

# GRAPHCORE WORKSHOP FOR TAMU

**October 3, 2023**

Alexander Tsyplikhin



# AGENDA

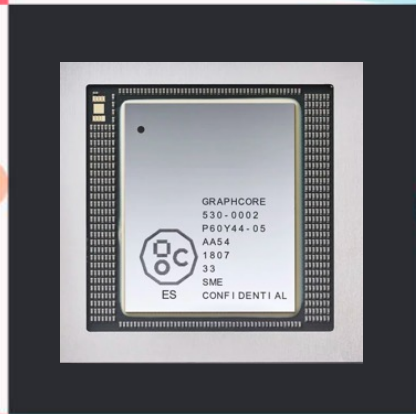
---

- Introduction to Graphcore, IPU, and Poplar
  - Hands-on: access the POD, enable the SDK, run an example
- TensorFlow2/Keras
  - Hands-on: Port a Keras script, leverage loop on device, replicate and run data-parallel, pipeline
- PyTorch
  - Hands-on: PopTorch example, DataLoader, options to optimize performance



# GRAPHCORE ENABLING MACHINE INTELLIGENCE

- Founded in 2016
- Technology: Intelligence Processor Unit (IPU)
- Team: ~500
- Offices: UK, US, China, Poland
- Raised >\$710M



SEQUOIA



ATOMICO



SOFINA

Microsoft



BMW i VENTURES

DELL



BOSCH

SAMSUNG

Merian  
GLOBAL INVESTORS

Amadeus  
Capital Partners

pitango  
VENTURE CAPITAL

draper esprit

Foundation  
CAPITAL



An aerial, top-down view of numerous Graphcore IPU (Intelligent Processing Units) accelerator cards. The cards are arranged in a grid-like pattern on a light yellow background. Each card features a distinctive heat shield design composed of several colored rectangular sections: dark blue, light blue, red, and beige. The word "GRAPHCORE" is embossed on the top surface of each card. The perspective is slightly angled, showing the front and top of the cards.

**GRAPHCORE IPU LETS INNOVATORS CREATE THE NEXT  
BREAKTHROUGHS IN MACHINE INTELLIGENCE**

# IPU – ARCHITECTURED FOR AI

Massive parallelism with ultrafast memory access

## Parallelism

Processors   
Memory 

## Memory Access

### CPU

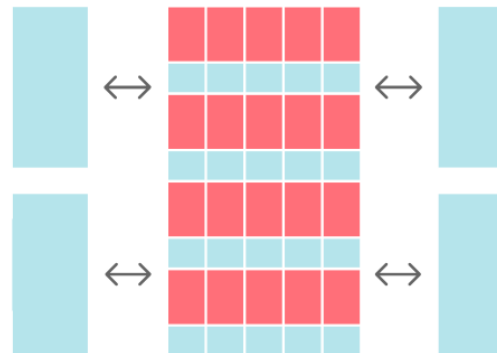
Designed for scalar processes



Off-chip memory

### GPU

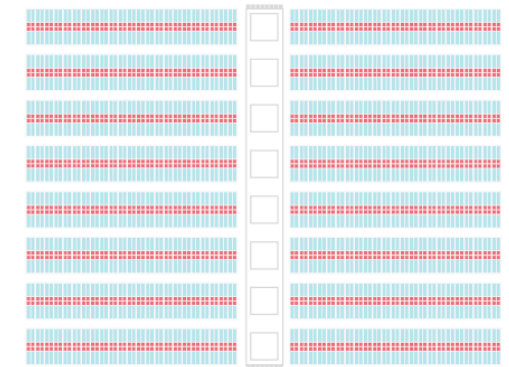
SIMD/SIMT architecture. Designed for large blocks of dense contiguous data



Model and data spread across off-chip and small on-chip cache, and shared memory

### IPU

Massively parallel MIMD. Designed for fine-grained, high-performance computing



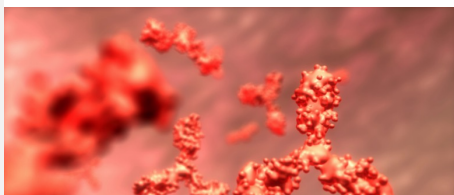
Model and data tightly coupled, and large locally distributed SRAM

# PROVEN IPU ADVANTAGE

## SELECT CASE STUDIES ACROSS MANY INDUSTRIES & FIELDS



HEALTHCARE



CASE STUDY : NLP >



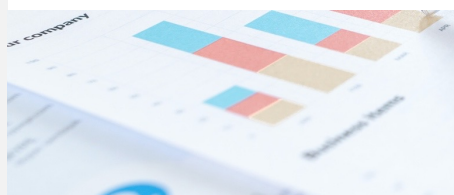
FINANCE – OPTION PRICING



CASE STUDY : SIM >



AI SaaS – TEXT ANALYTICS



CASE STUDY : NLP >



RESEARCH / BIG LABS



CASE STUDY >



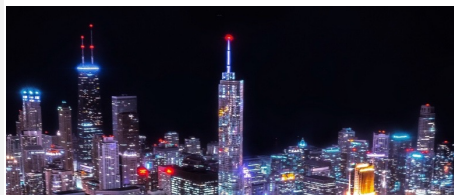
COMPUTATIONAL CHEMISTRY



CASE STUDY : GNN >

SENSORO

SMART CITY



CASE STUDY : CV >



TRACTABLE

FINANCE - INSURANCE



CASE STUDY : CV >



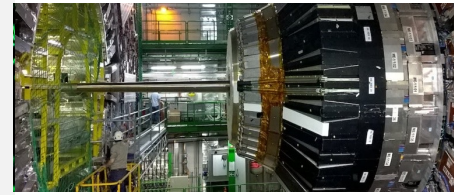
WEATHER FORECASTING



CASE STUDY : SIM >



HIGH ENERGY PHYSICS



CASE STUDY >



DYNAMIC GRAPHS



CASE STUDY : GNN >



# IPU COMPUTATIONAL ADVANTAGES

Heterogeneous gather/scatter operations. E.g. GNNs

Group and depthwise convolutions. E.g. ResNeXt, EfficientNet

Vector operations with low arithmetic intensity. E.g. Sparse matmuls

Dense as well as Sparse Matrix Multiplication. E.g. Transformers

Hardware accelerated Random Number Generation. E.g. Random Projections

Hard to vectorize workloads. E.g. DFT in Computational Chemistry

## References:

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

<https://www.graphcore.ai/posts/how-we-made-efficientnet-more-efficient>

<https://www.graphcore.ai/posts/delving-deep-into-modern-computer-vision-models>

<https://www.graphcore.ai/posts/training-neural-networks-in-low-dimensional-random-bases>

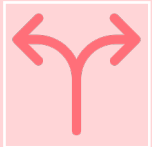
<https://www.graphcore.ai/posts/man-group-unlocks-massively-parallel-option-pricing-with-graphcore-ipu>



# WORKLOADS THAT CAN'T EASILY BE VECTORIZED



Workloads with while loops that continue until convergence is achieved e.g. ray tracing



Workloads where different compute paths are required depending on the inputs e.g. DFT or CRR model



Tree-based models with unbalanced trees



- > .github
- > .gradient
- > gdb
- > images
- > notebooks
- > pyscf\_ipu
- > qm1b
- > schnet\_9m
- .gitignore
- LICENSE
- README.md
- density\_functional\_theory.py
- generate.sh
- pyproject.toml
- requirements.txt
- requirements\_ipu.txt
- setup.sh
- test\_requirements.txt

# PySCF on IPU

[Installation guide](#) | [Example DFT Computations](#) | [Generating data](#) | [Training SchNet](#) | [QM1B dataset](#)

 Run on Gradient  pytest notebooks passing

Port of PySCF to Graphcore IPU.

## Limitations

- Restricted Kohn Sham DFT (based on [RKS](#), [KohnShamDFT](#) and [hf.RHF](#)).
- Number of atomic orbitals less than 70 `mol.nao_nr() <= 70`.
- Larger numerical errors due to `np.float32` instead of `np.float64`.
- Limited support for `jax.grad(.)`

## Installation

PySCF on IPU requires Python 3.8, [JAX IPU experimental](#), [TessellateIPU library](#) and [Graphcore Poplar SDK 3.2](#).

To run this package on a standard CPU machine (laptop or server), install the base Python requirements:

```
pip install -r requirements.txt
```

On IPU machines, please additionally use the IPU requirements file:

```
pip install -U pip
pip install -r requirements_ipu.txt
```

This will configure Graphcore research experimental JAX support in your python environment.

We recommend upgrading `pip` to the latest stable release when using the IPU requirements. This may be an optional step depending on the overall configuration of your python environment.

And finally, make our sub-packages available:

```
pip install -e .
```



## Files



- gdb
- images
- notebooks
  - DFT-dataset-generation.ipynb
  - ERI-visualisation-JK.ipynb
  - nanoDFT-demo.ipynb
- pyscf\_ipu
- qm1b
- schnet\_9m
- density\_functional\_theory.py
- generate.sh
- LICENSE
- pyproject.toml
- README.md
- requirements\_ipu.txt
- requirements.txt
- setup.sh
- test\_requirements.txt

## Kernel sessions

- nanoDFT-demo.ipynb

STOP MACHINE

Running

RESTART KERNEL

SAVE

RUN ALL

```

34
35 def plot_orbital(orbital, mol):
36     xyzfmt = f"{len(mol.atom)}\n\n" + mol.tostring()
37     v = py3Dmol.view(data=xyzfmt, style={"stick": {"radius": 0.06}, "sphere": {"radius": 0.2}})
38     v.addVolumetricData(cube_data(axes, orbital), "cube", build_transferfn(orbital))
39     return v

```

Try changing the `mo_index` variable to select the different molecular orbitals benzene.

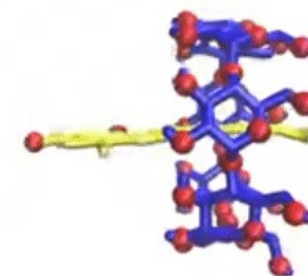
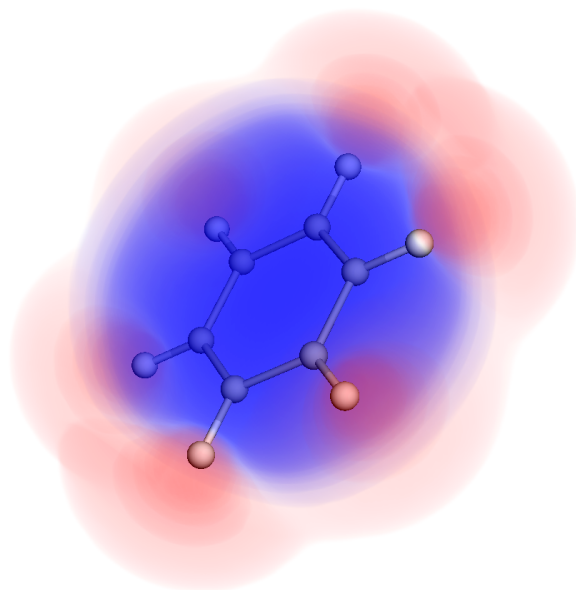
Run

10

```

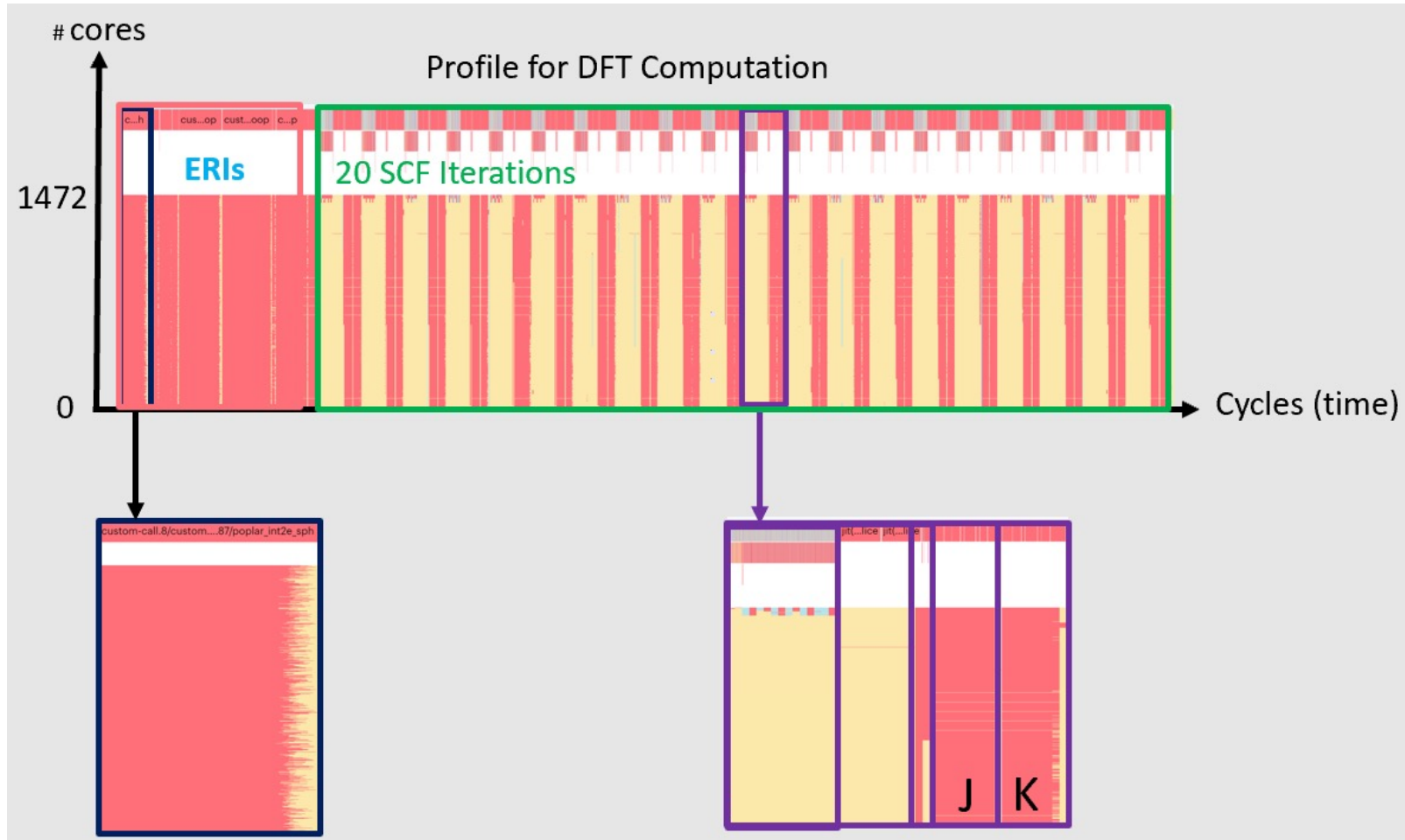
1 mo_index = 5
2
3 orbital = molecular_orbitals[:, mo_index]
4 mol_view = plot_orbital(orbital, mol)
5 mol_view.spin()

