</pre>	<pre>MaxPooling2D(pool_size=(2, 2)),</pre>				
Dropout(0.25),	Dropout(0.25),				
Conv2D(64, (3, 3), padding='same'),	Conv2D(64, (3, 3), padding='same'),				
Activation('relu'),	Activation('relu'),				
Conv2D(32, (3, 3)),	Conv2D(32, (3, 3)),				
Activation('relu'),	Activation('relu'),				
<pre>MaxPooling2D(pool_size=(2, 2)),</pre>	<pre>MaxPooling2D(pool_size=(2, 2)),</pre>				
Dropout(0.25),	Dropout(0.25),				
Flatten(),	Flatten(),				
Dense(512),	Dense(512),				
Activation('relu'),	Activation('relu'),				
Dense(10)	Dropout(0.5),				
Activation('softmax')	Activation ('softmax')				
<pre>model.compile(loss='categorical_crossentropy',</pre>	<pre>model.compile(loss='categorical_crossentropy',</pre>				
<pre>optimizer=tf.optimizers.SGD(learning_rate=0.016),</pre>	<pre>optimizer=tf.optimizers.SGD(learning_rate=0.016),</pre>				
<pre>metrics=['accuracy'])</pre>	<pre>metrics=['accuracy'])</pre>				
<pre>model.fit(ds_train, epochs=40)</pre>	<pre>model.fit(ds_train, epochs=40)</pre>				

.

MODEL GARDEN ML DOMAIN COVERAGE



GENERAL BENCHMARKS





PERFORMANCE RESULTS MLPERF & BENCHMARKS

We publish our MLPerf & benchmark performance on multiple models, platforms, & frameworks

- Published performance results with config/parameters
- > Allows customers to do their own analysis/comparison

Model	Variant	Platform	SDK Version	Framework	Dataset	Batch Size	Precision	Throughput (items/sec)
ResNet-50 v1.5		Bow Pod16	SDK2.6	TensorFlow1	ImageNet2012	3,520	16.16	44,047
ResNet-50 v1.5		Bow Pod16	SDK2.6	PyTorch	ImageNet2012	16,384	16.16	37,219
ResNet-50 v1.5		Bow Pod64	SDK2.6	TensorFlow1	ImageNet2012	5,120	16.16	152,522
ResNet-50 v1.5		Bow Pod128	SDK2.6	TensorFlow1	ImageNet2012	5,120	16.16	0
ResNet-50 v1.5		Bow Pod256	SDK2.6	TensorFlow1	ImageNet2012	10,240	16.16	450,701
EfficientNet-B4	G16-EfficientNet	Bow Pod16	SDK2.6	TensorFlow1	ImageNet2012	6,144	16.16	8,992
EfficientNet-B4	G16-EfficientNet	Bow Pod16	SDK2.6	PyTorch	ImageNet2012	1,024	16.32	6,111
EfficientNet-B4	G16-EfficientNet	Bow Pod64	SDK2.6	TensorFlow1	ImageNet2012	6,144	16.16	35,247
EfficientNet-B4	G16-EfficientNet	Bow Pod256	SDK2.6	TensorFlow1	ImageNet2012	6,144	16.16	117,240
ResNeXt101		Bow Pod16	SDK2.6	TensorFlow1	ImageNet2012	768	16.16	12,233
VIT	Pre-Training	Bow Pod16	SDK2.6	PyTorch	ImageNet1k	65,536	16.16	7,512
VIT	Pre-Training	Bow Pod64	SDK2.6	PyTorch	ImageNet1k	65,536	16.16	27,981
ViT	Fine-Tuning	Bow Pod16	SDK2.6	PyTorch	ImageNet1k	2,040	16.16	7,977



