



Programming IPUs for Scientific Computing

Alexander Tsyplikhin, Field AI Engineering

Tutorial with Argonne National Laboratory

IPU – Architectured For AI

Massive parallelism with ultrafast memory access

Parallelism

CPU

Designed for scalar processes

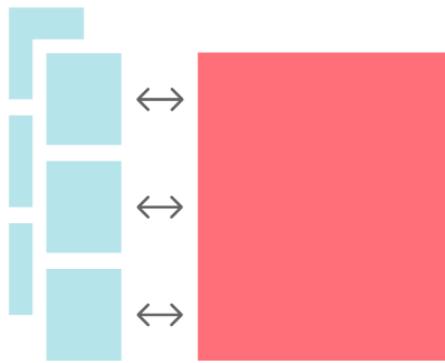
GPU

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

IPU

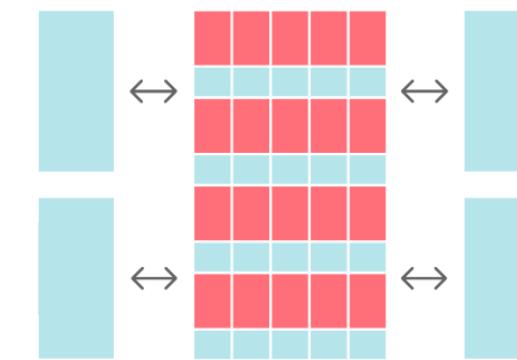
Massively parallel MIMD. Designed for fine-grained, high-performance computing

Processors 
Memory 

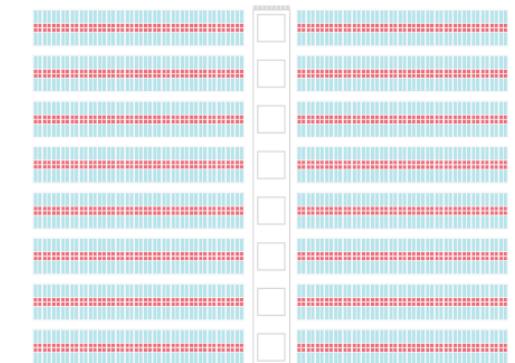


Memory Access

Off-chip memory

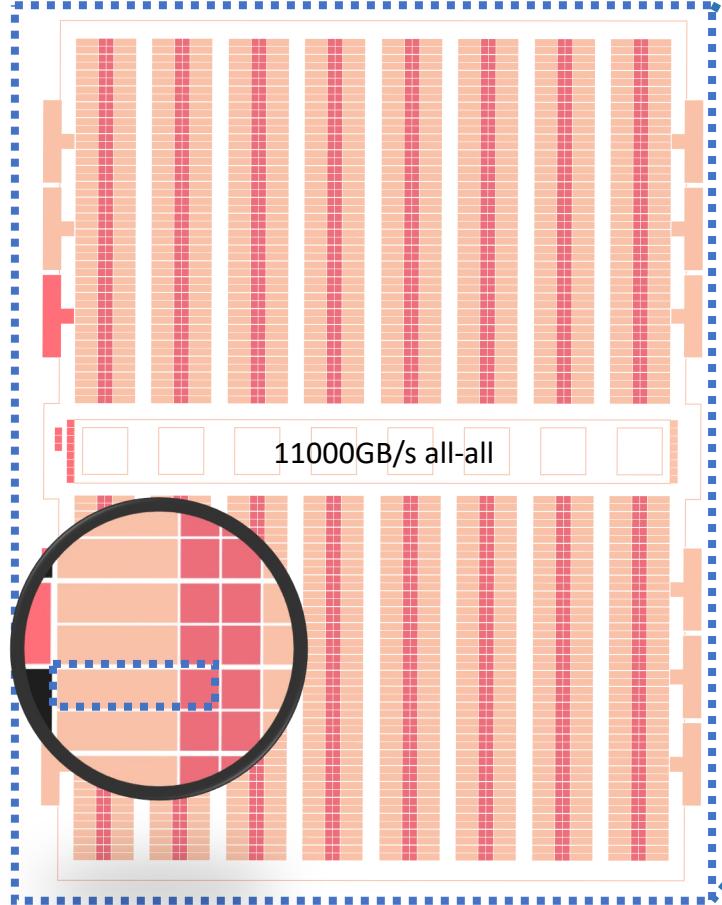


Model and data spread across off-chip and small on-chip cache, and shared memory



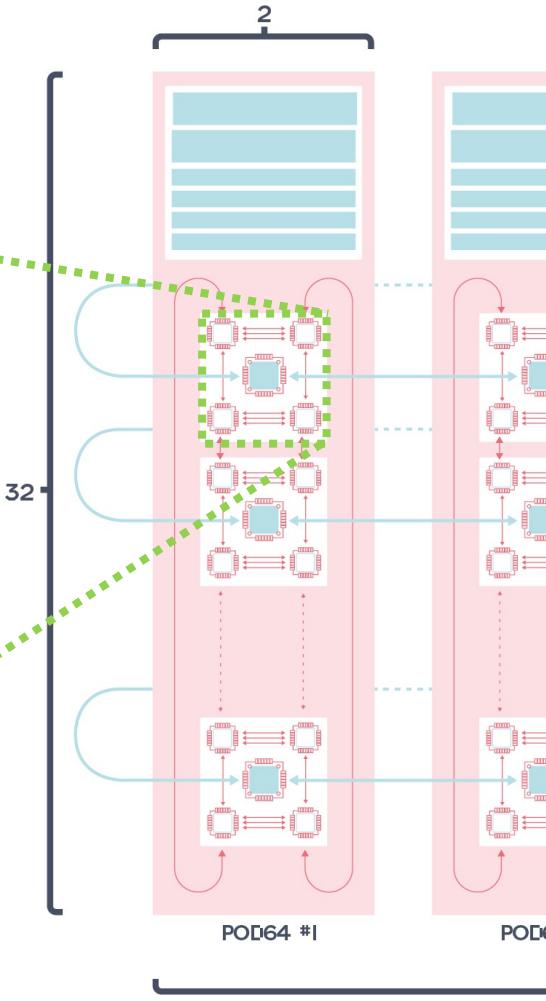
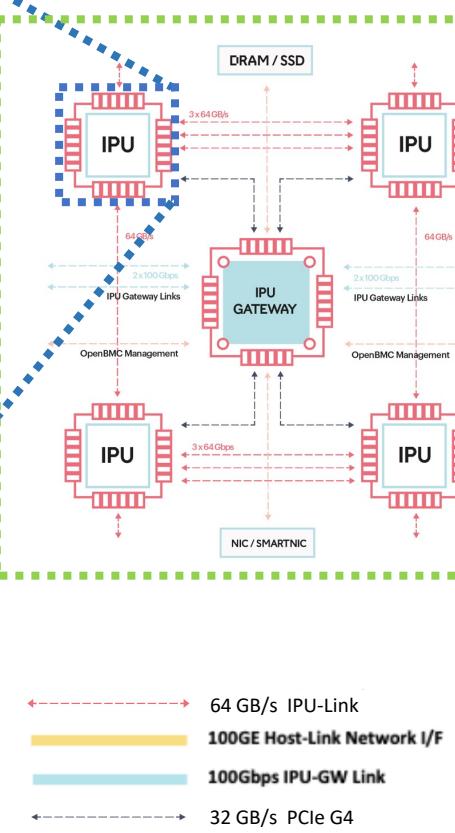
Model and data tightly coupled, and large locally distributed SRAM

