NEW APPROACHES FOR ENABLING MACHINE LEARNING INNOVATION

Alexander Tsyplikhin, Senior Al Engineer

GRAFHCORE





MACHINE INTELLIGENCE REPRESENTS A COMPLETELY NEW COMPUTE WORKLOAD

conv2 - 1x1

conv3 – 1x1 12 in, 128 out]

conv1 - 7x7

[4 in, 64 out]

conv4 – 1x1 [256 in, 1024 out] Massive parallelism

conv2 - 1x1

conv2 - 3x3

Sparsity in data structures Low precision compute

Model parameter re-use Static graph structure



Fully Connected [2048 in, 1000 out]

256 in. 128 out

IN FINANCIAL MODELLING

MODELS ARE OFTEN OVERSIMPLIFIED TO ACHIEVE COMPETITIVE LATENCY OR FIGHT NOISE



MODEL SIZE GROWING MASSIVELY



(^DC



LET'S CONSIDER 3 EXAMPLES



ALPHA ESTIMATION



Traditional method: linear regression

• Limited efficiency on real data

Markov Chain Monte Carlo (MCMC): **predictions are distributions**, not just point values

This helps to:

- Avoid overfitting
- Deal with noisy data
- Handle non-linearities

6

HAMILTONIAN MONTE CARLO

A hockey puck sliding over a surface without friction, being stopped at some point in time and then kicked again in a random direction.



THE PROBLEM WITH MCMC

MCMC has long been considered too computationally expensive.

Runs for hours or days. Not practical using traditional technology.



NOT A STAC BENCHMARK

NOTES:

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, calcs in sliding window=200) GPU (300W TDP) - 800 samples

LSTM FOR TIME SERIES ANALYSIS



Long short-term memory (LSTM) models: efficient for sequences of observations

Used for:

- Financial time series analysis
- Feature extraction
- Alpha estimation
- Fraud detection



RECURRENT NEURAL NETWORK





LSTM



THE PROBLEM WITH LSTM

Inference latencies are often too high, or throughput is too low



NOTES:

2 LSTM layers, each with 256 units, 200 time steps, 16 input dimensions, real data, mixed precision GPU: using TensorFlow with optimizations @ 300W TDP (Batch Sizes upto 1024)

8c

EFFICIENTNET FOR IMAGE ANALYSIS



https://paperswithcode.com/sota/image-classification-on-imagenet

Some hedge funds use alternative data, such as images.

SOTA computer vision models use group or depthwise convolutions:

- Out of the top 50 published results on ImageNet, 42 use group or depthwise convolutions
- ResNeXt and EfficientNet are the leading backbone architectures

8c

GROUP CONVOLUTION

Regular convolution

Group convolution





THE PROBLEM WITH EFFICIENTNET

Both training and inference for group convs can be slow on traditional technology

NOT A STAC BENCHMARK

NOTES:

EfficientNet-BO | Synthetic Data | throughput comparison using highest throughput | latency comparison using lowest latency GPU: 1x GPU (FP32) using TensorFlow & published Google reference. Batch Size 1-32 @ 300W TDP GPU results using public Google repo (https://github.com/tensorflow/tpu/tree/master/models/official/efficientnet/)

EFFICIENTNET-BO: TRAINING

7x higher throughput

NOTES:

EfficientNet-BO | Real Data (ImageNet)

GPU: 2x GPU (FP32) using TensorFlow @ 600W TDP (note: GPU throughput drops with modified version of EN-BO) GPU results using public Google repo (https://github.com/tensorflow/tpu/tree/master/models/official/efficientnet/)

WHAT IS HOLDING THESE MODELING TECHNIQUES BACK?