GRAPHCORE

https://www.graphcore.ai/performance-results

BOW POD ADVANTAGE

SPEED-UP OVER DGX-AIOO ACROSS RANGE OF MODELS | SCALING TO LARGER POD SYSTEMS



Bow Pod Platforms | SDK2.6 | ResNet-50 v1.5 Training Throughput TensorFlow | G16-EfficientNet-B4 TTT Tensorflow | BERT Large Ph1 Pre-Training (SL128) PopART Packed DGX A100 (A100-SXM4-80GB) | TensorFlow | Mixed Precision | <u>https://developer.nvidia.com/deep-learning-performance-training-inference</u>

MLPERF RESNET: BOW POD PERFORMANCE

BOW POD256 DELIVERS SIGNIFICANTLY BETTER PERFORMANCE THAN 8X DGX AIOO



ResNet Performance vs DGX A100 server

GRAPHCORE

MLPerf v2.0 Training Results | MLPerf ID: 2.0-2047, 2.0-2050, 2.0-2052, 2.0-2054, 2.0-2090, 2.0-2094, 2.0-2073 The MLPerf name and logo are trademarks. See www.mlperf.org for more information,

TRAINING BENCHMARKS



BERT-Large Phase 1 Pre-Training Throughput (SL128) Bow Pod16 | Bow Pod256 | PopART with Packing | SDK 2.6 | <u>https://www.graphcore.ai/performance-results</u> DGX A100 320GB (A100-SXM4-40GB) | Mixed Precision | published BERT results DGX A100 TensorFlow - <u>https://ngc.nvidia.com/catalog/resources/nvidia.bert for tensorflow/performance</u>



Cluster-Graph Convolutional Networks Training | Multiple Datasets | TensorFlow 2 | FP16 Bow-2000 | SDK 2.6 | <u>https://www.graphcore.ai/berformance-results</u> | (A100-SXM4-40GB) measured results | Throughput values are measured using the best combination of data structure and matmul implementation for each device and each dataset

Вс



MPNN-GIN Graph Isomorphism Network Training | mol-hiv dataset (host-generated) | TensorFlow 2 41K graphs, 25.5 nodes per graph, 9-dimensional input node features | 1.7M model parameters Bow-2000 | FP 16.16 | SDK 2.6 | <u>https://www.graphcore.ai/performance-results</u> | 1x A100 (A100-SXM4-40GB) | FP16 | Measured results

EfficienNet-B4 Time-To-Train

5X SHORTER TTT, ~IOX PERF/\$



EfficientNet-B4 Training Time To Train Performance Bow Pod Platforms | SDK2.6 Results | G16-EfficientNet-B4 Training DGX A100 (A100-SXM4-80GB) | TensorFlow | Mixed Precision | https://developer.nvidia.com/deep-learning-performance-training-inference

GRAPHCORE BOW PODI6 NVIDIA DGX-AIOO 640GB SERVER \$149,995 MSRP \$299,000 MSRP

INFERENCE BENCHMARKS



ResNet-50 v1.5 Inference | Highest throughput comparison 1x Bow-M2000 | FP 16.16 | SDK2.6 | Synthetic Data (host-generated) | <u>https://dwww.graphcore.ai/performance-results</u> 1x A100 (A100-SXM4-80GB) results using TensorRT 8.0 | INT8 | Synthetic Data (host-generated) 1x A100 (A100-SXM4-40GB) results using TensorFlow | Mixed Precision + XLA | Synthetic Data (host-generated) 1x A100 (A100-SXM4-40GB) results using TensorFlow | Mixed Precision + XLA | Synthetic Data (host-generated) Results published by NVIDIA | A100 TensorRT: <u>(https://developer.nvidia.com/deep-learning-performance-training-inference</u>) on 22 Mar 2022 A100 TensorFlow: <u>https://github.com/NVIDIA/beepLearningZxmaples/tree/master/TensorFlow/Classification/ConvNets/resnet50v1.5#inference-performance-nvidia-dax-a100-1x-a100-40gb</u>