BOW IPU
350TFLOPS(F16)
900MB SRAM
1472 IPU-Tiles
8832 independent instruction streams

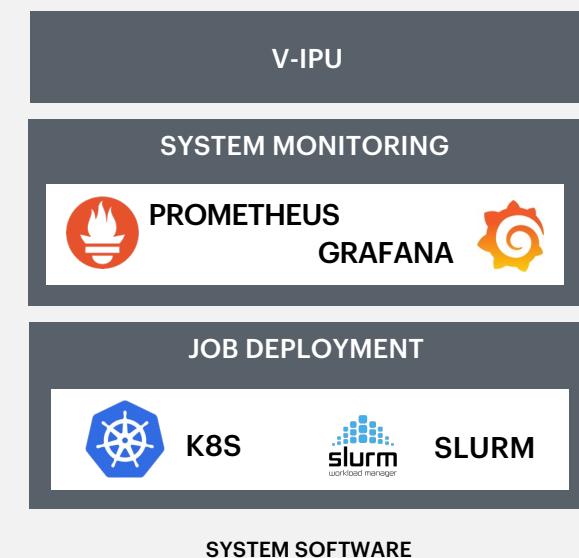
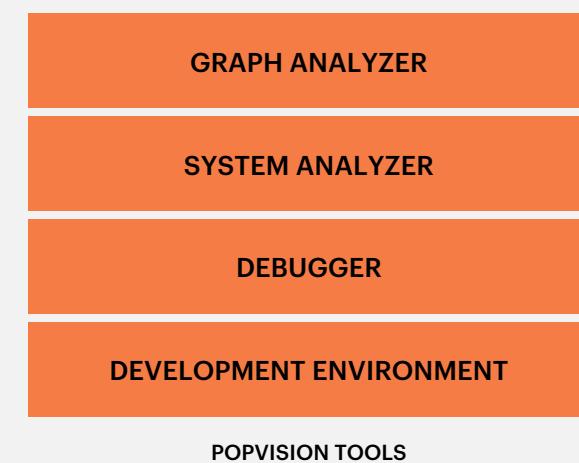
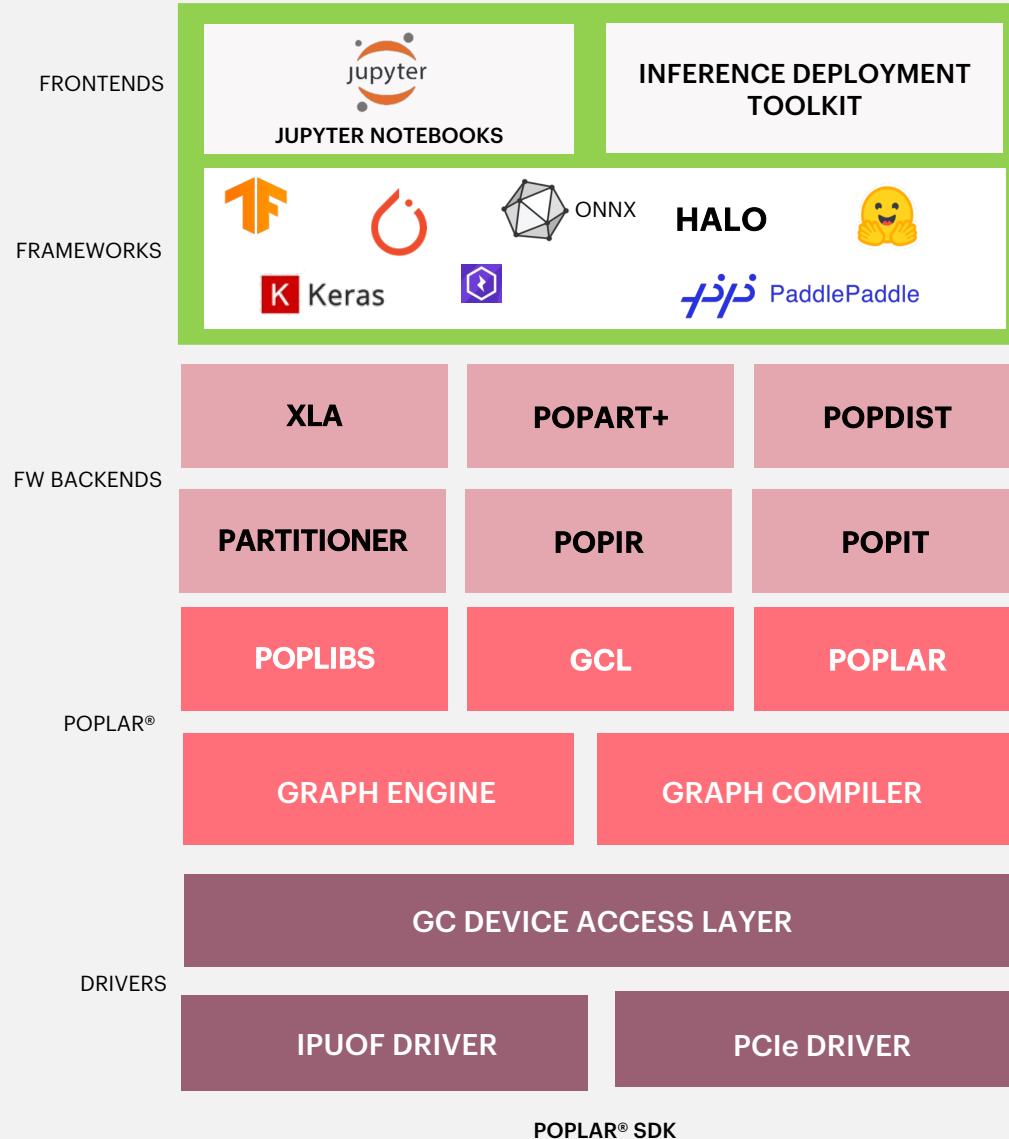
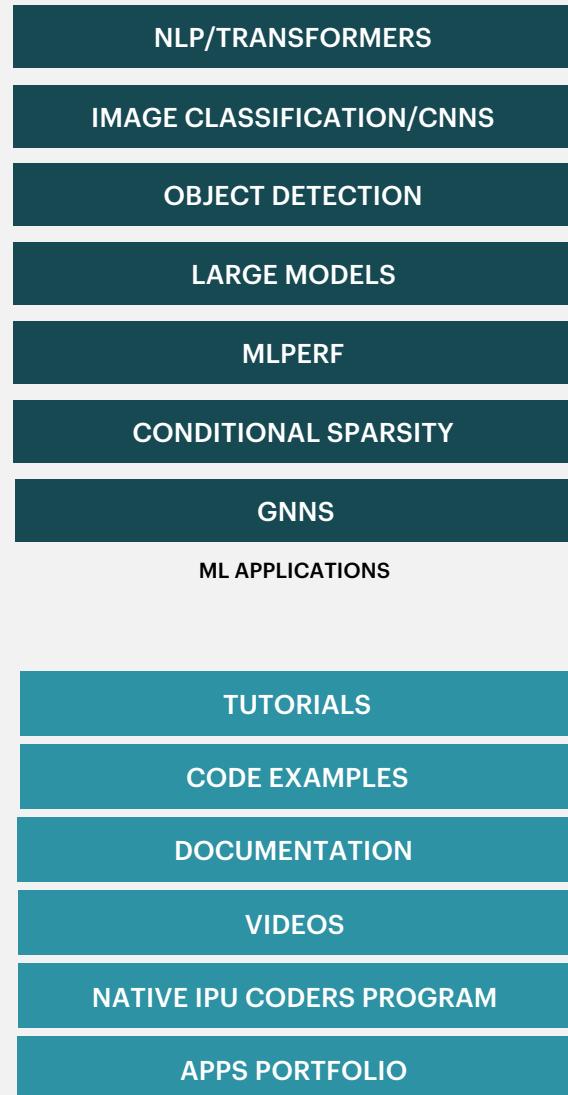