```



SHOW MORE

10



MIMD for  
 $E \in \mathbb{R}^{N^2 \times N^2}$

65 TB/s for  
 $J = E \cdot v \in \mathbb{R}^{N^2}$   
 $K = E^T \cdot v \in \mathbb{R}^{N^2}$

# BOW IPU

## IPU-Tiles™

1472 independent IPU-Tiles™ each with an IPU-Core™ and In-Processor-Memory™

## IPU-Core™

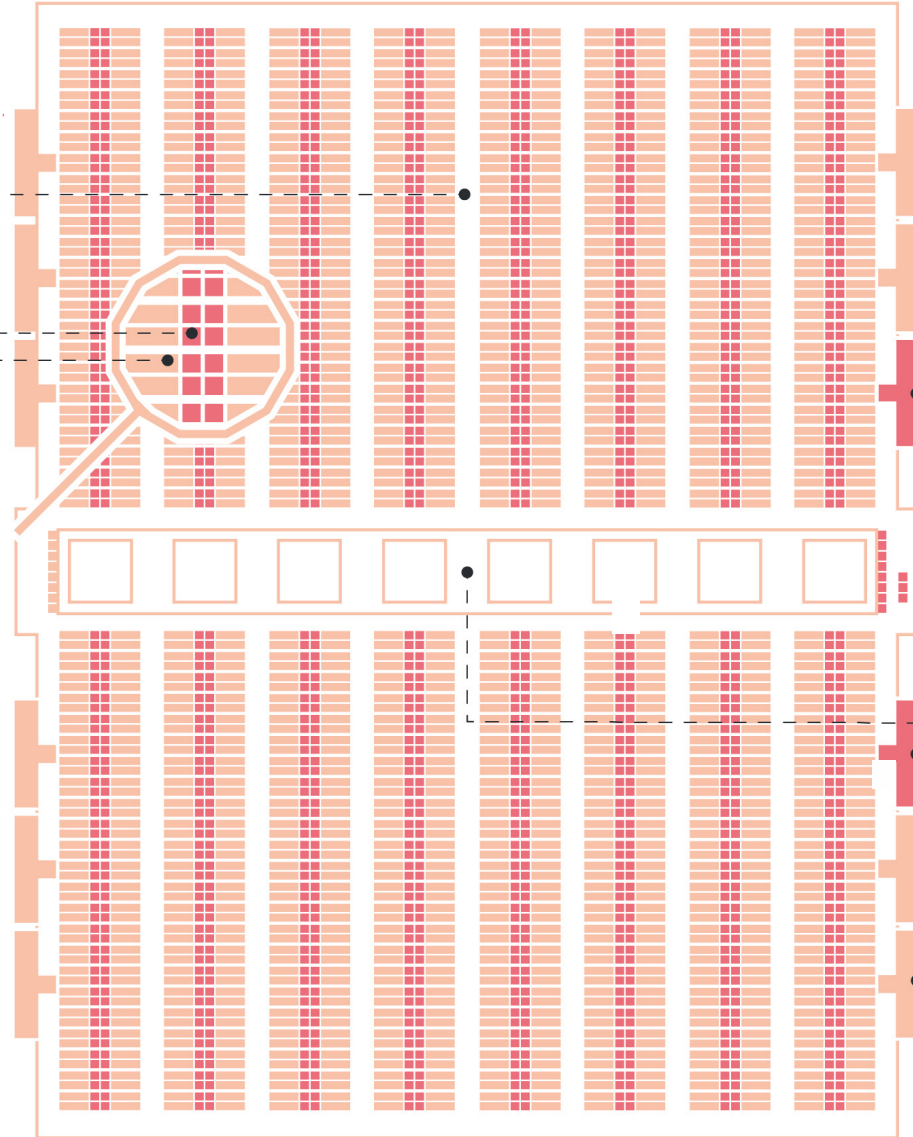
1472 independent IPU-Core™

8832 independent program threads executing in parallel

## In-Processor-Memory™

900MB In-Processor-Memory™ per IPU

65TB/s memory bandwidth per IPU



## IPU-Exchange™

11 TB/s all to all IPU-Exchange™  
Non-blocking, any communication pattern

## PCIe

PCI Gen4 x16  
64 GB/s bidirectional bandwidth to host

## IPU-Links™

10 x IPU-Links,  
320GB/s chip to chip bandwidth

# BOW IPU PROCESSOR

## Deep Trench Capacitor

Efficient power delivery  
Enables increase in operational performance

## Wafer-On-Wafer

Advanced silicon 3D stacking technology  
Closely coupled power delivery die  
Higher operating frequency and enhanced overall performance

## IPU-Tiles™

1472 independent IPU-Tiles™ each with an IPU-Core™ and In-Processor-Memory™

## IPU-Core™

1472 independent IPU-Core™  
8832 independent program threads executing in parallel

## In-Processor-Memory™

900MB In-Processor-Memory™ per IPU  
65.4TB/s memory bandwidth per IPU

## Solder Bumps

## IPU-Links™

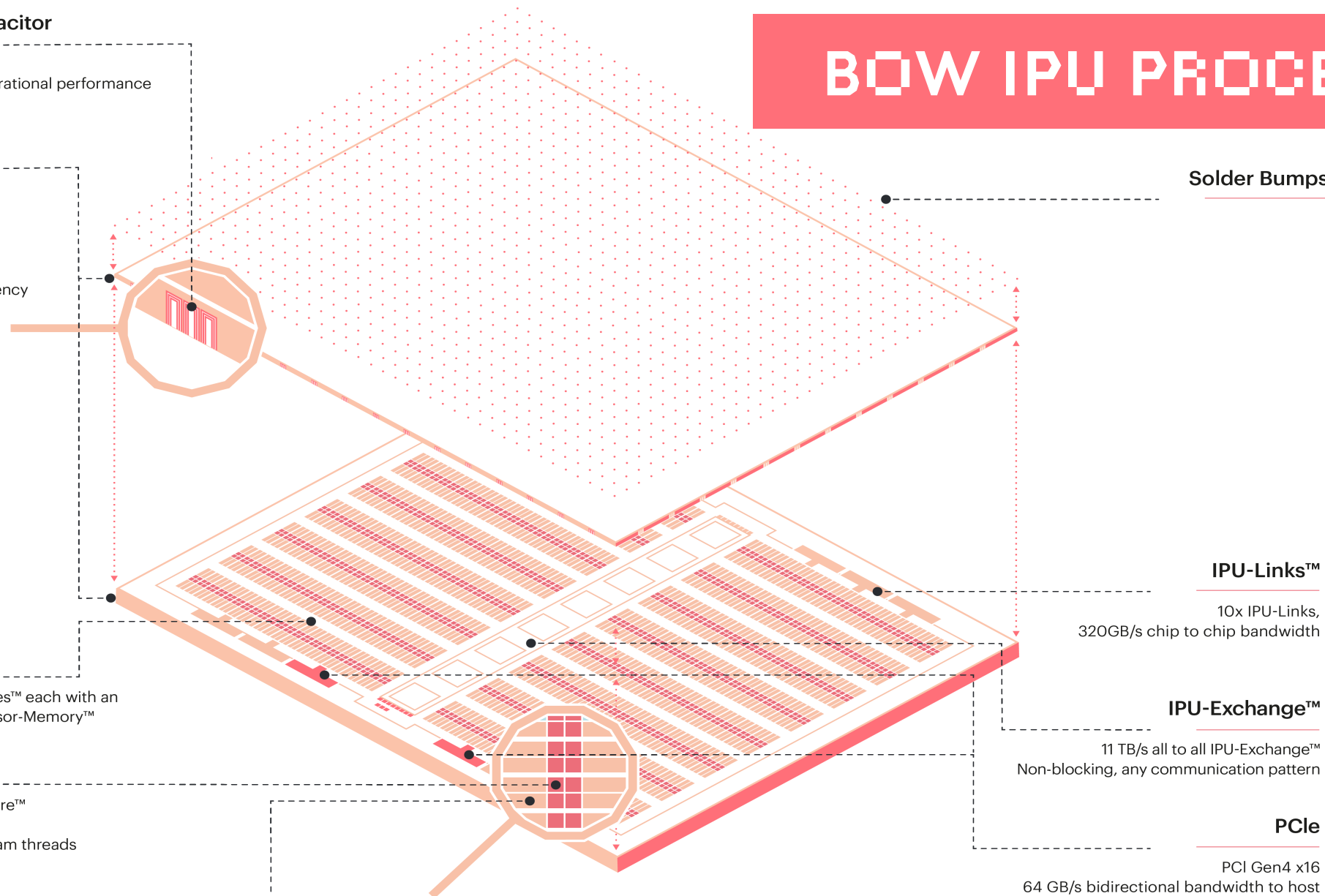
10x IPU-Links,  
320GB/s chip to chip bandwidth

## IPU-Exchange™

11 TB/s all to all IPU-Exchange™  
Non-blocking, any communication pattern

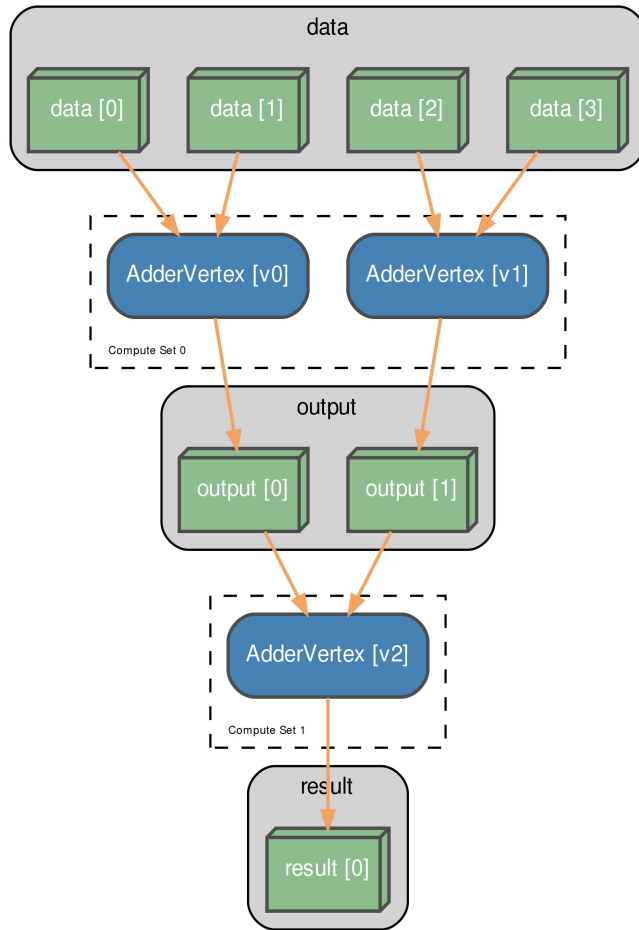
## PCIe

PCI Gen4 x16  
64 GB/s bidirectional bandwidth to host

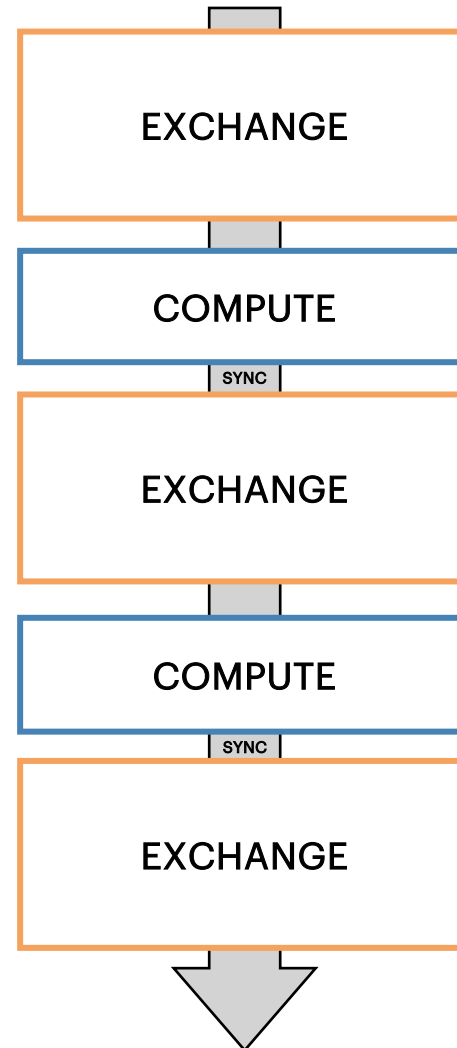


# EXECUTION MODEL

## COMPUTATIONAL GRAPH

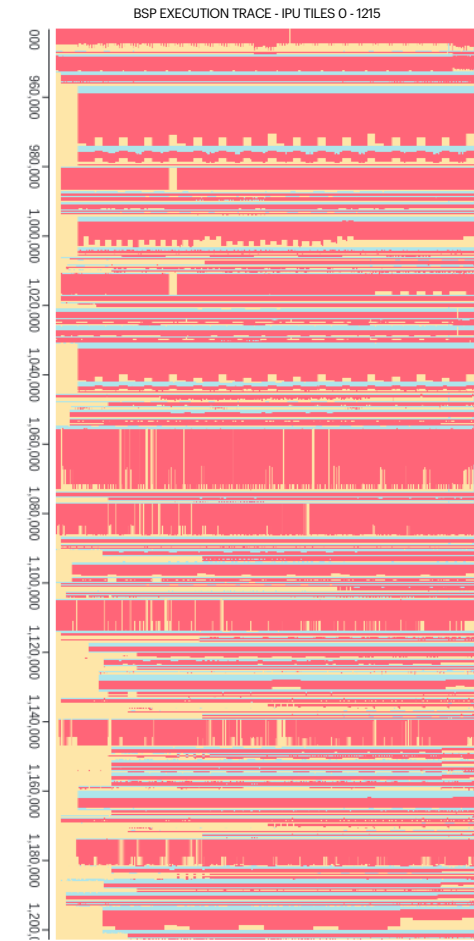


## BSP SCHEDULE



GRAPHCORE

## OPTIMIZED IPU EXECUTION



OUTPUT FROM POPVISION GRAPH ANALYSER

# BULK SYNCHRONOUS PARALLEL (BSP)

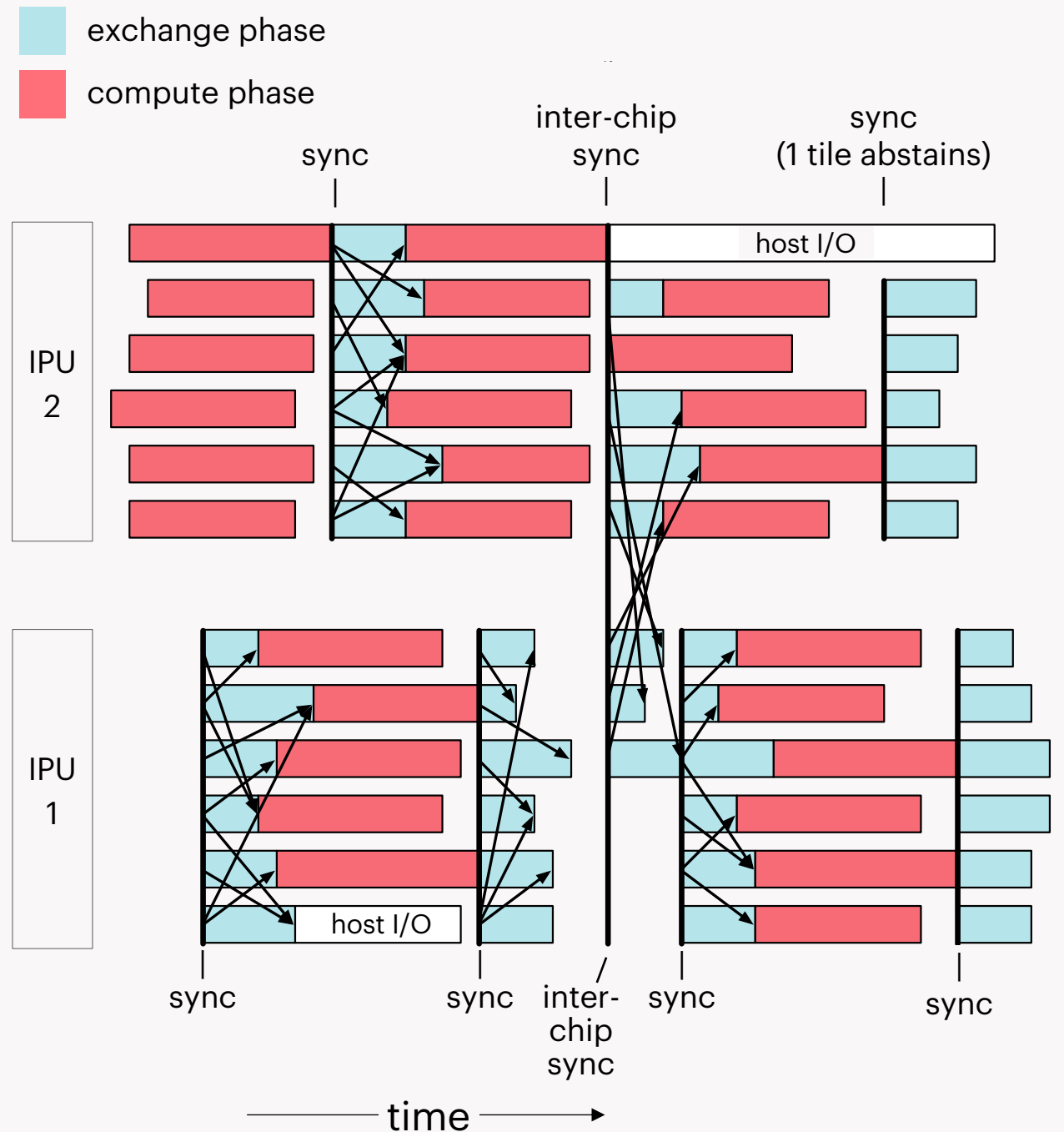
BSP software bridging model – massively parallel computing with no concurrency hazards

3 phases: compute, sync, exchange

Easy to program – no live-locks or dead-locks

Widely-used in parallel computing – Google, FB, ...