YOLO v4 Inference| Highest throughput comparison | Image Size 896 1x Bow-2000 | FP 16.16 | SDK 2.6 | PyTorch | Synthetic Data (host-generated) | <u>https://www.graphcore.ai/performance-results</u>, 1x A100 (A100-SXM4-40GB) results using PyTorch | FP16 | Synthetic Data(host-generated)



Unet (medical) Inference| Highest throughput comparison | 572x572x1 1x Bow-M2000 | FP 16.16 | SDK 2.6 | Synthetic Data (host-generated) | <u>https://www.graphcore.ai/performance-results</u> 1x A100 (A100-SXM4-80GB) results using TensorFlow2 | FP16 | Synthetic Data (host-generated)





Markov Chain Monte Carlo – Probabilistic model with TensorFlow Probability, representative of workload used by Carmot Capital Neural network with 3 fully-connected layers (num units in 1st layer=40, #dimensions in training set =22, #leapfrog steps=1000 Results multiplied by number of calcs in sliding window=200

1x Bow-2000 using TensorFlow | FP 32.32 | SDK 2.6 | 1600 samples | https://www.graphcore.ai/performance-results 1x A100 (A100-SXM4-40GB) | FP 32.32 | 1600 samples

CUSTOMER CASE STUDIES







HIGH ENERGY PHYSICS – PARTICLE DETECTOR



Model: <u>CookieNetAE</u> is a deep neural network designed to estimates the energy-angle dependent probability density function of electrons' energy.

Dataset: The training dataset 900,000 pairs of 128x128 images, a grainy input image and a smooth target image. The inference dataset consisted of 50,000 unique 128x128 images.

Task: Inference - Fast de-noising of detector images. Training – Quick "at-detector" re-training of model.

Goals:

- Improve inference throughput at detector
- Compress time-to-train

Results:

- Inference 7.7x throughput vs GPU
- Training 2.4x TTT improvement vs GPU



Machine Learning via direct attached accelerator for streaming data processing at high shot rate x-ray free electron lasers paper in review with <u>Frontiers in Physics</u>

GNN: CUSTOMER SUCCESS IOx FASTER FOR TGN GRAPH ML MODEL



Michael Bronstein Head of Graph ML Research @Twitter, DeepMind Professor of AI, University of Oxford

IPU Advantage for Dynamic Temporal Graph Networks



Using dynamic graph ML model, TGN, on IPU to achieve IOx faster time to result. Many real-world graphs are dynamic and evolve over time, including social networks, financial transactions, and recommender systems. In many cases, it is the dynamic behaviour of such systems that conveys important insights.

GNN: CUSTOMER SUCCESS 36X FASTER FOR SCHNET GRAPH ML MODEL (8 IPU VS 4 GPU)



IPU Advantage for molecular dynamics SchNet GNN



US Department of Energy National Lab, PNNL, training SchNet Graph Neural Network₁ with a 500k water clusters dataset² to predict the potential energy per cluster & seeing **36x faster time to result with IPU**

¹ JK. T. Schütt1, H. E. Sauceda, P.-J. Kindermans, A. Tkatchenko, and K.-R. Müller. "SchNet – A deep learning architecture for molecules and materials" *J. Chem. Phys.* **148**, 241722 (2018).

² Jenna A. Bilbrey, Joseph P. Heindel, Malachi Schram, Pradipta Bandyopadhyay, Sotiris S. Xantheas, and Sutanay Choudhury. "A look inside the black box: Using graph-theoretical descriptors to interpret a Continuous-Filter Convolutional Neural Network (CF-CNN) trained on the global and local minimum energy structures of neutral water clusters" *J. Chem. Phys.* **153**, 024302 (2020).