BOW-2000 IPU-Machine

4x IPU
256 GB DRAM



GRAPHCORE SOFTWARE





GRAPHCORE SOFTWARE ECOSYSTEM

WORLD CLASS DEVELOPER RESOURCES FOR IPU USERS

GRAPHCORE

Home About Products Industries Developer Blog Careers Get Started

[Quickstart Docs](#) → [Try on Paperspace](#) →

DOCUMENTATION
Explore our software documentation, user guides and technical notes.
[See more](#) →

TUTORIALS + CODE EXAMPLES
Hands-on code tutorials, simple application and feature examples.
[See more](#) →

GRAPHCORE GITHUB
Access our open source libraries, APIs, applications and code examples.
[See more](#) →

DOCKER HUB
Access a selection of Poplar SDK container images via Docker Hub.
[See more](#) →

MODEL GARDEN
Access a repository of deployable ML applications on the IPU.
[See more](#) →

HOW-TO VIDEOS
Watch practical how-to videos and demos by Graphcore engineers.
[See more](#) →

WEBINARS
Register for upcoming Graphcore webinars or watch on-demand.
[See more](#) →

RESEARCH PAPERS
Read publications from Graphcore's Research team and IPU innovators.
[See more](#) →

POPVISION™ TOOLS
Download PopVision to analyse IPU performance and utilisation.
[See more](#) →

WWW.GRAPHCORE.AI/DEVELOPER

GRAPHCORE

Graphcore Documents Version: Latest

Search docs

- Getting Started
- Software Documents
- Hardware Documents
- Technical Notes and White Papers
- Examples and Tutorials
- Document Updates
- Alphabetical List of All Documents
- Graphcore License Agreements

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

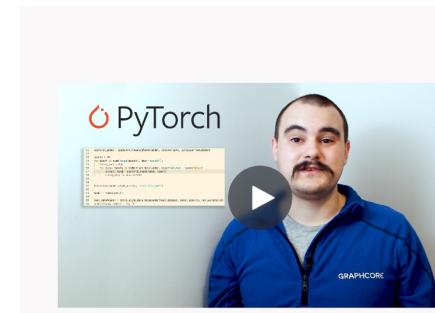
Technical notes and white papers on Graphcore technology

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 Bogdikiewicz 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.