First use of BSP inside a parallel processor



# BOW-2000 IPU MACHINE

IPU blade form factor delivering 1.4 PetaFLOPS AI Compute

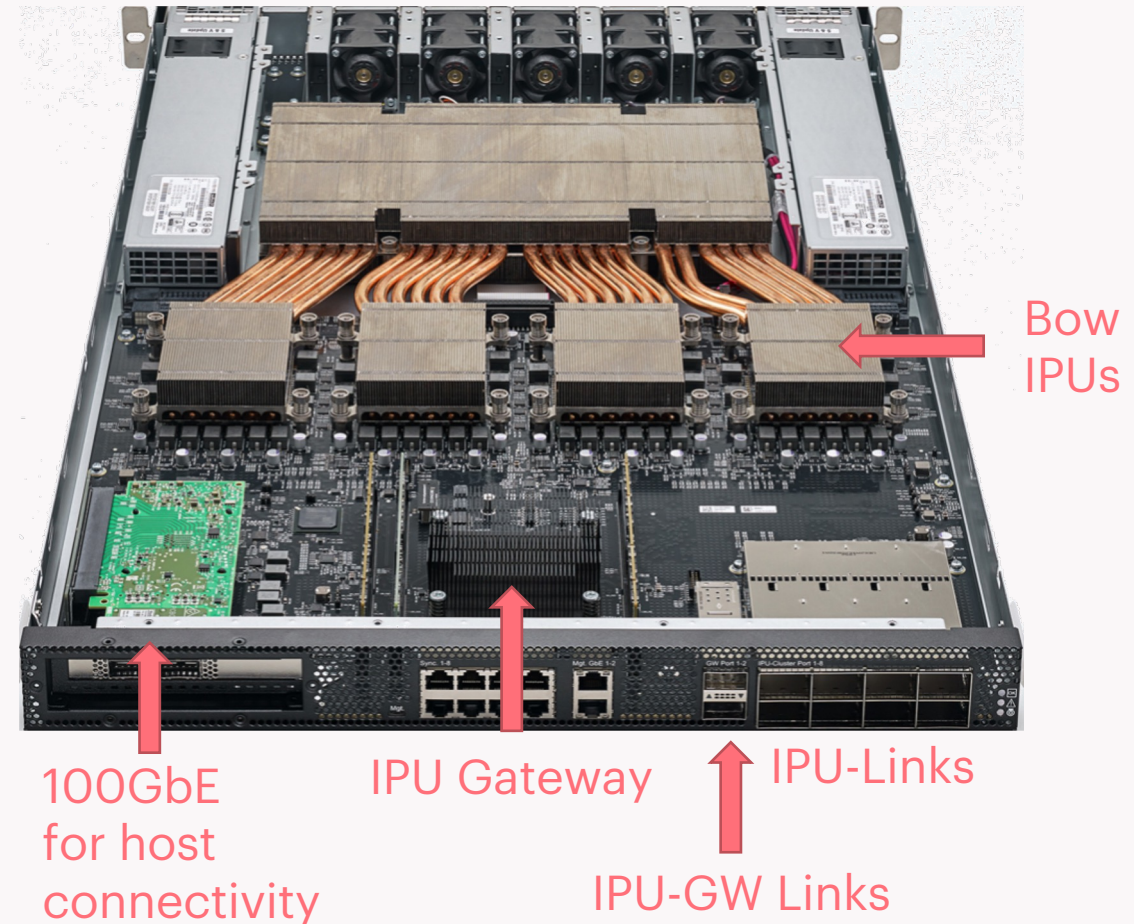
Disaggregated AI/ML accelerator platform

Excellent performance & TCO leveraging  
In-Processor memory & IPU-Exchange

IPU-Links scale to Bow Pod64

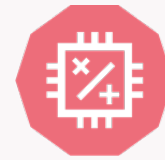
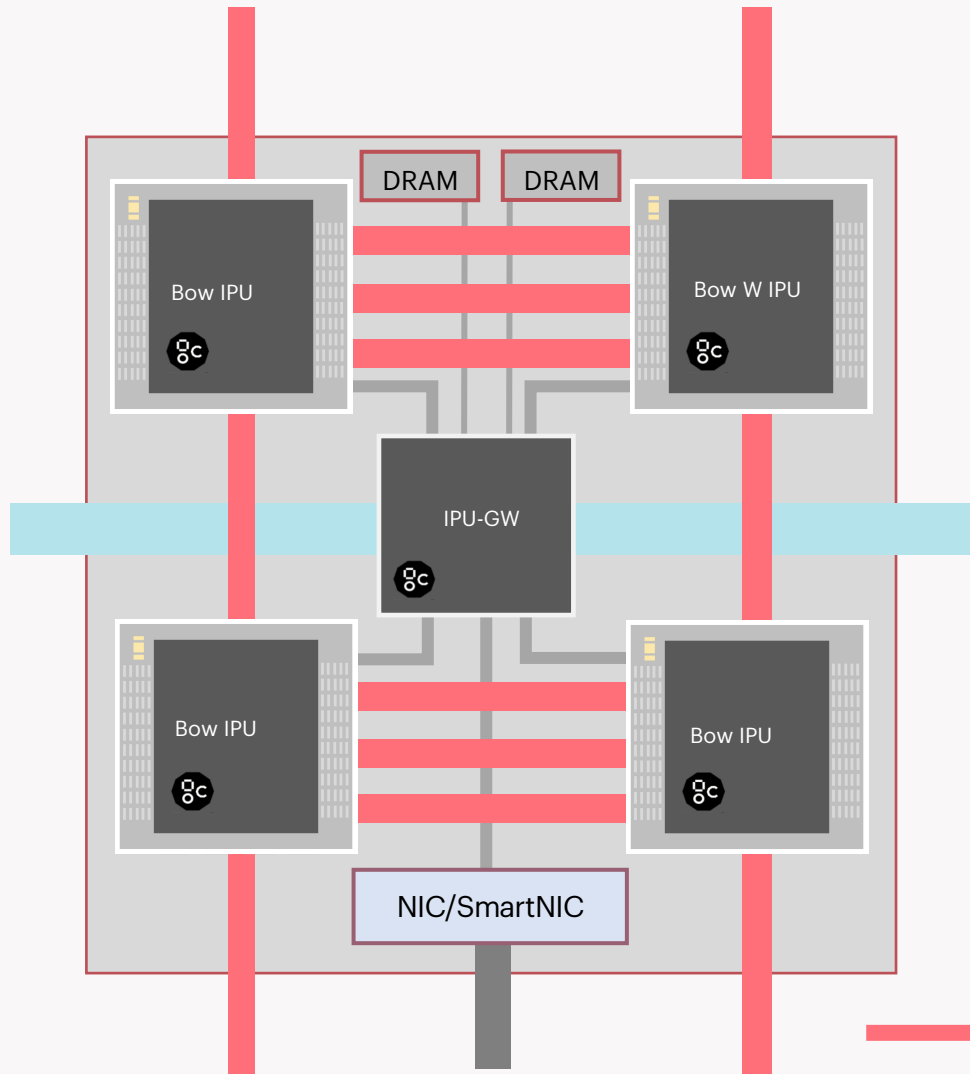
Expansion to Bow Pod256 and beyond  
with IPU-GW Links

BOW IPU-2000





# BOW-2000: THE BUILDING BLOCK OF LARGE PODS



COMPUTE

## 4x Bow IPUs

- 1.4 PFLOP<sub>16</sub> compute
- 5,888 processor cores
- > 35,000 independent parallel threads



DATA

## Exchange Memory





- 3.6GB In-Processor Memory @ 260 TB/s
- 128GB Streaming Memory DRAM (up to 256GB)



COMMUNICATIONS

## IPU-Fabric managed by IPU-GW

- Host-Link - 100GE to Poplar Server for standard data center networking
- IPU-Link - 2D Torus for intra-POD64 communication
- GW-Link - 2x 100Gbps Gateway-Links for rack-to-rack - flexible topology





-  x16 IPU-Link [64GB/s]
-  Host-Link Network I/F [100Gbps]
-  IPU-GW Link [100Gbps]
-  x8 PCIe G4 [32GB/s]



# BOW POD & DIRECT ATTACH

Simple design with server with 4 Bow-2000 and 16 Bow IPU



-  Host-Link 100GE network interface (QSFP, 1.0m)
-  1GbE Management (Cat5, 1.5m)
-  Sync-Link (Cat5, 0.15m)
-  IPU-Link (OSFP, 0.3m)

**HANDS-ON:**

**GET STARTED**

**RUN AN EXAMPLE**



**HANDOUT**

[bit.ly/tamu231003](https://bit.ly/tamu231003)





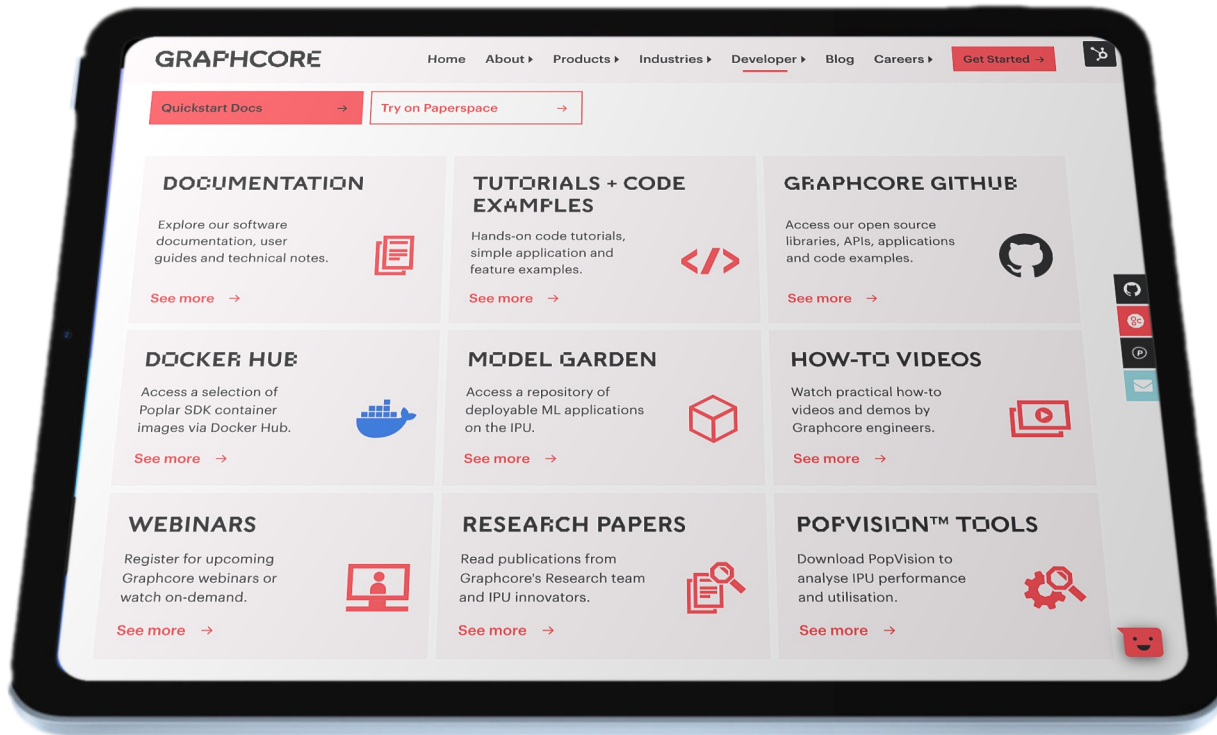
# IPU DEVELOPER ECOSYSTEM

GRAPHCORE

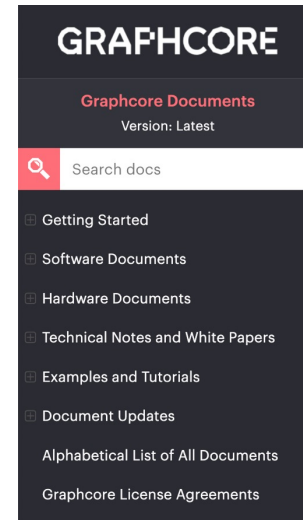


# GRAPHCORE SOFTWARE ECOSYSTEM

## WORLD CLASS DEVELOPER RESOURCES FOR IPU USERS



[WWW.GRAPHCORE.AI/DEVELOPER](http://WWW.GRAPHCORE.AI/DEVELOPER)



### GRAPHCORE DOCUMENTS

#### Getting Started

Background information and quick-start guides for Graphcloud and Pod systems

#### Software Documents

Documentation for the Poplar SDK and other software

#### Hardware Documents

Documentation for installing and using IPU-Machines and Pod systems

#### Technical Notes and White Papers

Technical notes and white papers on Graphcore technology

#### Document Updates

The latest news about new documents and examples

#### Examples and Tutorials

Tutorials and application examples for running on the IPU



### Getting started with PyTorch for the IPU

Running a basic model for training and inference

AI Customer Engineer, Chris Bogdiukiewicz introduces PyTorch for the IPU. With PopTorch™ - a simple Python wrapper for PyTorch programs, developers can easily run models, directly on Graphcore IPUs with a few lines of extra code.

[Get the Code](#) →

In this video, Chris provides a quick demo on running a basic model for both training and inference using a MNIST based example.

[Read the Guide](#) →



# OPEN SOURCE

[github.com/graphcore](https://github.com/graphcore)

- As part of our ethos to put power in the hands of AI developers, Graphcore open sourced in 2020
- PopLibs™, PopART, PyTorch & TensorFlow for IPU fully open source and available on GitHub
- Our code is public and open for code contributions from the wider ML developer community



A screenshot of the Graphcore GitHub organization page. The page shows the organization's profile with the Graphcore logo and name. Below the profile, there are navigation tabs for Repositories (6), Packages, People, and Projects. A prominent banner reads "Grow your team on GitHub" with a "Sign up" button. Below the banner is a search bar and filters for repository type and language. The main content area lists several repositories with their respective languages, licenses, and statistics. On the right side, there are sections for "Top languages" (C++ and Python) and "People" (indicating no public members).

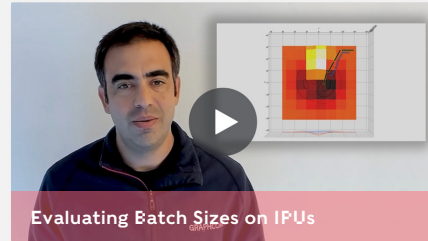
# VIDEO + GITHUB TUTORIALS

A comprehensive set of online developer training materials and educational content

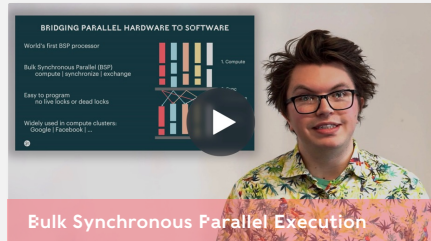


PyTorch

Getting started with PyTorch for the IPU



Evaluating Batch Sizes on IPU



BRIDGING PARALLEL HARDWARE TO SOFTWARE

World's first BSP processor

Bulk Synchronous Parallel (BSP) compute | synchronize | exchange

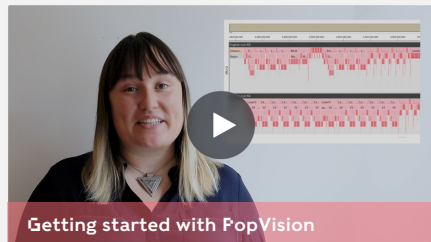
Easy to program, no low-level details

Widely used in compute clusters (Google | Facebook)

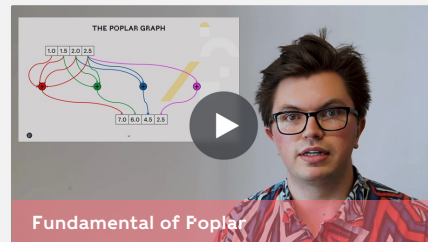
Bulk Synchronous Parallel Execution



Running PyTorch on the IPU: NLP

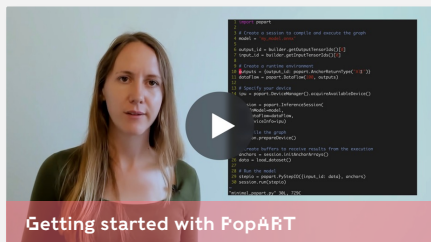


Getting started with PopVision

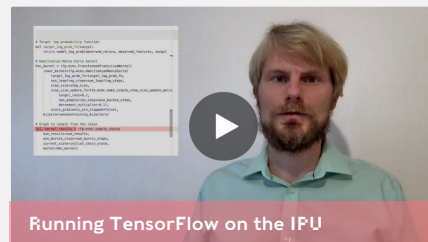


THE POPLAR GRAPH

Fundamental of Poplar



Getting started with PopART



Running TensorFlow on the IPU

## TUTORIALS

Learn how to create and run programs using Poplar and PopLibs with our hands-on programming tutorials.