LEGACY PROCESSOR ARCHITECTURES HAVE BEEN REPURPOSED FOR ML

CPU Apps and Web/ Scalar

GPU Graphics and HPC/ Vector

MACHINE LEARNING REQUIRES **2 KINDS OF COMPUTE INNOVATION**

	CPU	GPU	IPU
Parallelism	Suitable for scalar processes	SIMD/SIMT architecture. Suitable for large blocks of dense contiguous data	Massively parallel MIMD. High performance/efficiency as ML trends to sparsity & small kernels
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
	1 x	5x – 32x	320x
		Processor Memory	20

GRAPHCORE HAVE IMPLEMENTED THOSE INNOVATIONS IN THE IPU

CPU Apps and Web/ Scalar

GPU Graphics and HPC/ Vector IPU Artificial Intelligence/ Graph

LET'S RETURN TO OUR EXAMPLES

MCMC : NOW FEASIBLE

Alpha estimation with TensorFlow Probability

15x faster time to train

NOT A STAC BENCHMARK

NOTES:

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, calcs in sliding window=200) IPU: C2 card (300W TDP) results (SDK 1.2 GA) – 800 samples GPU (300W TDP) - 800 samples

LSTM INFERENCE FOR TIME SERIES ANALYSIS

NOT A STAC BENCHMARK

NOTES:

2 LSTM layers, each with 256 units, 200 time steps, 16 input dimensions, real data, mixed precision IPU: Graphcore C2 using TensorFlow and PopNN (SDK 1.2 GA) @ 300W TDP (Batch Sizes upto 1024) GPU: using TensorFlow with optimizations @ 300W TDP (Batch Sizes upto 1024)

EFFICIENTNET-BO INFERENCE

>15x higher throughput | >14x lower latency

Throughput (images/sec)

NOT A STAC BENCHMARK

ပြီင

NOTES:

EfficientNet-B0 | Synthetic Data | throughput comparison using highest throughput | latency comparison using lowest latency IPU: Graphcore C2 (SDK 1.2 GA) mixed-precision using TensorFlow Batch Size 1-12 @ 300W TDP GPU: 1x GPU (FP32) using TensorFlow & published Google reference. Batch Size 1-32 @ 300W TDP GPU results using public Google repo (https://github.com/tensorflow/tpu/tree/master/models/official/efficientnet/)

EFFICIENTNET-BO: TRAINING

7x higher throughput

NOT A STAC BENCHMARK

(gc

NOTES:

EfficientNet-BO | Real Data (ImageNet)

IPU: Graphcore 2x C2 (SDK 1.2 GA) mixed-precision using TensorFlow @ 600W TDP (EN-BO modified version uses Group Dim 16) GPU: 2x GPU (FP32) using TensorFlow @ 600W TDP (note: GPU throughput drops with modified version of EN-BO) GPU results using public Google repo (https://github.com/tensorflow/tpu/tree/master/models/official/efficientnet/)

IPU's novel architecture is already enabling us to explore new techniques that have been inefficient or simply not possible before."

"

We were able to train one of our proprietary probabilistic models In 4.5 minutes instead of 2 hours. That's 26x faster. "

GRAFHCORE IPU-M2000TM INTRODUCTION

COLOSSUS MK2

the world's 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

COLOSSUS Mk2 PERFORMANCE

IPU-MACHINE M2000

1 PetaFlop IPU compute 2.8Tbps IPU-Fabric[™]

M2000 PRODUCT CONFIGURATIONS

STANDALONE

IPU-POD₆₄

IPU-POD

COMMUNICATIONS

1010100101 0100101011 0101010101 0010111010 1011001010

DATA

 \bullet

COMPUTE

IT IS NOW POSSIBLE FOR FINANCIAL INNOVATORS TO CREATE THE NEXT GENERATION OF MACHINE INTELLIGENCE

Inference

- Make decisions with lower latency
- Use models of greater complexity

Training

- Iterate faster when experimenting with model architectures
- Use probabilistic models to avoid overfitting and account for noise

THANK YOU

Alexander Tsyplikhin alext@graphcore.ai