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

[Get the Code](#) →

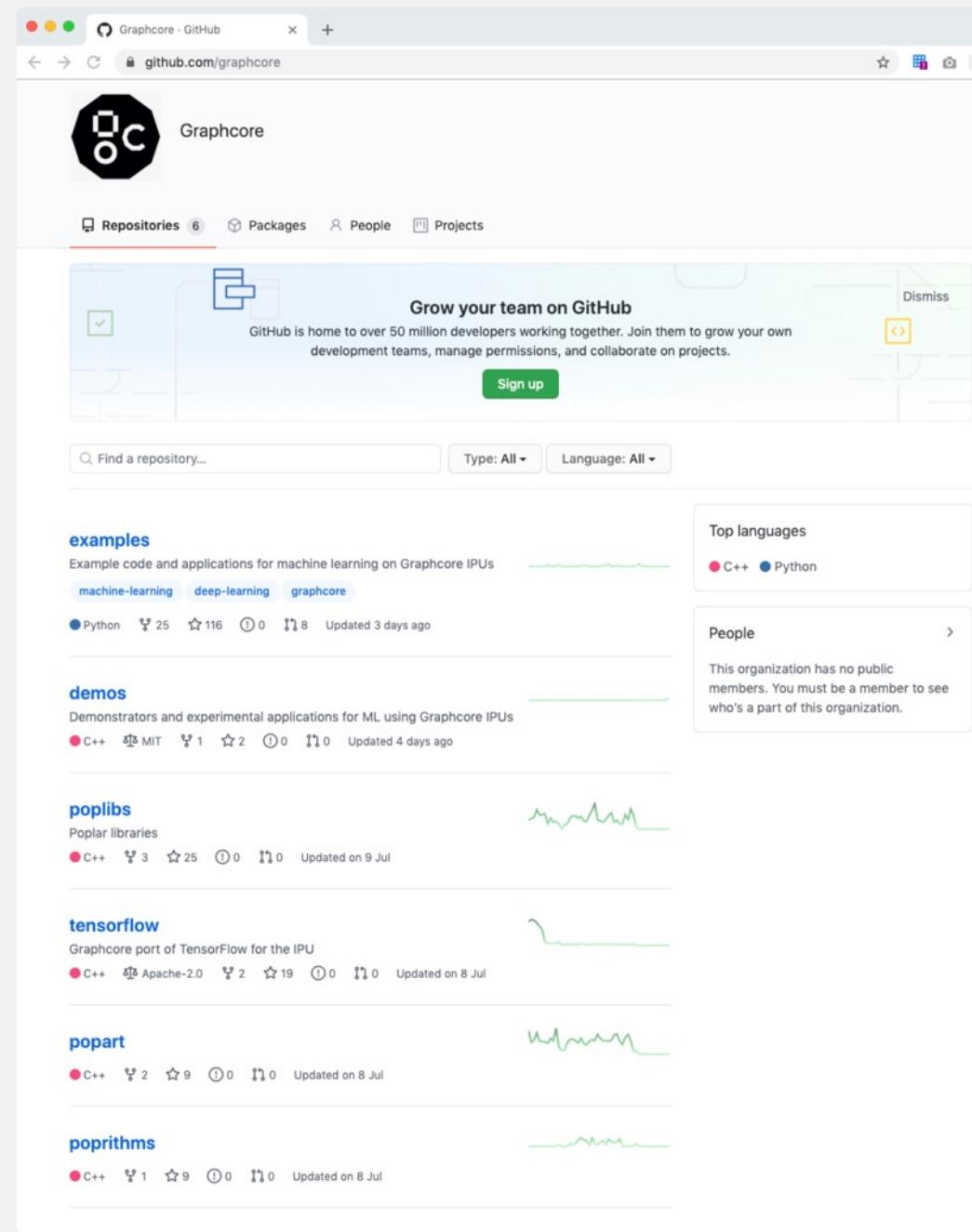
[Read the Guide](#) →



OPEN SOURCE

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



ENHANCED MODEL GARDEN

The screenshot shows the Enhanced Model Garden interface. On the left, there's a sidebar with navigation links like 'Resources > Model Garden' and sections for 'MODEL GARDEN' (with a back arrow), 'LIBRARY' (with a search bar and filters for Type, Framework, and Category), and a detailed view of a model card.

FEATURED MODELS

GPS++ INFERENCE	DISTRIBUTED KGE - TRANSE (256) TRAINING	GPT-J 6B FINE-TUNING
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.	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.	GPT-J 6B fine-tuned using the GLUE MNLI dataset leveraging the Hugging Face Transformer library.
View Repository	View Repository	View Repository

LIBRARY

Search:

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

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 MNLI dataset leveraging the Hugging Face Transformer library. View Repository
DISTILBERT TRAINING DistilBERT is a small, fast, cheap and light View Repository	MAE TRAINING Implementation of MAE computer vision model in View Repository	FROZEN IN TIME TRAINING Implementation of Frozen in Time on the IPU in View Repository

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

NEW FILTER/SEARCH CAPABILITY

DIRECT ACCESS TO GITHUB

PAPERSPACE NOTEBOOK LINKS



PyTorch

GPU

```
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 (de
    parser.add_argument('--test-batch-size', type=int, default=8, help='batch size for testing
    parser.add_argument('--epochs', type=int, default=10, help='number of epochs to train (de
    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, tra
        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,
        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)))
```

IPU

```
_, ind = torch.argmax(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) * 100.0
return accuracy

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 (de
    parser.add_argument('--test-batch-size', type=int, default=8, help='batch size for testing
    parser.add_argument('--epochs', type=int, default=10, help='number of epochs to train (de
    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
    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, tra
        batch_size=args.batch_size, shuffle=True, drop_last=True)
    test_data = poptorch.DataLoader(opts,
        torchvision.datasets.MNIST('mnist_data/', train=False, download=True, tra

    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)))
```

Keras

```
gpu_cnn_keras.py ↔ ipu_cnn_keras.py tf_keras

    as tf
    keras.layers import *

Keras

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)
```

GPU

IPU

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

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)
```

IPU

CODELET DEFINITIONS

The fields of the vertex specify its inputs, outputs and internal data.

```
class AdderVertex : public Vertex {  
public:  
    Input<float> x;  
    Input<float> y;  
    Output<float> z;  
    float bias;  
  