Programs and Variables

Using PopLibs

Writing Vertex Code

Profiling Output

Basic Machine Learning Example

Matrix-Vector Multiplication

Matrix-Vector Multiplication Optimisation

Simple PyTorch for the IPU

NEW

### Tutorial 1: programs and variables

Copy the file `tut1_variables/start_here/tut1.cpp` to your working directory and open it in an editor. The file contains the outline of a C++ program including some Poplar library headers and a namespace.

#### Graphs, variables and programs

All Poplar programs require a `Graph` object to construct the computation graph. Graphs are always created for a specific target (where the target is a description of the hardware being targeted, such as an IPU). To obtain the target we need to choose a device.

The tutorials use a simulated target by default, so will run on any machine even if it has no Graphcore hardware attached. On systems with accelerator hardware, the header file `poplar/DeviceManager.hpp` contains API calls to enumerate and return `Device` objects for the attached hardware.

Simulated devices are created with the `IPUModel` class, which models the functionality of an IPU on the host. The `createDevice` function creates a new virtual device to work with. Once we have this device we can create a `Graph` object to target it.

- Add the following code to the body of `main`:

```
// Create the IPU Model device
IPUModel ipuModel;
Device device = ipuModel.createDevice();
Target target = device.getTarget();

// Create the Graph object
Graph graph(target);
```

Any program running on an IPU needs data to work on. These are defined as variables in the graph.

- Add the following code to create the first variable in the program:

### Tutorial 5: a basic machine learning example

This tutorial contains a complete training program that performs a logistic regression on the MNIST data set, using gradient descent. The files for the demo are in `tut5_m1`. There are no coding steps in the tutorial. The task is to understand the code, build it and run it. You can build the code using the supplied makefile.

Before you can run the code you will need to run the `get_mnist.sh` script to download the MNIST data.

The program accepts an optional command line argument to make it use the IPU hardware instead of a simulated IPU.

As you would expect, training is significantly faster on the IPU hardware.

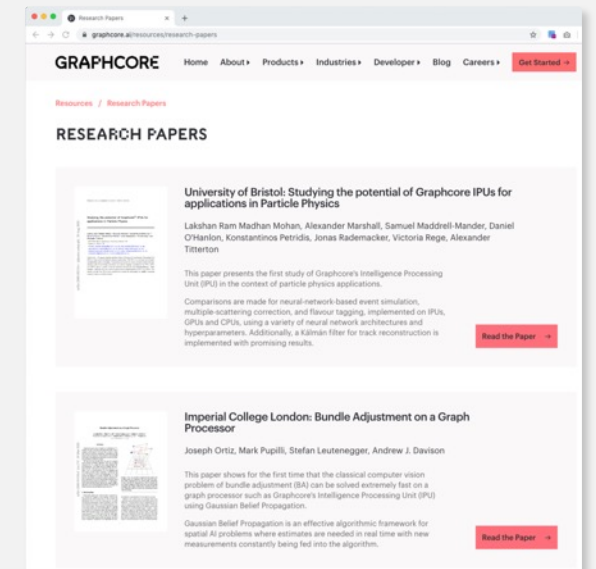
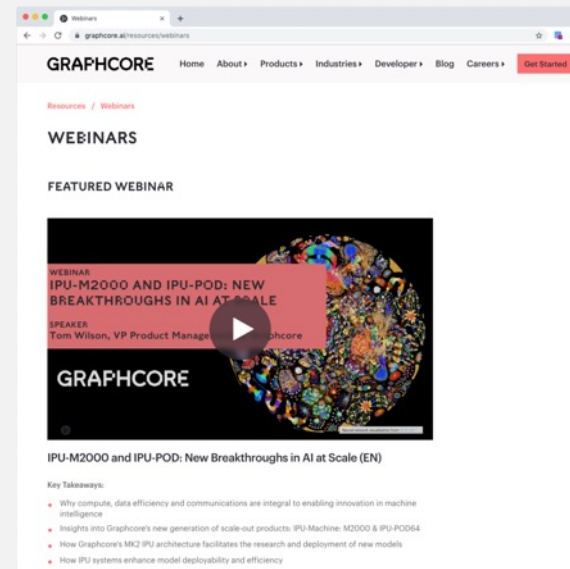
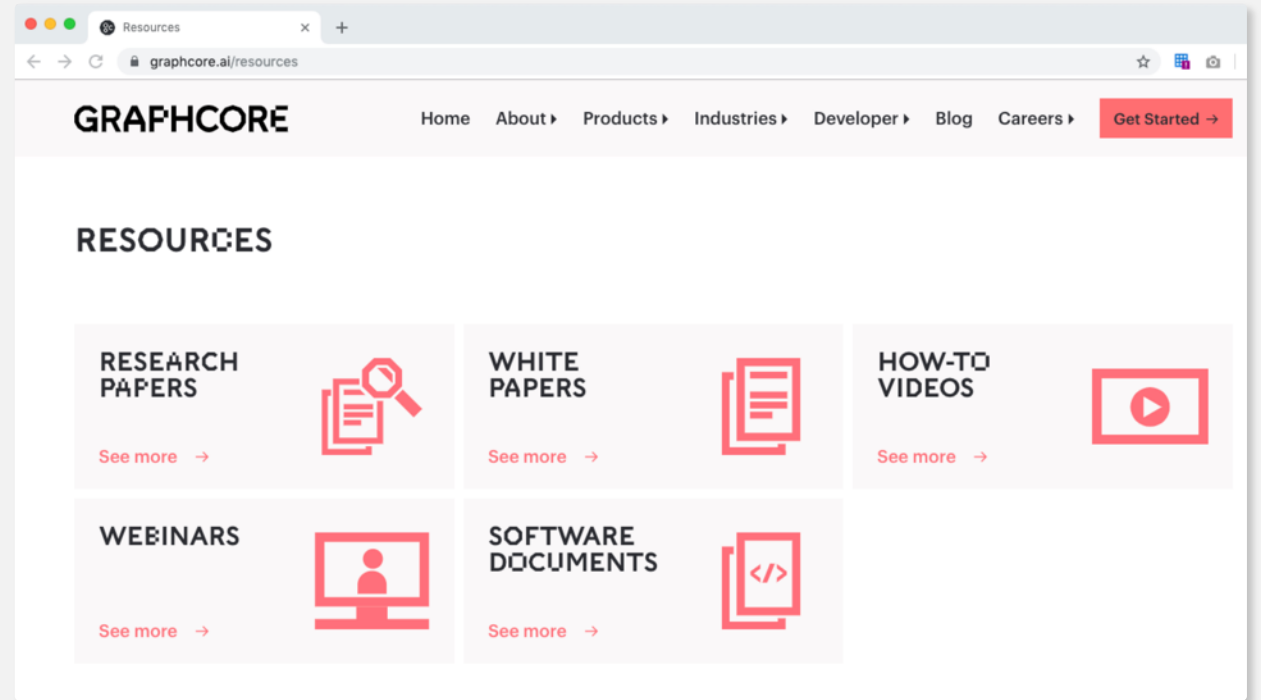
Copyright (c) 2018 Graphcore Ltd. All rights reserved.



# RESOURCES CENTRE

[graphcore.ai/resources](https://graphcore.ai/resources)

- Central source of research papers, white papers, videos, on-demand webinars and documentation
- Product resources for ML Engineers & IT / Infrastructure Managers now available





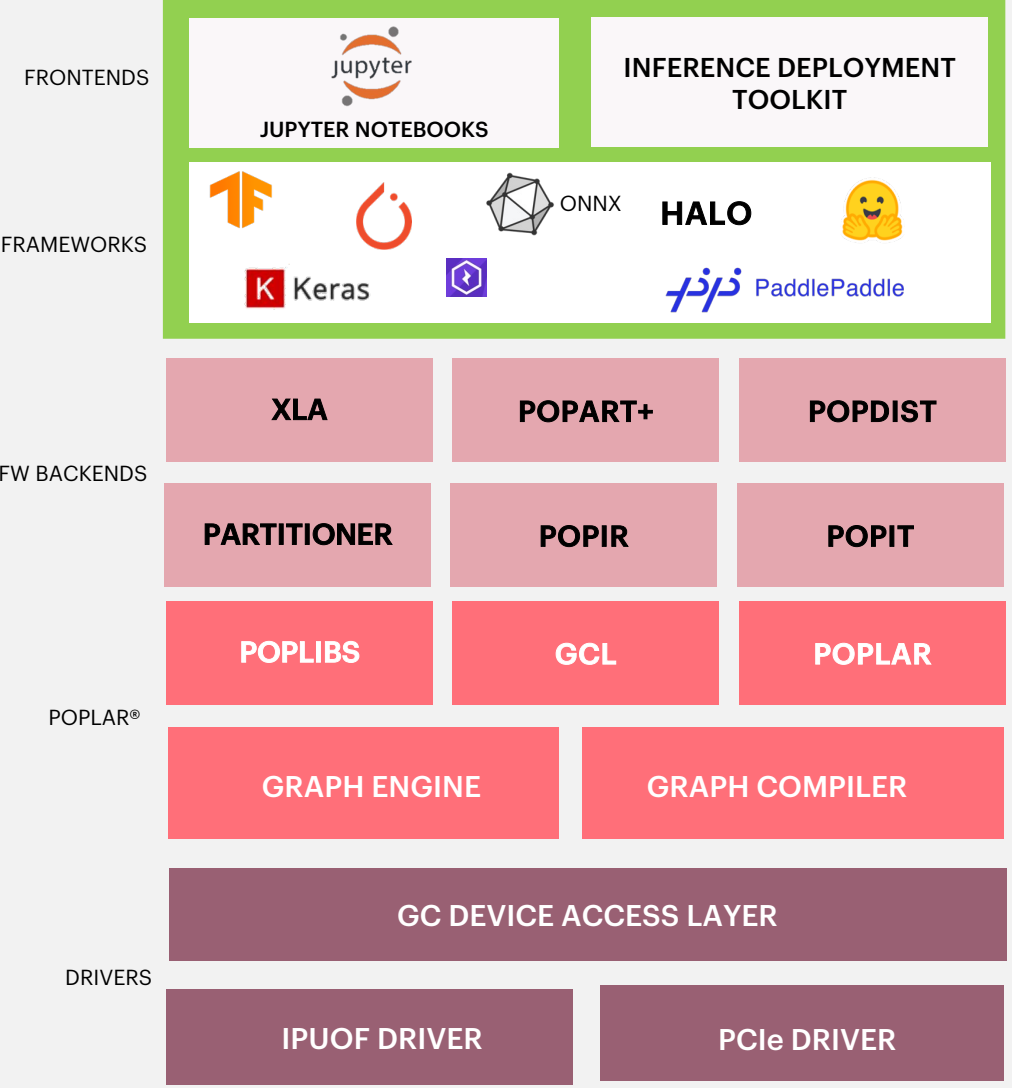
# GRAPHCORE SOFTWARE

- NLP/TRANSFORMERS
- IMAGE CLASSIFICATION/CNNS
- OBJECT DETECTION
- LARGE MODELS
- MLPERF
- CONDITIONAL SPARSITY
- GNNS

ML APPLICATIONS

- TUTORIALS
- CODE EXAMPLES
- DOCUMENTATION
- VIDEOS
- NATIVE IPU CODERS PROGRAM
- APPS PORTFOLIO

DEVELOPER ECOSYSTEM



POPLAR® SDK

- GRAPH ANALYZER
- SYSTEM ANALYZER
- DEBUGGER
- DEVELOPMENT ENVIRONMENT

POPVISION TOOLS

- V-IPU
- SYSTEM MONITORING
- PROMETHEUS
- GRAFANA

- JOB DEPLOYMENT
- K8S
- SLURM

SYSTEM SOFTWARE





# ENHANCED MODEL GARDEN

Resources > Model Garden

## MODEL GARDEN

PERFORMANT MODELS

**LIBRARY**

Type:

- Paperspace
- New
- Benchmarked
- Training
- Inference

Framework:

- PyTorch
- TensorFlow 1
- TensorFlow 2
- Hugging Face
- PopART
- PaddlePaddle
- Poplar

Category:

- Natural Language Processing >
- Computer Vision >
- Speech Processing >
- GNN
- Multimodal >
- AI for Simulation >
- Recommender >
- Probabilistic Modelling
- Reinforcement Learning
- Other

Search:

**GPS++ INFERENCE**

A hybrid GNN/Transformer for Molecular Property Prediction inference using IPUs trained on the PCQM4Mv2 dataset. Winner of the Open Graph Benchmark Large-Scale Challenge.

[Try on Paperspace](#) [View Repository](#)

**DISTRIBUTED KGE - TRANSE (256) TRAINING**

Knowledge graph embedding (KGE) for link-prediction training on IPUs using Poplar with the WikiKG90Mv2 dataset. Winner of the Open Graph Benchmark Large-Scale Challenge.

[View Repository](#)

**GPT-J 6B FINE-TUNING**

GPT-J 6B fine-tuned using the GLUE MLI dataset leveraging the Hugging Face Transformer library.

[View Repository](#)

**STABLE DIFFUSION TEXT-TO-IMAGE INFERENCE**

The popular latent diffusion model for generative AI with support for text-to-image on IPUs using Hugging Face Optimum.

[Try on Paperspace](#) [View Repository](#)

**STABLE DIFFUSION IMAGE-TO-IMAGE INFERENCE**

The popular latent diffusion model for generative AI with support for image-to-image on IPUs using Hugging Face Optimum.

[Try on Paperspace](#) [View Repository](#)

**STABLE DIFFUSION INPAINTING INFERENCE**

The popular latent diffusion model for generative AI with support for inpainting on IPUs using Hugging Face Optimum.

[Try on Paperspace](#) [View Repository](#)

**GPS++ TRAINING**

A hybrid GNN/Transformer for training Molecular Property Prediction using IPUs on the PCQM4Mv2 dataset. Winner of the Open Graph Benchmark Large-Scale Challenge.

[Try on Paperspace](#) [View Repository](#)

**GPS++ INFERENCE**

A hybrid GNN/Transformer for Molecular Property Prediction inference using IPUs trained on the PCQM4Mv2 dataset. Winner of the Open Graph Benchmark Large-Scale Challenge.

[Try on Paperspace](#) [View Repository](#)

**DISTRIBUTED KGE - TRANSE (256) TRAINING**

Knowledge graph embedding (KGE) for link-prediction training on IPUs using Poplar with the WikiKG90Mv2 dataset. Winner of the Open Graph Benchmark Large-Scale Challenge.

[View Repository](#)

**DISTRIBUTED KGE - TRANSE (256) INFERENCE**

Knowledge graph embedding (KGE) for link-prediction inference on IPUs using Poplar with the WikiKG90Mv2 dataset. Winner of the Open Graph Benchmark Large-Scale Challenge.

[View Repository](#)

**DISTRIBUTED KGE - TRANSE (256) TRAINING**

Knowledge graph embedding (KGE) for link-prediction training on IPUs using PyTorch with the WikiKG90Mv2 dataset. Winner of the Open Graph Benchmark Large-Scale Challenge.