BOW REAL WORLD PERFORMANCE GAINS

BOW POD VS IPU-POD CLASSIC: UP TO 40% IMPROVEMENT



30w Pod64 vs IPU-POD64 | G16-EfficientNet-B4 Training/TensorFlow/ | Vision Transfomer (ViT) Training/PyTorch | GPT2-Large Training/PyTorch | Conformer SpeechToText Automatic Speech Recognition/PyTorch

Bow Pod16 vs (PU-PoD16 | ResNet-50 v1.5 Training/TensorFlow | ResNetX-101 Training/TensorFlow | Mini DALL-E Training/TensorFlow2 | BERT-Large Ph1 Pre-Training (SL128)/PopART | BERT-Base Ph1 Pre-Trai

Comparing Preliminary Results (Pre-SDK2



Accelerating HPC with Al using IPU

European Centre for Medium-range Weather Forecasting

Graphcore IPU trained an ECMWF weather forecasting model **5x faster than a leading GPU** (and potentially up to 50x faster than CPUs)



Mar 09, 2022 💧 Al, HPC

AI FOR SIMULATION: HOW GRAPHCORE IS HELPING TRANSFORM TRADITIONAL HPC



Written By: Alex Titterton

Atos

Cedric Bourrasset Head of HPC AI at Atos



"Graphcore plays a central role in Atos' Think AI solution, helping customers take advantage of the many benefits that AI is bringing to HPC – whether that's delivering **faster and more accurate simulations**, **improving cost efficiency**, or opening up **new areas of research and commercial applications**."

f 🔰

"Without any changes to the model definition or its hyperparameters, a single IPU processor was able to train this MLP model **5x faster than an A100 GPU** and some 50x faster than either the MLP model or the traditional parametrisation scheme performed on a CPU."

"Graphcore's IPU, designed from the ground up for machine intelligence, is the ideal platform on which to build, explore and grow the next generation of machine learning-driven solvers, emulators and simulations."

More information can be found in the full technical blog post: <u>https://www.graphcore.ai/posts/ai-for-simulation-how-graphcore-is-helping-transform-traditional-hpc</u>

https://agupubs.onlinelibrary.wiley.com/doi/full/10.1029/2021MS002477

ARGONNE APPLICATIONS





CANDLE UNO

CANcer distributed Learning Environment (CANDLE)

Model: CANDLE Uno for precision medicine related to cancer care

Dataset: dataset with 20 million training samples

Task: train model to predict tumor-drug response

Goals:

- Reduce time-to-train
- Demonstrate efficient scaling

Results: Increased training throughput by 6x @ higher precision





BraggNN Model Details



Peak Location with Sub-Pixel Accuracy



Framework: PyTorch

Input: An 11x11 (single channel) image generated on the fly by cropping a larger detector image

Output: Normalized location of the peak in (y, z) coordinates

Loss: Mean square error between predicted location and target location

GPU Performance Reference: ~800 seconds for 500 epochs with minibatch size: 512

BC GRAPHCORE **Goal**: Reduce the time to train to < 2 minutes & show excellent scaling over multiple IPUs

Zhengchun L. et al. "BraggNN: Fast X-ray Bragg Peak Analysis Using Deep Learning", 2021

DATAFLOW MODIFICATION FOR FASTER TRAINING



DATAFLOW MODIFICATION FOR FASTER TRAINING



DATAFLOW MODIFICATION FOR FASTER TRAINING



38

BRAGGNN SUMMARY



- Out of the box, training on <u>MK2</u> IPU is **3-4x** faster than V100
- Flexibility of Graphcore Architecture allowed migration of image scaling and data augmentation to IPU
- Final version of BraggNN is >5x faster on IPU compared to A100 and >11x faster than V100



"GOOD" COMPUTER





ANNOUNCING: THE 'GOOD' COMPUTER



Over 10 **Exa-Flops** of AI floating point compute from 8,192 roadmap IPUs

3D Wafer-on-Wafer logic stack

Up to **4** PB of memory with bandwidth of over **10** PB/s

Enabling AI models to be developed with 500 Tn parameters

Fully supported by Poplar[®] SDK

GRAPHCORE

THANK YOU

Richard Bohl richardb@graphcore.ai



42