    bool compute() {  
        *z = x + y + bias;  
        return true;  
    }  
}
```

Each codelet is defined as a C++ class that inherits from the **Vertex** class.

The compute method specifies the vertex execution behaviour.

BULK SYNCHRONOUS PARALLEL (BSP)

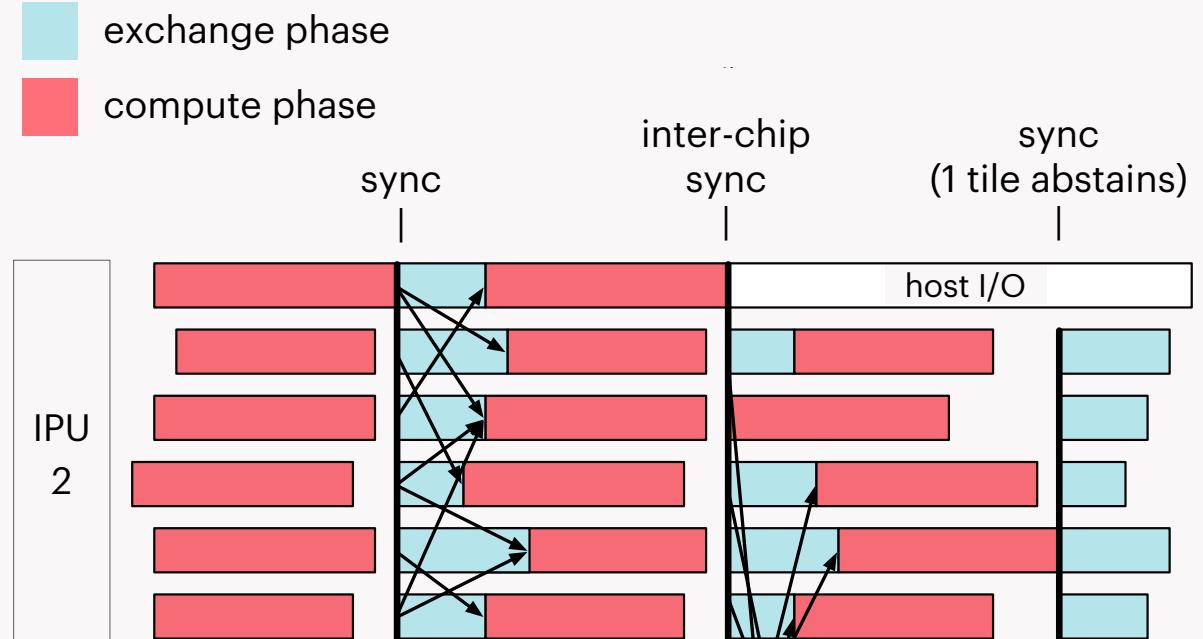
BSP software bridging model – massively parallel computing with no concurrency hazards

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



time →

POPVISION™ TOOLS

GRAPH ANALYSER

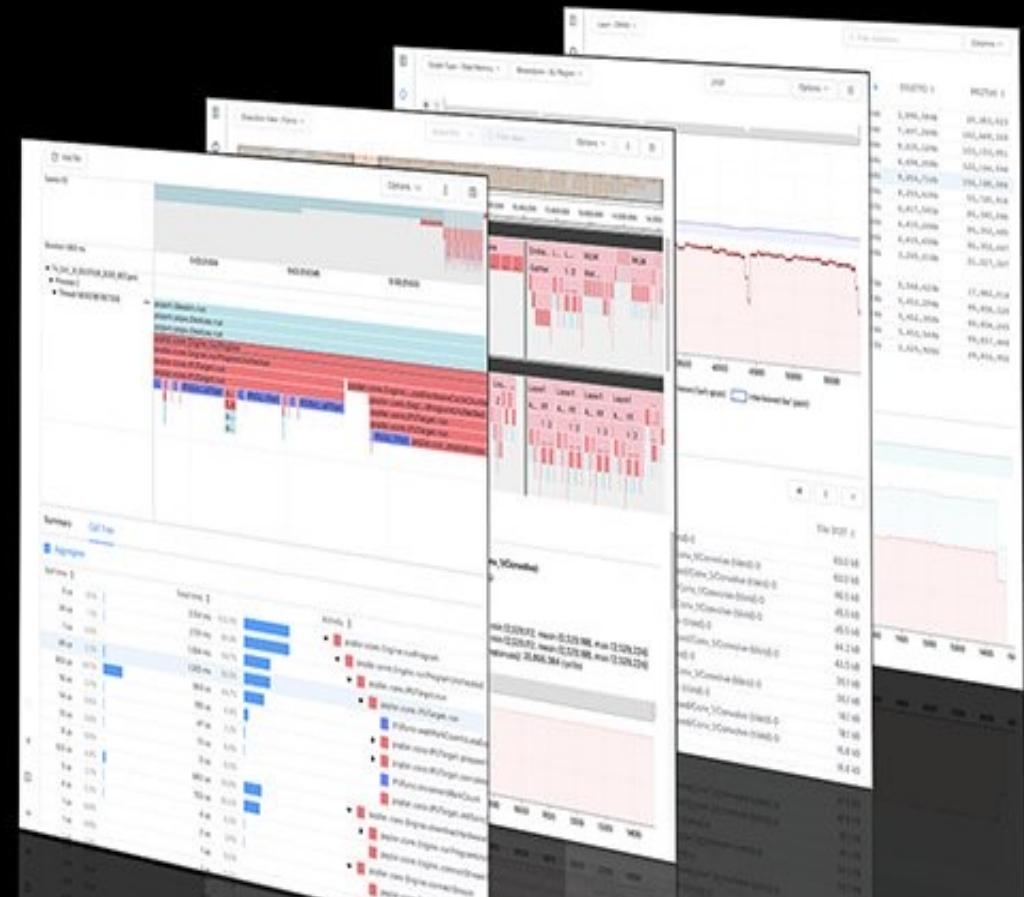
Useful for analysing and optimising the memory use and execution performance of ML models on the IPU

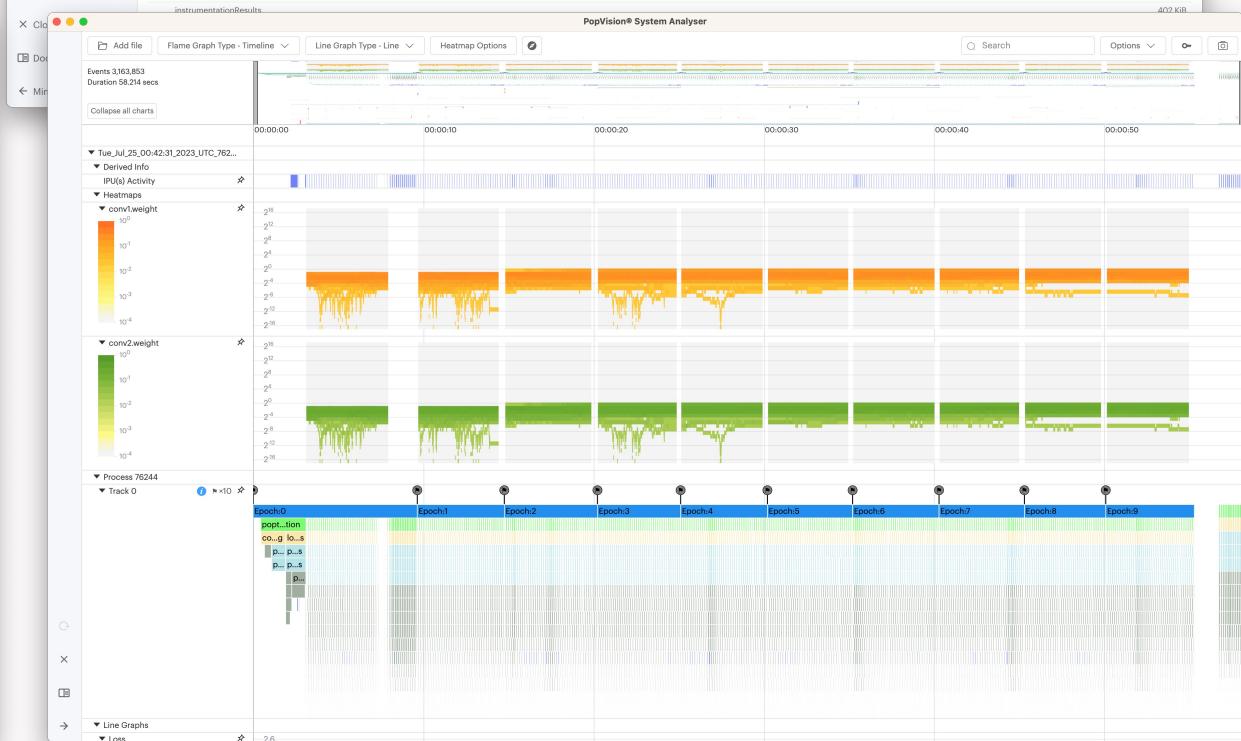
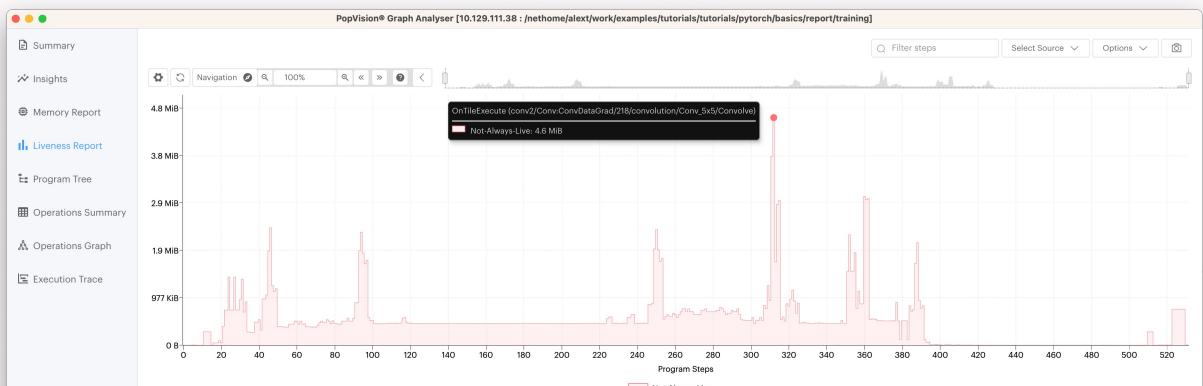
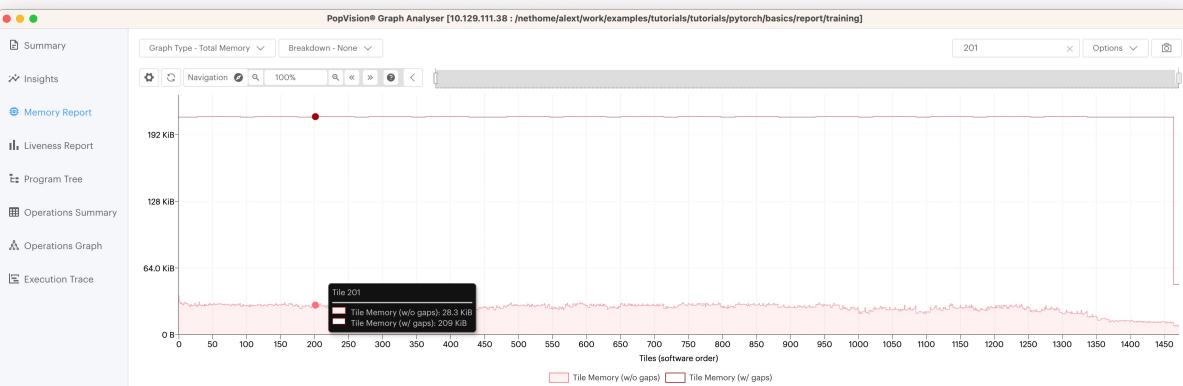
SYSTEM ANALYSER

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





TUTORIALS

PyTorch

github.com/graphcore/examples/tree/master/tutorials/tutorials/pytorch/basics

github.com/graphcore/examples/tree/master/tutorials/tutorials/pytorch/mixed_precision

github.com/graphcore/examples/tree/master/tutorials/tutorials/pytorch/efficient_data_loading

github.com/graphcore/examples/tree/master/tutorials/tutorials/pytorch/pipelining

TensorFlow/Keras

github.com/graphcore/examples/tree/master/tutorials/tutorials/tensorflow2/keras

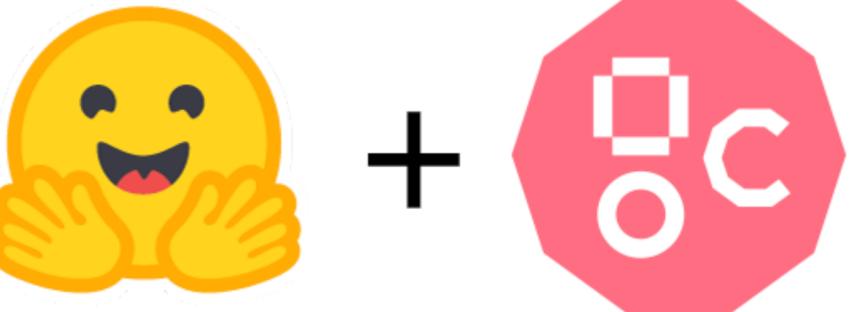
PopVision

[tutorials/tensorflow2/infeed_outfeed](https://github.com/graphcore/tutorials/tree/main/tutorials/tensorflow2/infeed_outfeed)

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