[View Repository](#)

**GPT-J 6B FINE-TUNING**

GPT-J 6B fine-tuned using the GLUE MLI dataset leveraging the Hugging Face Transformer library.

[View Repository](#)

**DISTILBERT TRAINING**

DistilBERT is a small, fast, cheap and light Transformer-based language model.

**MAE TRAINING**

Implementation of MAE computer vision model in PyTorch for the IPU based on the open-MModel.

**FROZEN IN TIME TRAINING**

Implementation of Frozen in Time on the IPU in PyTorch for image-to-image inpainting.

PUBLIC ACCESS TO WIDE VARIETY OF MODELS, READY TO RUN ON IPU

NEW FILTER/SEARCH CAPABILITY

DIRECT ACCESS TO GITHUB

PAPERSPACE NOTEBOOK LINKS

<https://www.graphcore.ai/resources/model-garden>



# Access IPU-ready notebooks in seconds

## COMPUTER VISION

**Stable Diffusion 2**  
Text => Image Inference

**Stable Diffusion**  
Text => Image Inference

**Stable Diffusion**  
Image => Image Inference

**Stable Diffusion**  
Text Guided In-Painting Inference

**Training ViT HF model**  
ViT - Fine Tuning

**Image Classification**  
ViT - Fine Tuning

**Object Detection**  
YOLO v4 - Inference

**NEW**

## NATURAL LANGUAGE PROCESSING (NLP)

**BERT Fine Tuning**  
BERT-Large - Fine Tuning

**Fast sentiment analysis BERT |**  
RoBERTa - Inference

**Named Entity Recognition**  
BERT - Inference

**Multiple choice task**  
RoBERTa - Fine Tuning

**Question-Answering task**  
RoBERTa - Fine Tuning

**Instruction Tuned LLM**  
Dolly 2.0 - Inference

**NEW**

**Flan-T5 is all you need**  
Fine-Tuning & Inference

**NEW**

**Chatbot Open Source LLM**  
Llama 2 - Inference

**NEW**

**Summarization task**  
T5 Small - Fine Tuning

**Text Classification task**  
RoBERTa - Fine Tuning

**Token Classification task**  
BERT

**Translation task**  
BART-Base

**Text Entailment**  
GPT-J - Fine Tuning

**Text Generation**  
GPT-J

**Chatbot using OpenAssistant**  
Pythia 12B - Inference

**NEW**

**SQuAD & MNLI**  
DeBERTa - Inference

**NEW**

## GNN

**Training Large Graphs**  
Cluster-GCN - Training

**Training Dynamic Graphs**  
TGN - Training

**Predicting molecular properties**  
SchNet - Training

**Predicting molecular properties**  
GIN - Training

**Predicting molecular properties**  
GPS++ (OGB-LSC) Train | Inf

**Link Prediction training**  
Dist KGE (OGB-LSC) - Train

**Link Prediction training**  
NBFNet - Training

**Molecular Prediction**  
MolFeat - FT & Inf

**Molecular Modelling with Graphium**  
GCN/GIN - Train & Inf

**NEW**

## SPEECH PROCESSING

**Fine-Tuning a wav2vec**  
wav2vec - Fine Tuning

**Running ASR**  
wav2vec - Inference

**Speech Transcription on IPU**  
Whisper - FT & Inf

**NEW**

## TIME SERIES

**Multi-horizon financial forecasting**  
DeepLOB Seq2Seq - Train & Inf

**NEW**

**Multi-horizon financial forecasting**  
DeepLOB Attention - Train & Inf

**NEW**



# USEFUL ENV VARIABLES



# LOGGING

Logging messages can be generated when your program runs. This is controlled by the environment variables described below. For more detailed information see the docs:

<https://docs.graphcore.ai/projects/poplar-user-guide/en/latest/env-vars.html>

POPLAR\_LOG\_LEVEL: Enable logging for Poplar

POPLAR\_LOG\_DEST: Specify the destination for Poplar logging (“stdout”, “stderr” or a file name)

“OFF”	No logging information. The default.
“ERR”	Only error conditions will be reported.
“WARN”	Warnings when, for example, the software cannot achieve what was requested (for example, if the convolution planner can’t keep to the memory budget, or Poplar has determined that the model won’t fit in memory but the debug.allowOutOfMemory option is enabled).
“INFO”	Very high level information, such as PopLibs function calls.
“DEBUG”	Useful per-graph information.
“TRACE”	The most verbose level. All useful per-tile information.

# CREATE EXECUTION PROFILE

```
POPLAR_ENGINE_OPTIONS='{ "autoReport.all": "true", "autoReport.directory": "./report" }'
```

- The PopVision Graph Analyser uses report files generated during compilation and execution by the Poplar SDK.
- These files can be created using POPLAR\_ENGINE\_OPTIONS.
- In order to capture the reports needed for the PopVision Graph Analyser you only need to set POPLAR\_ENGINE\_OPTIONS='{ "autoReport.all": "true" }' before you run a program. By default this will enable instrumentation and capture all the required reports to the current working directory.

# EXECUTABLE CACHE

If you often run the same models you might want to enable executable caching to save time:

POPTORCH:

- You can do this by either setting the POPTORCH\_CACHE\_DIR environment variable or by calling poptorch.Options.enableExecutableCaching.

TENSORFLOW:

- You can use the flag `--executable_cache_path` to specify a directory where compiled files will be placed. Fused XLA/HLO graphs are hashed with a 64-bit hash and stored in this directory.

## Warning

The cache directory might grow large quickly. Poplar doesn't evict old models from the cache and, depending on the number and size of your models and the number of IPUs used, the executables might be quite large.

It is your responsibility to delete the unwanted cache files.



# SYNTHETIC-DATA

```
TF_POPLAR_FLAGS= "--use_synthetic_data --synthetic_data_initializer=random"
```

Used for measuring the IPU-only throughput and disregards any host/CPU activity.

# GRAPHCORE COMMAND LINE TOOLS

- gc-docker*** Allows you to use IPU devices in Docker containers using the Docker container engine.
- gc-flops*** Allows you to benchmark the number of floating point operations per second on one or more IPU processors.
- gc-info*** Determines what IPU cards are present in the system.
- gc-inventory*** Lists device IDs, physical parameters and firmware version numbers.
- gc-links*** Displays the status and connectivity of each of the IPU-Links that connect IPUs. See also *IPU-Link channel mapping* for connectivity in an IPU Server containing C2 cards.
- gc-monitor*** Monitors IPU activity on shared systems.
- gc-reset*** Resets IPU devices.
- gc-exchangetest*** Allows you to test the internal exchange fabric in an IPU.
- gc-exchangewritetest*** Tests direct writes to the IPU's tile memory via the host.
- gc-gwlinkstraffictest*** Tests GW-Links on multi-rack IPU-POD systems.
- gc-hostsynclatencytest*** Reports the latency of transfers between the host machine and the IPUs (in both directions).
- gc-hosttraffictest*** Allows you to test the data transfer between the host machine and the IPUs (in both directions).
- gc-iputraffictest*** Allows you to test the data transfer between IPUS.
- gc-memorytest*** Tests all the memory in an IPU, reporting any tiles that fail.
- gc-podman*** Allows you to use IPU devices in Docker containers using the Podman container engine.
- gc-powertest*** Tests power consumption and temperature of the IPU processors.



# TF2/KERAS ON IPU

LSTM Encoder Decoder



# KERAS ON IPU

IPU optimized Keras `Model` and `Sequential` with the following features:

- On-device training loop for reduction of communication overhead.
- Gradient accumulation for simulating larger batch sizes.
- Automatic data-parallelisation of the model when placed on a multi-IPU device.



## Keras

```
import tensorflow as tf
from tensorflow.keras.layers import *
```

## GPU

```
(x_train, y_train), (x_test, y_test) = tf.keras.datasets.cifar10.load_data()
x_train = x_train.astype('float32') / 255.0
y_train = tf.keras.utils.to_categorical(y_train, 10)
ds_train = tf.data.Dataset.from_tensor_slices((x_train, y_train)).batch(64, drop_remainder=True)
```

```
model = tf.keras.Sequential([
    Conv2D(32, (3, 3), padding='same', input_shape=x_train.shape[1:]),
    Activation('relu'),
    Conv2D(32, (3, 3)),
    Activation('relu'),
    MaxPooling2D(pool_size=(2, 2)),
    Dropout(0.25),
    Conv2D(64, (3, 3), padding='same'),
    Activation('relu'),
    Conv2D(32, (3, 3)),
    Activation('relu'),
    MaxPooling2D(pool_size=(2, 2)),
    Dropout(0.25),
    Flatten(),
    Dense(512),
    Activation('relu'),
    Dropout(0.5),
    Dense(10),
    Activation('softmax')
])
```

```
model.compile(loss='categorical_crossentropy',
              optimizer=tf.optimizers.SGD(learning_rate=0.016),
              metrics=['accuracy'])
```

```
model.fit(ds_train, epochs=40)
```

```
import tensorflow as tf
from tensorflow.keras.layers import *
+ from tensorflow.python import ipu
```

## IPU

```
+ cfg = ipu.config.IPUConfig()
+ cfg.auto_select_ipus = 1
+ cfg.configure_ipu_system()
+ with ipu.ipu_strategy.IPUStrategy().scope():
    (x_train, y_train), (x_test, y_test) = tf.keras.datasets.cifar10.load_data()
    x_train = x_train.astype('float32') / 255.0
    y_train = tf.keras.utils.to_categorical(y_train, 10)
    ds_train = tf.data.Dataset.from_tensor_slices((x_train, y_train)).batch(64, drop_remainder=True)
```

```
model = tf.keras.Sequential([
    Conv2D(32, (3, 3), padding='same', input_shape=x_train.shape[1:]),
    Activation('relu'),
    Conv2D(32, (3, 3)),
    Activation('relu'),
    MaxPooling2D(pool_size=(2, 2)),
    Dropout(0.25),
    Conv2D(64, (3, 3), padding='same'),
    Activation('relu'),
    Conv2D(32, (3, 3)),
    Activation('relu'),
    MaxPooling2D(pool_size=(2, 2)),
    Dropout(0.25),
    Flatten(),
    Dense(512),
    Activation('relu'),
    Dropout(0.5),
    Dense(10),
    Activation('softmax')
])
```

```
model.compile(loss='categorical_crossentropy',
              optimizer=tf.optimizers.SGD(learning_rate=0.016),
              metrics=['accuracy'])
```

```
model.fit(ds_train, epochs=40)
```

# TF2/KERAS TUTORIALS

Continued in the repositories below (follow the READMEs)

[github.com/graphcore/examples/tree/master/tutorials/tutorials/tensorflow2/keras](https://github.com/graphcore/examples/tree/master/tutorials/tutorials/tensorflow2/keras)

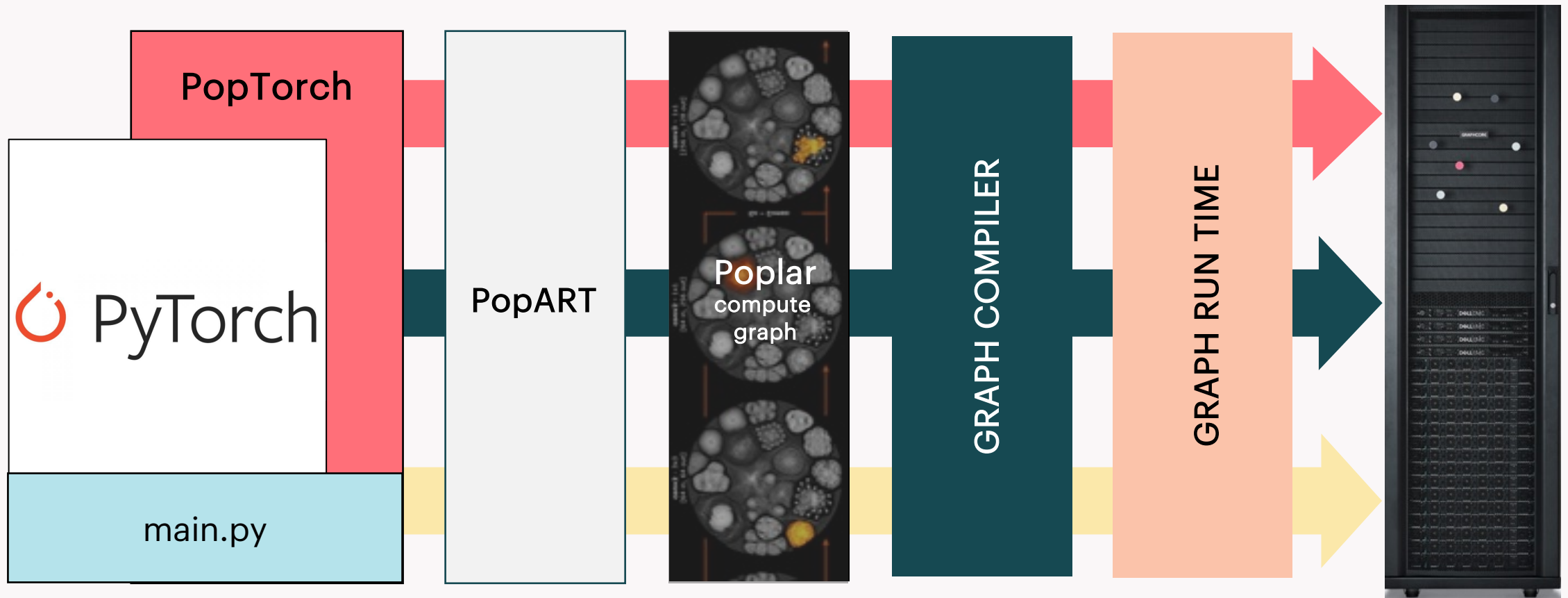


# INTRO TO POPTORCH

GRAPHCORE



# WHAT IS POPTORCH?





# WHAT IS POPTORCH?

- PopTorch is a set of extensions for PyTorch to enable PyTorch models to run on Graphcore's IPU hardware.
- PopTorch supports both inference and training. To run a model on the IPU you wrap your existing PyTorch model in either a PopTorch inference wrapper or a PopTorch training wrapper.
- You can provide further annotations to partition the model across multiple IPUs. Using the user-provided annotations, PopTorch will use PopART to parallelise the model over the given number of IPUs.
- Additional parallelism can be expressed via a replication factor which enables you to data-parallelise the model over more IPUs.

# PYTORCH FOR IPU



Define a model within  
PyTorch



Create an IPU execution  
wrapper around the  
model and run as normal



PopTorch uses  
PyTorch dispatcher to  
trace the model



Compile the graph in  
PopART and then run on  
one or more IPU

# GETTING STARTED: TRAINING A MODEL



# TRAINING A MODEL

## 1. Import packages

PopTorch is a separate package from PyTorch, and must be imported.

## 2. Load dataset using torchvision.datasets and poptorch.DataLoader

In order to make data loading easier and more efficient, PopTorch offers an extension of `torch.utils.data.DataLoader` class:

`poptorch.DataLoader` class is specialised for the way the underlying PopART framework handles batching of data.

## 3. Define model and loss function using torch API

The only difference here from pure PyTorch is the loss computation, which has to be part of the forward function. This is to ensure the loss is computed on the IPU and not on the CPU, and to give us as much flexibility as possible when designing more complex loss functions.



# TRAINING A MODEL

## 4. Prepare training

Instantiate compilation and execution options, these are used by PopTorch's wrappers such as `poptorch.DataLoader` and `poptorch.trainingModel`.

## 5. Train the model

Define the optimizer using PyTorch's API.

Use `poptorch.trainingModel` wrapper, to wrap your PyTorch model. This wrapper will trigger the compilation of our model, using TorchScript, and manage its translation to a program the IPU can run. Then run your training loop.