[tutorials/popvision/system_analyser_instrumentation](https://github.com/graphcore/tutorials/tree/main/tutorials/popvision/system_analyser_instrumentation)

README.md



Optimum Graphcore

😊 Optimum Graphcore is the interface between the 😊 Transformers library and [Graphcore IPUs](#). It provides a set of tools enabling model parallelization and loading on IPUs, training, fine-tuning and inference on all the tasks already supported by 😊 Transformers while being compatible with the 😊 Hub and every model available on it out of the box.

What is an Intelligence Processing Unit (IPU)?

Quote from the Hugging Face [blog post](#):

IPUs are the processors that power Graphcore's IPU-POD datacenter compute systems. This new type of processor is designed to support the very specific computational requirements of AI and machine learning. Characteristics such as fine-grained parallelism, low precision arithmetic, and the ability to handle sparsity have been built into our silicon.

Instead of adopting a SIMD/SIMT architecture like GPUs, Graphcore's IPU uses a massively parallel, MIMD architecture, with ultra-high bandwidth memory placed adjacent to the processor cores, right on the silicon die.

This design delivers high performance and new levels of efficiency, whether running today's most popular models, such as BERT and EfficientNet, or exploring next-generation AI applications.



Contributors 36



+ 25 contributors

Languages



- Python 52.1%
- Jupyter Notebook 47.8%
- Makefile 0.1%



☰ README.md

How to use Optimum Graphcore

To immediately use a model on a given input (text, image, audio, ...), we support the `pipeline` API:

```
->>> from transformers import pipeline
+>>> from optimum.graphcore import pipeline

# Allocate a pipeline for sentiment-analysis
->>> classifier = pipeline('sentiment-analysis', model="distilbert-base-uncased-finetuned-sst-2-english")
+>>> classifier = pipeline('sentiment-analysis', model="distilbert-base-uncased-finetuned-sst-2-english")
>>> classifier('We are very happy to introduce pipeline to the transformers repository.')
[{'label': 'POSITIVE', 'score': 0.9996947050094604}]
```

It is also super easy to use the `Trainer` API:

```
-from transformers import Trainer, TrainingArguments
+from optimum.graphcore import IPUConfig, IPUTrainer, IPUTrainingArguments

-training_args = TrainingArguments(
+training_args = IPUTrainingArguments(
    per_device_train_batch_size=4,
    learning_rate=1e-4,
+    # Any IPUConfig on the Hub or stored locally
+    ipu_config_name="Graphcore/bert-base-ipu",
+)
+
+## Loading the IPUConfig needed by the IPUTrainer to compile and train the model on IPUs
+ipu_config = IPUConfig.from_pretrained(
+    training_args.ipu_config_name,
)

# Initialize our Trainer
-trainer = Trainer(
+trainer = IPUTrainer(
    model=model,
+    ipu_config=ipu_config,
    args=training_args,
```



☰ README.md

Supported models

The following model architectures and tasks are currently supported by 😊 Optimum Graphcore:

	Pre-Training	Masked LM	Causal LM	Seq2Seq LM (Summarization, Translation, etc)	Sequence Classification	Token Classification	Ques Answer
BART	✓		✗	✓	✓		✗
BERT	✓	✓	✗		✓	✓	✓
ConvNeXt	✓						
DeBERTa	✓	✓			✓	✓	✓
DistilBERT	✗	✓			✓	✓	✓
GPT-2	✓		✓		✓	✓	
GroupBERT	✓	✓	✗		✓	✓	✓
HuBERT	✗				✓		
LXMERT	✗						✓
RoBERTa	✓	✓	✗		✓	✓	✓
T5	✓			✓			
ViT	✗						
Wav2Vec2	✓						
Whisper	✗			✓			

If you find any issue while using those, please open an issue or a pull request.



huggingface / optimum-graphcore

Type ⌘ to search

Code Issues 9 Pull requests 8 Actions Projects 1 Security Insights

Code

main + ⌂ Go to file t

.github docs examples audio-classification image-classification language-modeling multiple-choice question-answering speech-pretraining speech-recognition summarization text-classification token-classification translation README.md notebooks images packed_bert

optimum-graphcore / examples / Add file ⌂

jimypbr Update examples requirements for sdk3.3 (#434) ✓ 5c4a5c6 · last week History

Name	Last commit message	Last commit date
..		
audio-classification	Update examples requirements for sdk3.3 (#434)	last week
image-classification	Update examples requirements for sdk3.3 (#434)	last week
language-modeling	Bump transformers to 4.29.2 (#389)	last month
multiple-choice	Bump transformers to 4.29.2 (#389)	last month
question-answering	Bump transformers to 4.29.2 (#389)	last month
speech-pretraining	Update examples requirements for sdk3.3 (#434)	last week
speech-recognition	Update examples requirements for sdk3.3 (#434)	last week
summarization	Bump transformers to 4.29.2 (#389)	last month
text-classification	Bump transformers to 4.29.2 (#389)	last month
token-classification	Bump transformers to 4.29.2 (#389)	last month
translation	Bump transformers to 4.29.2 (#389)	last month
README.md	Update README.md	last year

Files

examples / nlp / gpt2 / pytorch / inference_gpt2.py

Code Blame 205 lines (181 loc) · 8.03 KB

Raw ⌂ ⌄ ⌅ ⌆

master + ⌂

Go to file t

gpt2/pytorch config config.json config_large.json config_medium.json config_test.json config_xl.json custom_ops data model run tasks tests tokenizer LICENSE Makefile README.md arguments.py benchmarks.yml inference_gpt2.py ipu_options.py requirements.txt text_generate_gpt2.py tools.py train_gpt2.py gpt3_175B/popxl gpt3_2.7B/popxl