# PyTorch

GPU

```
_, ind = torch.max(predictions, 1)
# provide labels only for samples, where prediction is available (during the training, not
predictions.size()[0]:]
torch.eq(ind, labels)).item() / labels.size(0)

if __name__ == '__main__':
    parser = argparse.ArgumentParser(description='MNIST training in PopTorch')
    parser.add_argument('--batch-size', type=int, default=8, help='batch size for training (default: 8)')
    parser.add_argument('--test-batch-size', type=int, default=8, help='batch size for testing (default: 8)')
    parser.add_argument('--epochs', type=int, default=10, help='number of epochs to train (default: 10)')
    parser.add_argument('--lr', type=float, default=0.05, help='learning rate (default: 0.05)')

    args = parser.parse_args()

    training_data = torch.utils.data.DataLoader(
        torchvision.datasets.MNIST('mnist_data/', train=True, download=True, transform=None,
        batch_size=args.batch_size, shuffle=True, drop_last=True)

    test_data = torch.utils.data.DataLoader(
        torchvision.datasets.MNIST('mnist_data/', train=False, download=True, transform=None,
        batch_size=args.batch_size, shuffle=True, drop_last=True)

    model = Network()
    training_model = TrainingModelWithLoss(model)
    optimizer=optim.SGD(model.parameters(), lr=args.lr)

    # Run training
    for _ in range(args.epochs):
        for data, labels in training_data:
            preds, losses = training_model(data, labels)
            optimizer.zero_grad()
            losses.backward()
            optimizer.step()

    # Run validation
    sum_acc = 0.0
    with torch.no_grad():
        for data, labels in test_data:
            output = model(data)
            sum_acc += accuracy(output, labels)
    print("Accuracy on test set: {:.2f}%".format(sum_acc / len(test_data)))
```

```
_, ind = torch.max(predictions, 1)
# provide labels only for samples, where prediction is available (during the training, not
labels = labels[-predictions.size()[0]:]
accuracy = torch.sum(torch.eq(ind, labels)).item() / labels.size(0)
return accuracy
```

IPU

```
if __name__ == '__main__':
    parser = argparse.ArgumentParser(description='MNIST training in PopTorch')
    parser.add_argument('--batch-size', type=int, default=8, help='batch size for training (default: 8)')
    parser.add_argument('--test-batch-size', type=int, default=8, help='batch size for testing (default: 8)')
    parser.add_argument('--epochs', type=int, default=10, help='number of epochs to train (default: 10)')
    parser.add_argument('--lr', type=float, default=0.05, help='learning rate (default: 0.05)')
    parser.add_argument('--device-iterations', type=int, default=50, help='device iterations (default: 50)')

    args = parser.parse_args()

    + opts = poptorch.Options().deviceIterations(args.device_iterations)
    + training_data = poptorch.DataLoader(opts,
        torchvision.datasets.MNIST('mnist_data/', train=True, download=True, transform=None,
        batch_size=args.batch_size, shuffle=True, drop_last=True)

    + test_data = poptorch.DataLoader(opts,
        torchvision.datasets.MNIST('mnist_data/', train=False, download=True, transform=None,
        batch_size=args.batch_size, shuffle=True, drop_last=True)

    model = Network()
    training_model = TrainingModelWithLoss(model)
    optimizer=optim.SGD(model.parameters(), lr=args.lr)
    + training_model = poptorch.trainingModel(training_model, opts, optimizer=optimizer)
    + inference_model = poptorch.inferenceModel(model)

    # Run training
    for _ in range(args.epochs):
        for data, labels in training_data:
            preds, losses = training_model(data, labels)

    + # Detach the training model so that the same IPU could be used for validation
    + training_model.detachFromDevice()

    # Run validation
    sum_acc = 0.0
    with torch.no_grad():
        for data, labels in test_data:
            output = inference_model(data)
            sum_acc += accuracy(output, labels)
    print("Accuracy on test set: {:.2f}%".format(sum_acc / len(test_data)))
```

# POPTORCH.OPTIONS

- The compilation and execution on the IPU can be controlled using `poptorch.Options`
- Full list of options available here: <https://docs.graphcore.ai/projects/poptorch-user-guide/en/latest/overview.html#options>

- Some examples:

(i) **deviceIterations**

This option specifies the number of batches that is prepared by the host (CPU) for the IPU. The higher this number, the less the IPU has to interact with the CPU, for example to request and wait for data, so that the IPU can loop faster. However, the user will have to wait for the IPU to go over all the iterations before getting the results back. The maximum is the total number of batches in your dataset, and the default value is 1.

(ii) **replicationFactor**

This is the number of replicas of a model. We use replicas as an implementation of data parallelism. To achieve the same behavior in pure PyTorch, you'd wrap your model with `torch.nn.DataParallel`, but with PopTorch, this is an option.



# POPTORCH TUTORIALS

Continued in the repositories below (follow the READMEs)

[github.com/graphcore/examples/tree/master/tutorials/tutorials/pytorch/basics](https://github.com/graphcore/examples/tree/master/tutorials/tutorials/pytorch/basics)

[github.com/graphcore/examples/tree/master/tutorials/tutorials/pytorch/mixed\\_precision](https://github.com/graphcore/examples/tree/master/tutorials/tutorials/pytorch/mixed_precision)

[github.com/graphcore/examples/tree/master/tutorials/tutorials/pytorch/efficient\\_data\\_loading](https://github.com/graphcore/examples/tree/master/tutorials/tutorials/pytorch/efficient_data_loading)

[github.com/graphcore/examples/tree/master/tutorials/tutorials/pytorch/pipelining](https://github.com/graphcore/examples/tree/master/tutorials/tutorials/pytorch/pipelining)





# GRAPHIUM FOR IPU

Graphium integrates state-of-the-art Graph Neural Network (GNN) architectures and a user-friendly API, enabling the easy construction and training of custom GNN models.




A POWERFUL AND  
FLEXIBLE OPEN-  
SOURCE PYTHON  
LIBRARY FOR TRAINING  
MOLECULAR GNNS AT  
SCALE

ANNOUNCEMENT | TECHNICAL BLOG | GETTING STARTED

Jul 05, 2023  
**GRAPHIUM: AN IPU-READY PYTHON LIBRARY FOR TRAINING MOLECULAR GNNS AT SCALE**  
Written By:  
Dominique Beaini

Jul 05, 2023  
**MULTITASK MOLECULAR MODELLING WITH GRAPHIUM ON THE IPU**  
Written By:  
Sam Maddrell-Mander

RUN GRAPHIUM ON IPU WITH PAPERSPACE JUPYTER NOTEBOOK

<b>Graphium</b>	<b>Domain:</b> Molecules
Multitask Molecular Modelling on the IPU	<b>Tasks:</b> Multitask
	<b>Model:</b> GCN/GIN/GINE
	<b>Datasets:</b> QM9, Zinc, Tox21
	<b>Workflow:</b> Training, validation, inference
	<b>Execution time:</b> 20 mins



# PyG

## PyG is the ultimate library for Graph Neural Networks

Build graph learning pipelines with ease



GRAPHCORE



HARVARD MEDICAL SCHOOL

pyg.org

“The suitability of IPU for running GNNs and the kind of performance advantage that Graphcore and its customers have demonstrated is really helping to accelerate the uptake of this exciting model class”

**Matthias Fey – PyG creator & founder of Kumo.ai**

# PYTORCH GEOMETRIC FOR IPU

## ANNOUNCEMENT | TECHNICAL BLOG | GETTING STARTED

Apr 05, 2023 | Poplar, PyTorch, GNN

**GRAPHCORE USERS CAN NOW BUILD AND RUN GNNs WITH PYTORCH GEOMETRIC**

Written By: **Blazej Banaszewski**

Apr 05, 2023 | Developer, PyTorch, GNN

**ACCELERATING PYG ON IPUS: UNLEASH THE POWER OF GRAPH NEURAL NETWORKS**

Written By: **Blazej Banaszewski, Adam Sanders, Akash Swamy & Mihai Poplar**

Apr 05, 2023 | Developer, GNN, Paperspace  
**GETTING STARTED WITH PYTORCH GEOMETRIC (PYG) ON GRAPHCORE IPUS**

Written By: **Adam Sanders and Arianna Saracino**

```
import torch
import torch.nn.functional as F
from torch_geometric.nn import GCNConv

class GCN(torch.nn.Module):
    def __init__(self, in_channels, out_channels):
        super().__init__()
        torch.manual_seed(1234)
        self.conv = GCNConv(in_channels, out_channels, add_self_loops=False)

    def forward(self, x, edge_index, edge_weight=None):
        x = F.dropout(x, p=0.5, training=self.training)
        x = self.conv(x, edge_index, edge_weight).relu()
        return x

model = GCN(dataset.num_features, dataset.num_classes)
model.train()
optimizer = optim.Adam(model.parameters(), lr=0.001)
print("Training on CPU")
for epoch in range(1, 5):
    optimizer.zero_grad()
```

## RUN GNN MODELS IN PYG ON PAPERSPACE JUPYTER NOTEBOOKS

**Training Dynamic Graphs**

TGN Training



Run on Gradient

**Training Large Graphs**

Cluster-GCN Training



Run on Gradient

**Predicting molecular properties**

SchNet Training



Run on Gradient

**Predicting molecular properties**

GIN Training



Run on Gradient

**Link Prediction training**

NBFNet Training



Run on Gradient

**Molecular Modelling with Graphium**

GCN/GIN Training & Inf



Run on Gradient

NEW

# POPVISION™ TOOLS

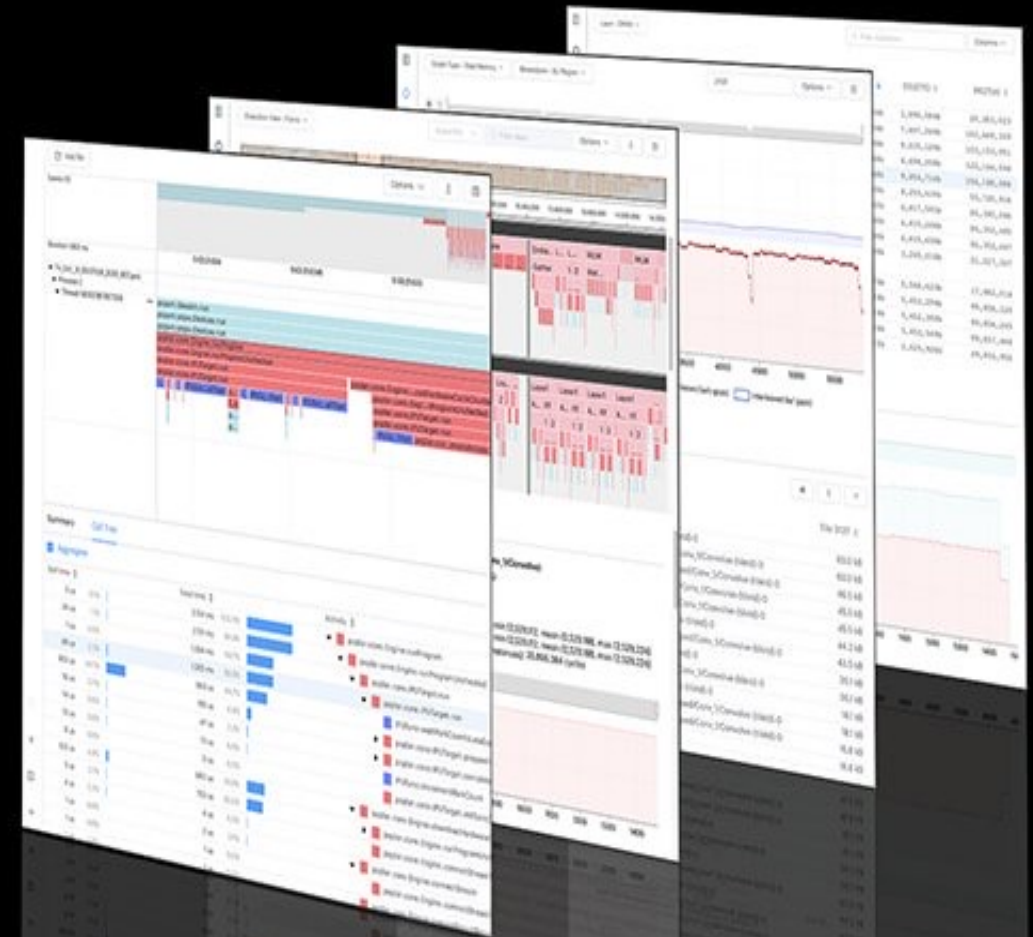
POPLAR™ POPVISION TOOLS

## GRAPH ANALYSER

Useful for analysing and optimising the memory use and execution performance of ML models on the IPU

## SYSTEM ANALYSER

Graphical view of the timeline of host-side application execution steps



“Our team was very impressed by the care and effort Graphcore has clearly put into the PopVision graph and system analysers. It’s hard to imagine getting such a helpful and comprehensive profiling of the code elsewhere, so this was really a standout feature in our IPU experience.”



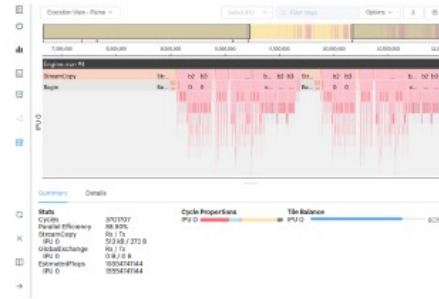
Dominique Beaini, Valence Discovery, a leader in AI-first drug design

# POPVISION TOOLS



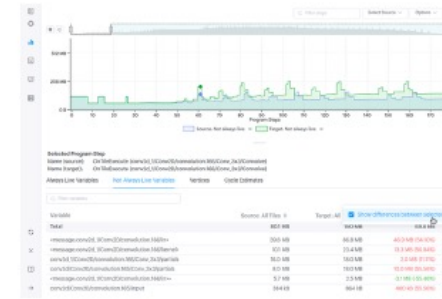
## IPU MEMORY ANALYSIS

Capture memory information from your ML models when executed on IPUUs. Inspect variable placement, size and liveness throughout the execution.



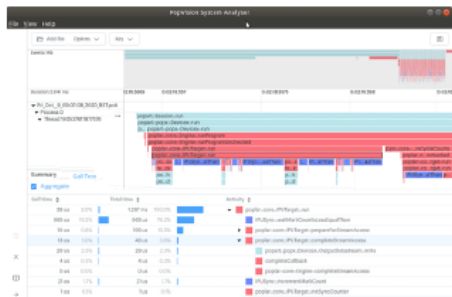
## EXECUTION TRACE REPORT

View the output of instrumenting a Poplar program, capturing cycle counts for each step. See execution statistics, tile balance, cycle proportions and compute-set details.



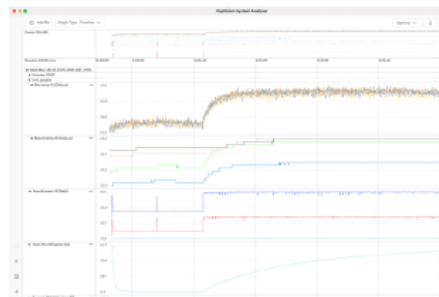
## REPORT COMPARISONS

Open two reports at once to compare their memory, execution, liveness and operations. Visualise where efficiencies can be made with different model parameters.



## HOST EXECUTION ANALYSIS

Understand the execution of IPU-targeted software on your host system processors. Identify any bottlenecks between CPUs and IPUUs across a visual interactive timeline.



## GRAPH DATA

Plot graph data of any numerical data points from the host or IPU processor systems, such as board temperature, power consumption and IPU utilisation.



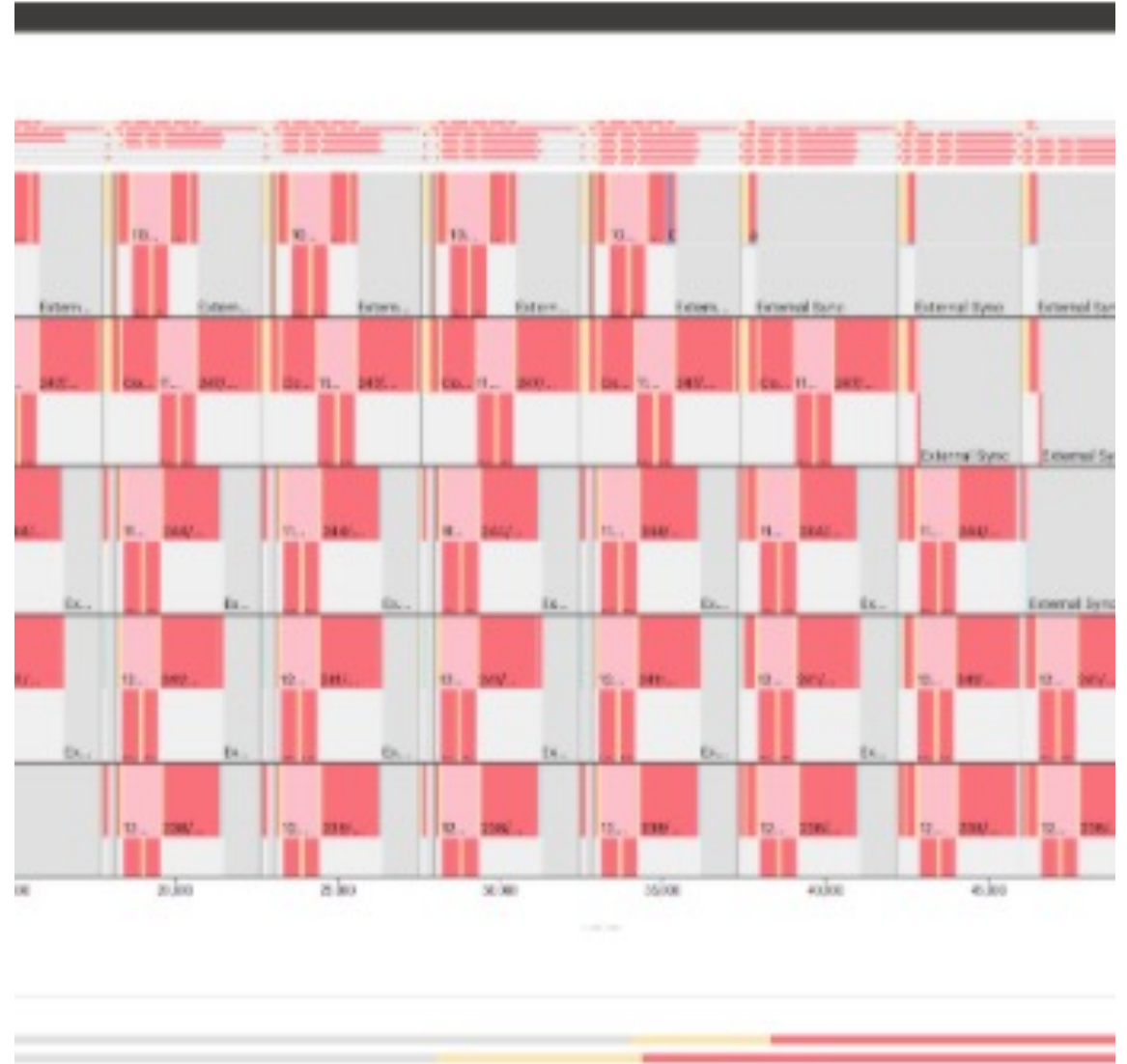
## LOCAL + REMOTE REPORTS

Ability to open reports either on your local machine, or remotely on the host machine. The Graph Analyser also supports local and remote report access.



# POPVISION GRAPH ANALYSER

- You can use the PopVision Graph Analyser tool to debug IPU programs and generate reports on compilation and execution of the program.
- This tool can be downloaded from [graphcore.ai/downloads](https://graphcore.ai/downloads)
- There is a built-in help system within the tool for any questions you might have about producing and analysing reports.



# PopVision Graph Analyser

## Getting started with PopVision™

Intro to the PopVision™ Graph Analyser



Getting started video available on the developers portal



Several new features including:

- A new file format for the graph and execution profile, resulting in a 50% file size reduction
- Enhanced PopLibs debug information

## Liveness Report

The debug information shown for a variable now displays enhanced information. For each variable that has debug information, you can now see the PopLibs API that created it, its arguments and its outputs.

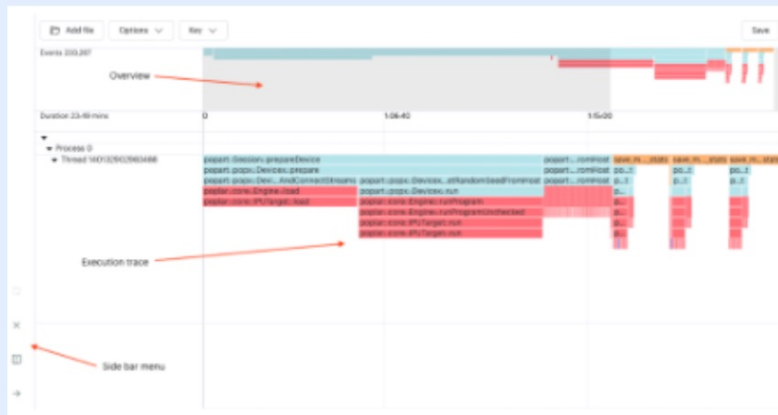
Enhanced debug information has been added to program steps. Program steps show Poplar and PopLibs debug information such as which PopLibs API created that program step, its arguments and its outputs.

Check out the integrated help or visit our developer portal for more information

# PopVision System Analyser

The PopVision System Analyser allows developers to understand the execution of programs running on the host processor which control the IPU(s). The System Analyser shows the interaction between the host and the IPU(s) so that developers can understand where the bottlenecks are in the execution of their applications.

The PopVision System Analyser visualises the information collected by the PopVision Trace Instrumentation Library which is part of the Poplar SDK.



Show the execution of the software on the host processor enabling users to identify bottlenecks in execution between CPU & IPU(s).

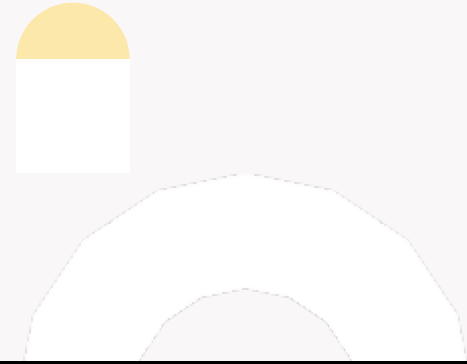


Provide profile insights as you scale models to multiple CPUs / IPUs.

Visit our developer portal for more information and the latest documentation:

<https://www.graphcore.ai/developer>

# CREATE PROFILE



```
(3.3.0+1403_poptorch) alext@bpod16:~/work/examples/tutorials/tutorials/pytorch/basics$ POPLAR_ENGINE_OPTIONS='{ "autoReport.all": "true", "autoReport.directory": "./report", "debug.allowOutOfMemory": "true", "profiler.includeFlopEstimates": "true" }' python walkthrough.py
epochs: 0% | 0/5 [00:00<?, ?it/s]
22:25:58.232 [poptorch:cpp] [warning] [DISPATCHER] Type coerced from Long to Int for tensor id 12 | 0/3750 [00:00<?, ?it/s]
Graph compilation: 100% | 100/100 [00:39<00:00]
epochs: 100% | 5/5 [01:21<00:00, 16.21s/it]
/nethome/alex/venvs/poplar_sdk-ubuntu_20_04-3.3.0+1403-208993bbb7/3.3.0+1403_poptorch/lib/python3.8/site-packages/torch/nn/modules/module.py:1802: UserWarning: Positional args are being deprecated, use kwargs instead. Refer to https://pytorch.org/docs/master/generated/torch.nn.Module.html#torch.nn.Module.state_dict for details.
  warnings.warn(
Graph compilation: 100% | 100/100 [00:19<00:00]
Eval accuracy: 89.25%
Graph compilation: 100% | 100/100 [00:09<00:00]
tensor([[ -4.4754, -0.9084, -3.2791, -5.8905, -3.4738, -1.4554, -1.8643, -6.6169,
          -2.1170, -5.5300]])
IPU predicted class: Trouser
CPU predicted class: Trouser
```





Summary

Insights

Memory Report

Liveness Report

Program Tree

Operations Summary

Operations Graph

Execution Trace

Reload Report

Close Report

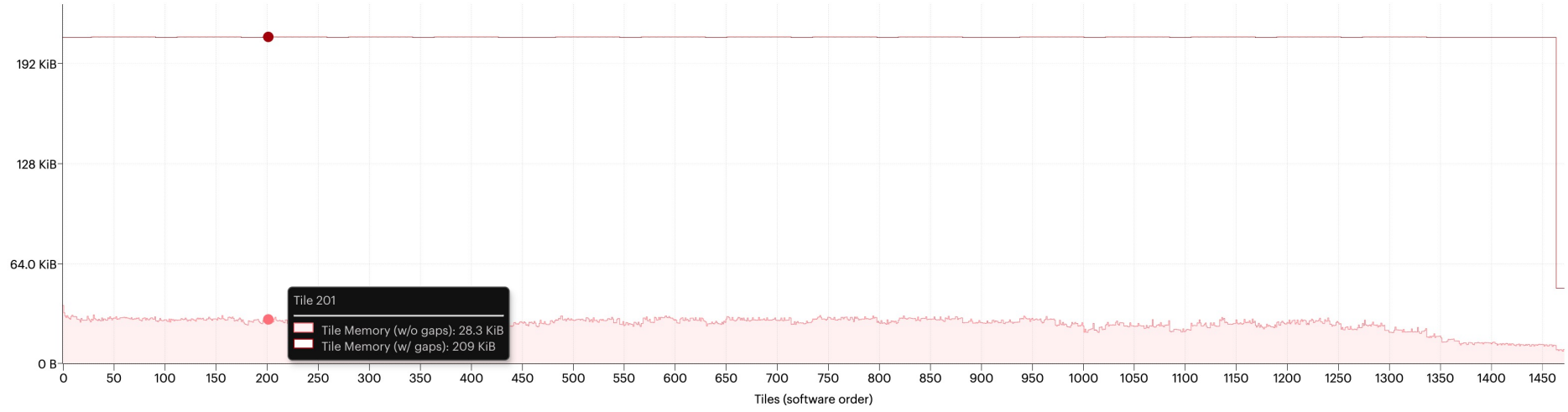
Documentation

Minimise

Graph Type - Total Memory Breakdown - None

201 Options

Navigation 100%



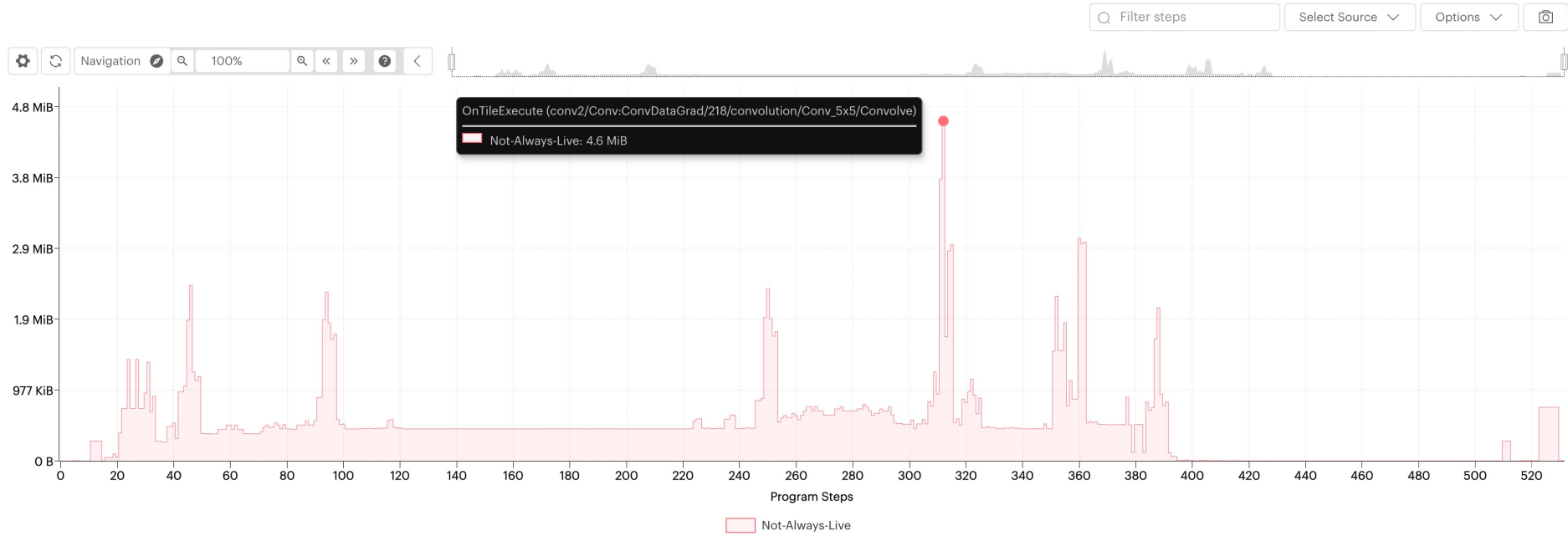
Tile Memory (w/o gaps) Tile Memory (w/ gaps)

Details Compute Sets Vertices Exchanges Variables

	Tile 201
Physical Tile Id	522
Memory Including Gaps	209 KiB
Memory Excluding Gaps	28.3 KiB
By Memory Region	
Non-interleaved	27.4 KiB
Interleaved	960 B
Overflowed	0 B
By Data Type	
Not Overlapped	
Variables	1.6 KiB
Constants	414 B
Control Code	2.4 KiB
Vertex Code	8.8 KiB
Internal Exchange Code	7.2 KiB
Host Exchange Code	640 B
Global Exchange Code	0 B
Vertex Instances	1.4 KiB



- Summary
- Insights
- Memory Report
- Liveness Report
- Program Tree
- Operations Summary
- Operations Graph
- Execution Trace



**Selected Program Steps**

Step 312 OnTileExecute (conv2/Conv:ConvDataGrad/218/convolution/Conv\_5x5/Convolve)

- Always-Live Variables
- Not-Always-Live Variables
- Vertices
- Cycle Estimates
- FLOP Estimates

Filter variables

Variable	All Tiles
<b>Total</b>	<b>29.3 MiB</b>
vertexCode	11.8 MiB
internalExchangeCode	8.5 MiB
controlCode	3.0 MiB
vertexInstanceState	1.9 MiB
hostExchangeCode	924 KiB
stack	535 KiB
instrumentationResults	402 KiB
▶ Accl1_fc1.weight	380 KiB
▶ fc1.weight	380 KiB
vertexFieldData	241 KiB
copyDescriptor	223 KiB
hostExchangeParketHeader	198 KiB

- Reload Report
- Close Report
- Documentation
- Minimise



Summary

Insights

Memory Report

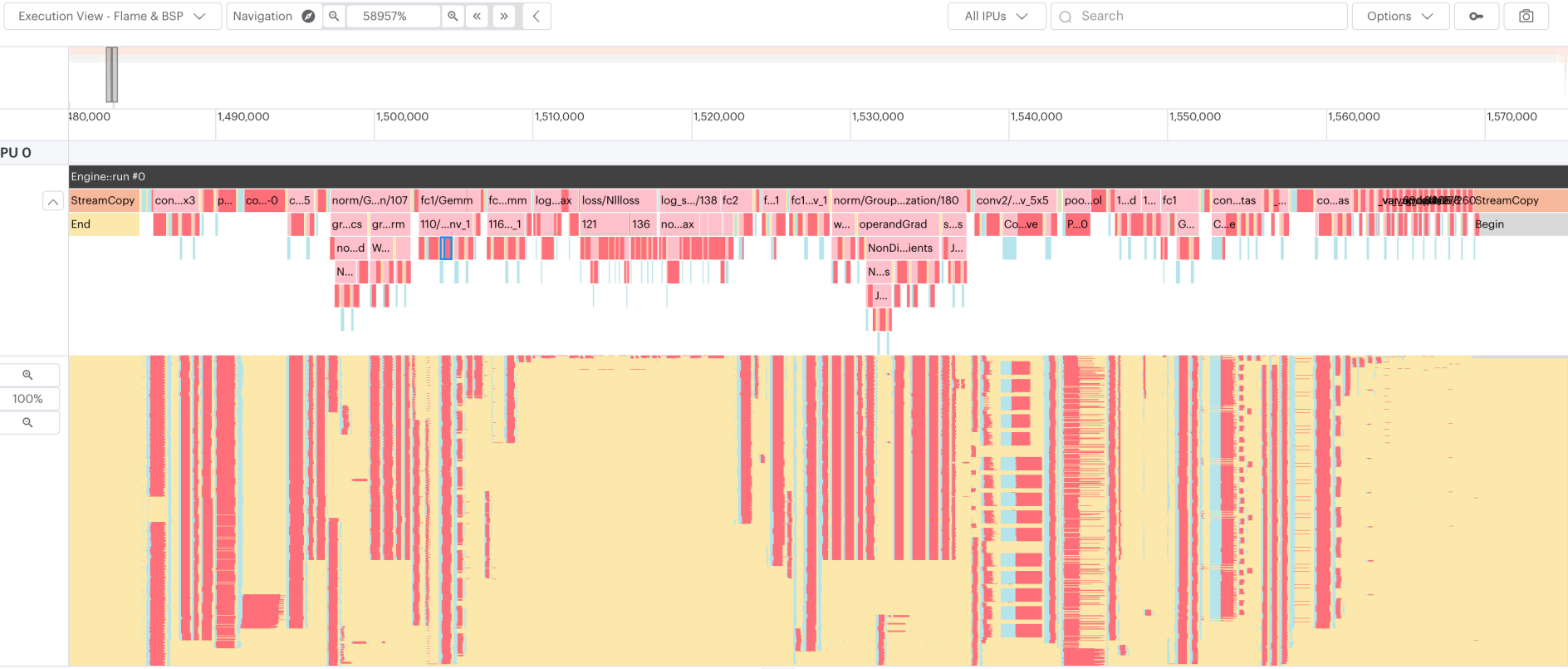
Liveness Report

Program Tree

Operations Summary

Operations Graph

Execution Trace



Summary

Details

Total Cycles 94,526 Total FLOPs 31,844,509

IPU	Proportions	Tile Utilisation	Stream Copy (Rx/Tx)	Global Exchange (Rx/Tx)
0		59.16%	0 B / 0 B	0 B / 0 B

! One or more step types are hidden. These step types are excluded from the tile balance, tile utilisation, and cycle proportions shown and so the statistics may not be representative. Step types can be shown or hidden using the view options dropdown.

Reload Report

Close Report

Documentation

Minimise



# POPVISION TUTORIALS

Continued in the repositories below (follow the READMEs)

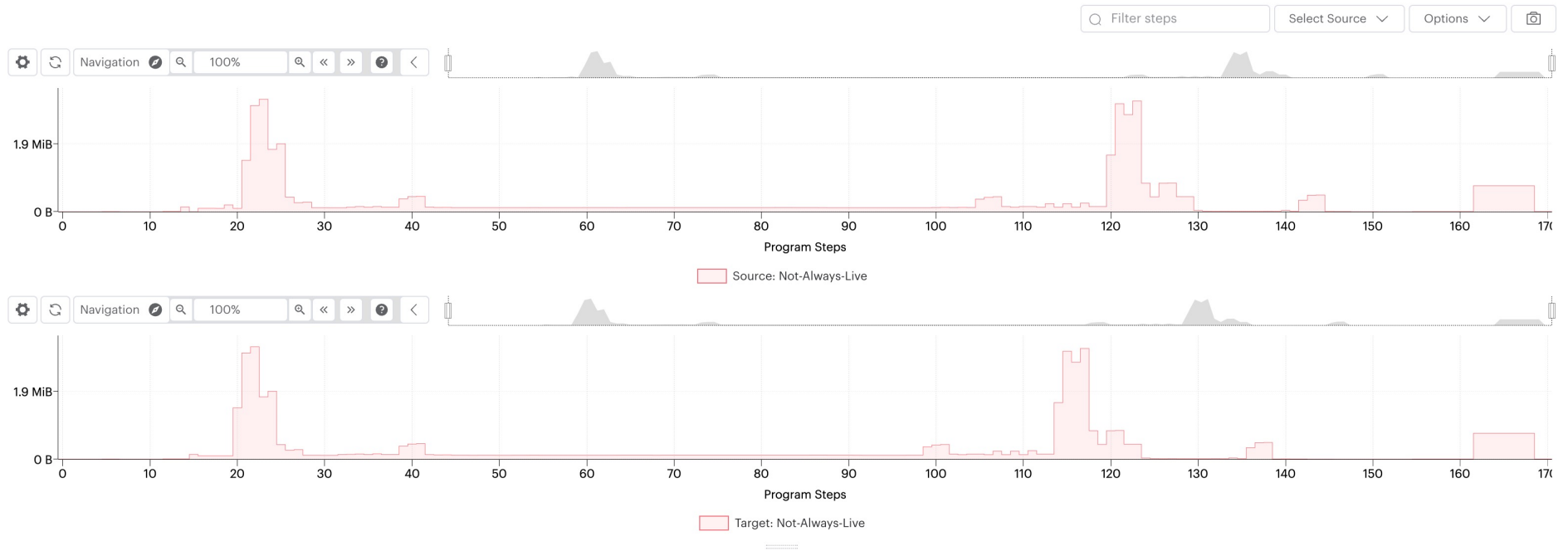
[tutorials/tensorflow2/infeed\\_outfeed](#)

[tutorials/pytorch/pipelining](#)

[tutorials/popvision/system\\_analyser\\_instrumentation](#)



- Summary
- Memory Report
- Liveness Report
- Program Tree
- Operations Summary
- Execution Trace
  
- Swap Reports
- Reload Report
- Close Report
- Documentation
- Minimise



**Selected Program Step**

Name (source): None selected. Click on the graph to select a source program step  
 Name (target): None selected. Click on the graph to select a target program step

[Always-Live Variables](#)   [Not-Always-Live Variables](#)   [Vertices](#)   [Cycle Estimates](#)

Filter variables

Variable	Source: All Tiles	Target: All Tiles	Diff: All Tiles
<b>Total</b>	<b>11.3 MiB</b>	<b>11.2 MiB</b>	<b>-113 KiB</b>
vertexCode	6.0 MiB	6.0 MiB	-16 B (-0.00%)
controlCode	1.5 MiB	1.5 MiB	-37.7 KiB (-2.41%)
internalExchangeCode	1.2 MiB	1.1 MiB	-47.6 KiB (-3.99%)
hostExchangeCode	595 KiB	556 KiB	-39.2 KiB (-6.58%)
stack	540 KiB	529 KiB	-11.0 KiB (-2.04%)
vertexInstanceState	527 KiB	527 KiB	20 B (0.00%)
while/sequential/dense/MatMul/dot.35/rhs	392 KiB	—	-392 KiB (-100.00%)
sequential/dense/MatMul/dot/rhs	—	392 KiB	392 KiB (100.00%)
instrumentationResults	180 KiB	180 KiB	4 B (0.00%)
hostExchangePacketHeader	114 KiB	137 KiB	23.7 KiB (17.25%)



Summary

Insights

Memory Report

Liveness Report

Program Tree

Operations Summary

Operations Graph

Execution Trace

Reload Report

Close Report

Documentation

Minimise

Execution View - Flame

Navigation

260775%

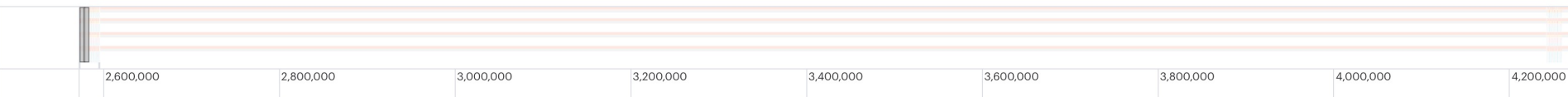
All IPUs

Search

Options

Options

Options



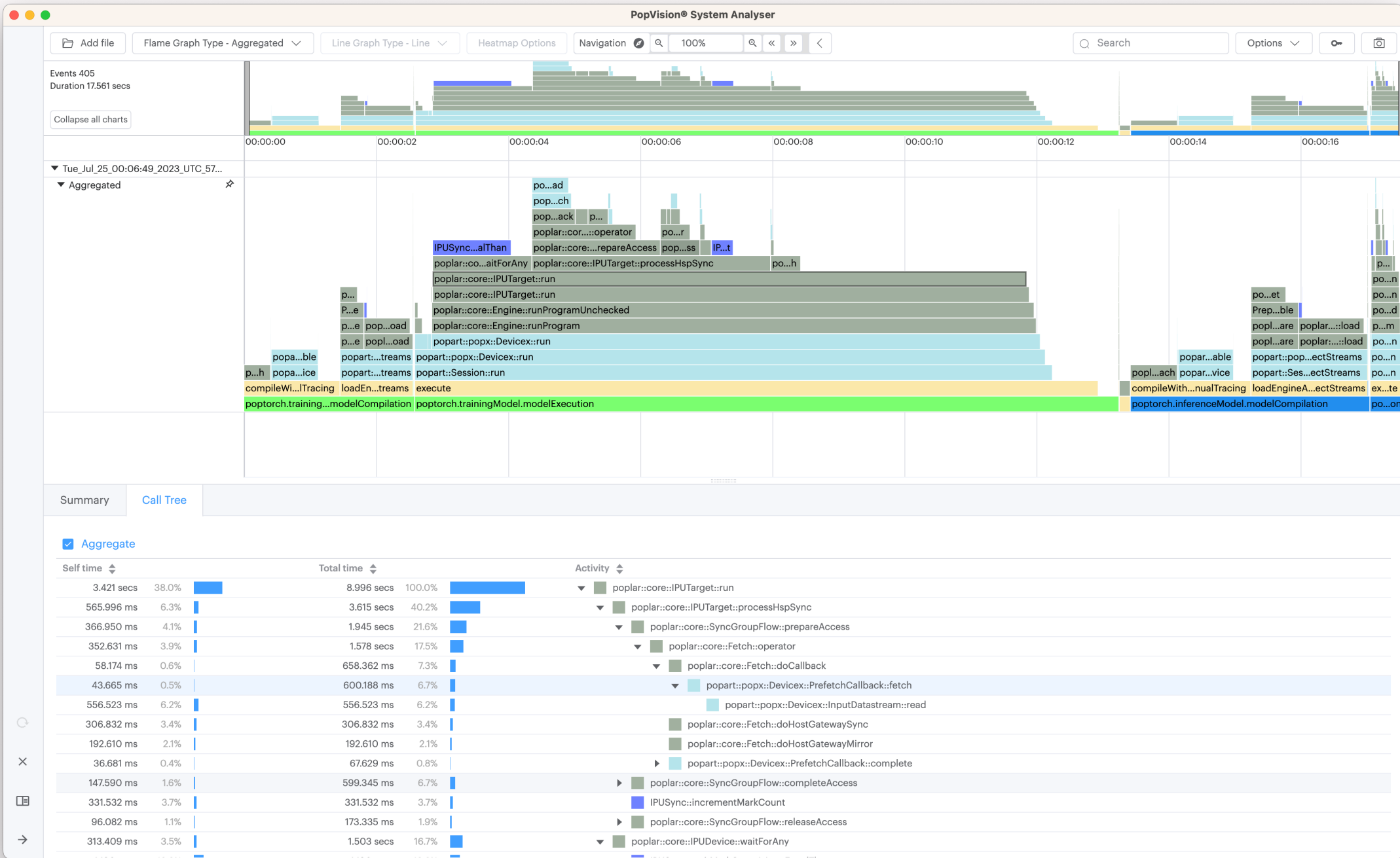
Summary

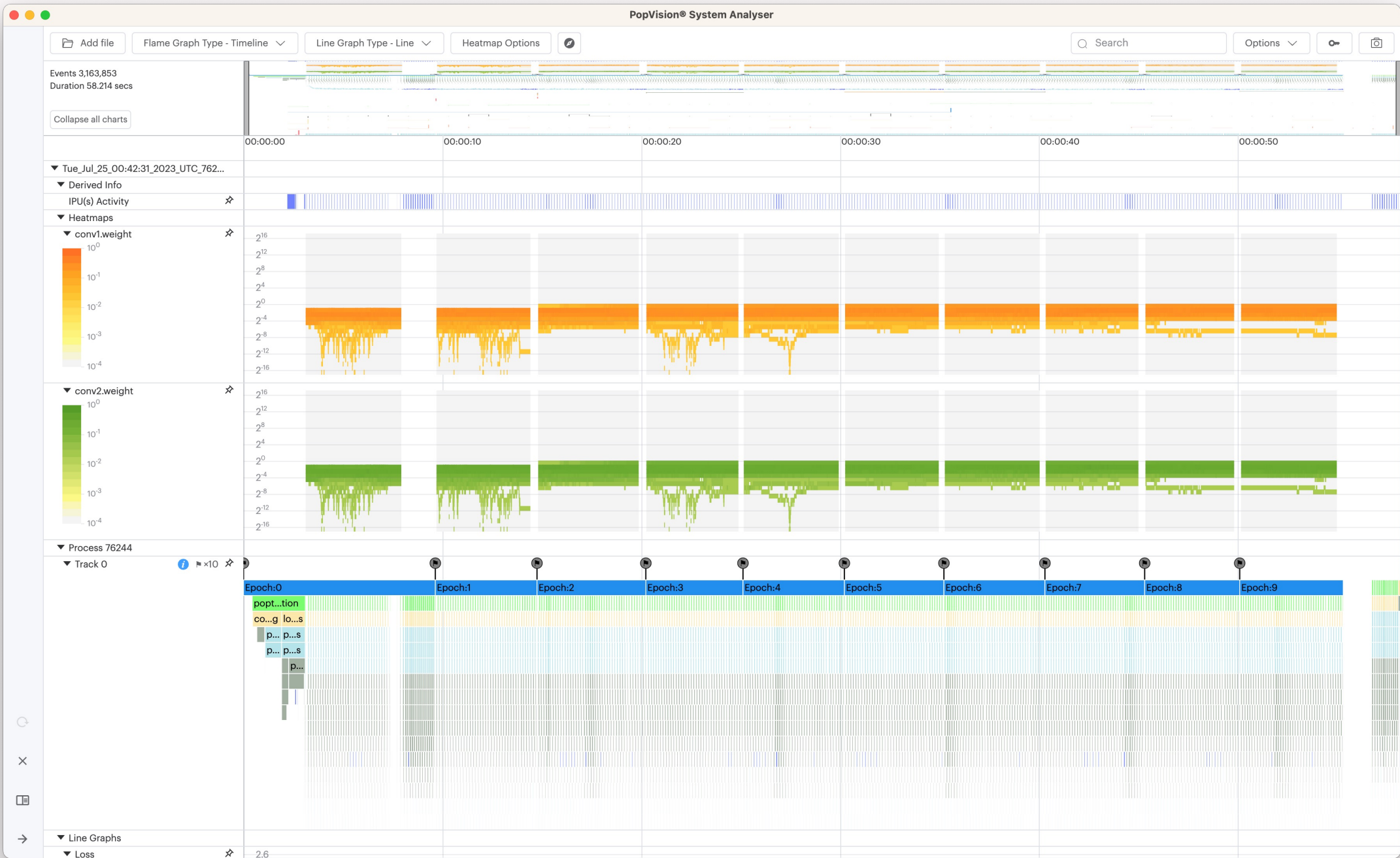
Details

Total Cycles 1,700,636 Total FLOPs 2,227,345,013

IPU	Proportions	Tile Utilisation	Stream Copy (Rx/Tx)	Global Exchange (Rx/Tx)
0		56.02%	1.1 MiB / 172 KiB	4.1 MiB / 4.2 MiB
1		54.73%	0 B / 0 B	5.0 MiB / 5.3 MiB
2		26.01%	0 B / 0 B	1.8 MiB / 1.4 MiB
3		18.01%	768 B / 320 B	640 KiB / 640 KiB











Request an IPU allocation for your research today!  
[allocations.access-ci.org/prepare-requests-overview](https://allocations.access-ci.org/prepare-requests-overview)

## Prepare requests: overview

We've organized your options for requesting access to advanced research computing resources into a set of opportunities designed to support needs ranging from new or entry-level exploration to the largest-scale computational experiments. We welcome you to find the opportunity that aligns with your best estimate of your resource needs. And don't worry about starting too small. As you clarify your needs, you can upgrade to a larger-scale opportunity when you're ready.

There are four opportunities, which are described below and **compared side-by-side in this table**:

- Explore ACCESS
- Discover ACCESS

**Prepare requests: overview**

- Explore ACCESS
- Discover ACCESS
- Accelerate ACCESS
- Maximize ACCESS
- Comparison Table





# THANK YOU

Alexander Tsyplikhin  
alex@graphcore.ai

Request an IPU allocation for your research today!  
[allocations.access-ci.org/prepare-requests-overview](https://allocations.access-ci.org/prepare-requests-overview)