```
62
63     class GPT2Wrapper(nn.Module):
64         def __init__(self, args, model_config):
65             super().__init__()
66             self.args = args
67             if args.checkpoint_input_dir: # load pretrained model checkpoint
68                 self.model = GPT2LMHeadModel.from_pretrained(args.checkpoint_input_dir)
69             else: # init model
70                 self.config = model_config
71                 self.model = GPT2LMHeadModel(config=self.config)
72
73             for layer in self.model.transformer.h:
74                 gpt2_attn = OptimizedGPT2Attention(self.model.config, layer_idx=layer.attn.layer_idx)
75                 gpt2_attn.load_state_dict(layer.attn.state_dict())
76                 layer.attn = gpt2_attn
77
78             if args.embedding_serialization_factor > 1:
79                 serialized_lmhead = SerializedLinear(
80                     self.model.config.n_embd,
81                     self.model.config.vocab_size,
82                     args.embedding_serialization_factor,
83                     bias=False,
84                     mode=poptorch.MatMulSerializationMode.OutputChannels,
85                 )
86                 serialized_lmhead.load_state_dict(self.model.lm_head.state_dict())
87                 self.model.lm_head = serialized_lmhead
88                 self.model.tie_weights()
89
90             logger("----- Device Allocation -----")
91             logger("Embedding --> IPU 0")
92             self.model.transformer.wte = poptorch.BeginBlock(self.model.transformer.wte, "wte", ipu_id=0)
93             self.model.transformer.wpe = poptorch.BeginBlock(self.model.transformer.wpe, "wpe", ipu_id=1)
94             outline_attribute(self.model.transformer.ln_f, "LayerNorm")
95
96             layer_ipu = _get_layer_ipu(args.layers_per_ipu)
97             for index, layer in enumerate(self.model.transformer.h):
98                 ipu = layer_ipu[index]
99                 self.model.transformer.h[index] = poptorch.BeginBlock(layer, f"Encoder{index}", ipu_id=ipu)
100                logger(f"Layer {index}<2 --> IPU {ipu}")
101
102                logger(f"LM_head --> IPU 0")
103                self.model.lm_head = poptorch.BeginBlock(self.model.lm_head, ipu_id=0)
104
105        def forward(self, input_ids):
106            transformer_outputs = self.model.transformer(input_ids=input_ids)
107            hidden_states = transformer_outputs[0]
108            lm_logits = self.model.lm_head(hidden_states)
```

GPT2 training and inference

} Optimized attention

} Serialized embedding

} Pipeline-parallel

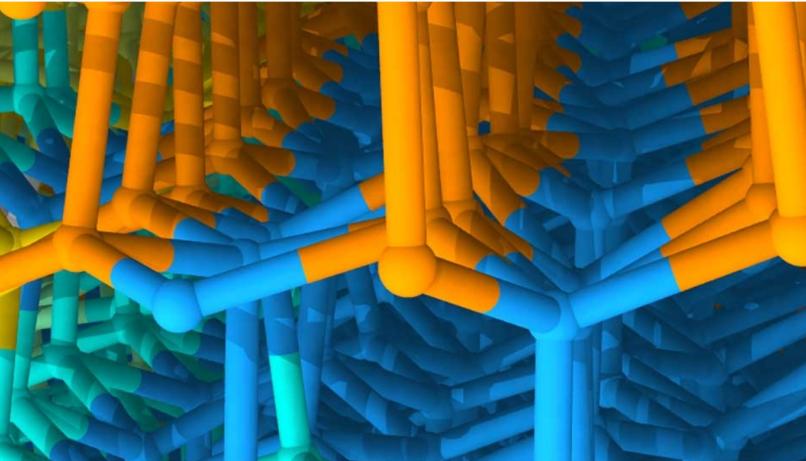
APPLY AND JOIN TODAY

Argonne Leadership Computing Facility

HOME / SCIENCE

Director's Discretionary Allocation Program

The ALCF Director's Discretionary program provides "start up" awards to researchers working to achieve computational readiness for a major allocation award.



Molecular dynamics simulations based on machine learning help scientists learn about the movement of the boundary between ice grains (yellow/green/cyan) and the stacking disorder that occurs when hexagonal (orange) and cubic (blue) pieces of ice freeze together. Image: Henry Chan and Subramanian Sankaranarayanan, Argonne National Laboratory

Apply at alcf.anl.gov/science/directors-discretionary-allocation-program

general ▾

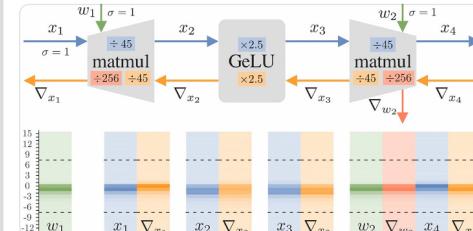
charlieb 6:05 AM
Pleased to share with you all some new work from the Graphcore research team! 🎉

Our paper *Unit Scaling* introduces a new method for low-precision number formats, making FP16 We've managed to train BERT in these formats for the first time without loss scaling.

- You can find our blog post here: <https://www.graphcore.ai/posts/simple-fp16-and-fp8-training>
- Paperspace notebook (try it yourself!): <https://ipu.dev/qXfm2a>
- Arxiv paper: <https://arxiv.org/abs/2303.11257>

(& we were also featured on Davis Blalock's popular [ML newsletter](#) this week) (edited)

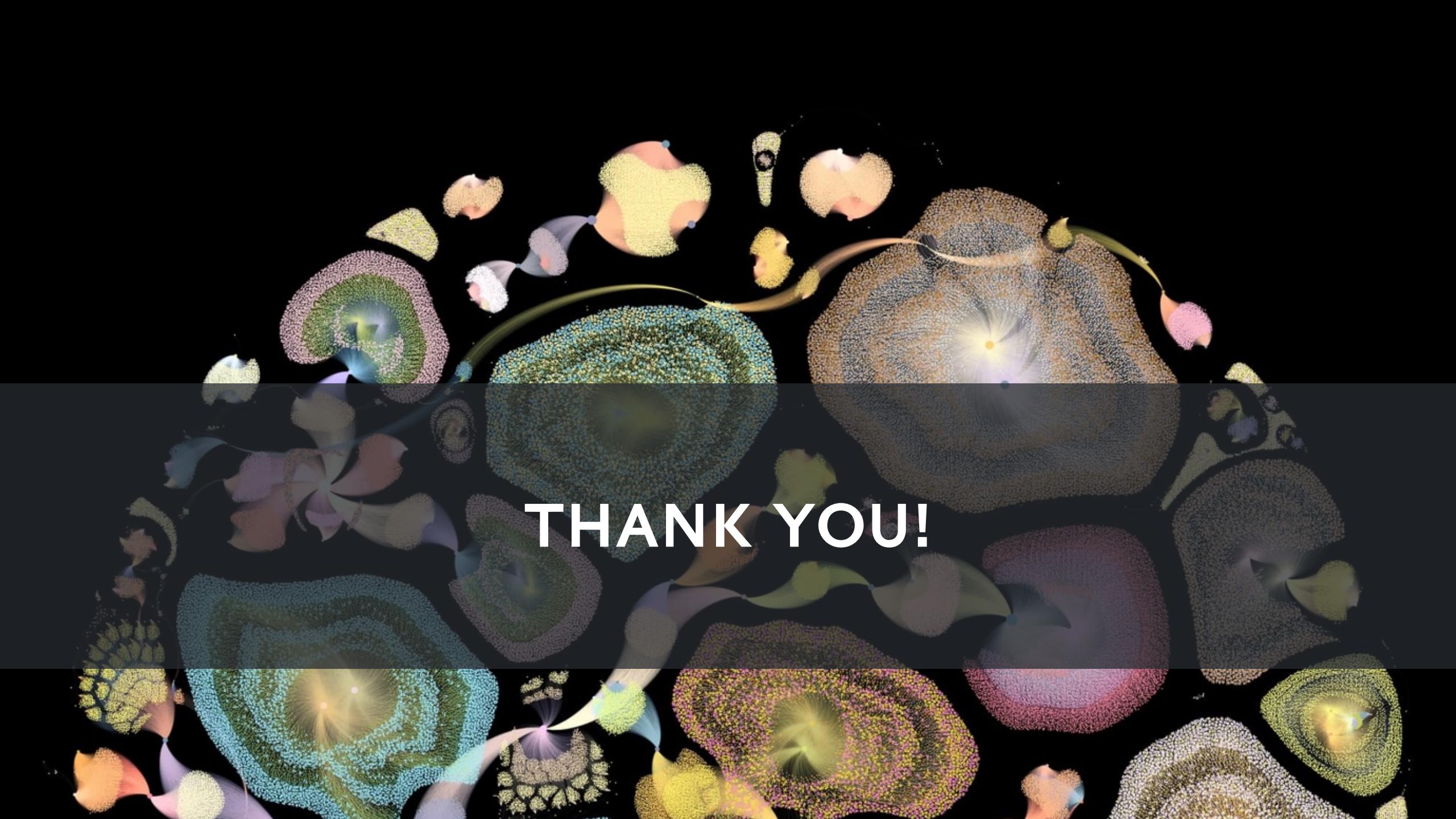
graphcore.ai
Simple FP16 and FP8 training with unit scaling
Unit Scaling is a new low-precision machine learning method able to train language models in FP16 and FP8 without loss scaling. (69 kB) ▾



arXiv.org
Unit Scaling: Out-of-the-Box Low-Precision Training
We present unit scaling, a paradigm for designing deep learning models that simplifies the use of low-precision number formats. Training in FP16 or the recently proposed FP8 formats offers substantial efficiency gains, but can lack sufficient range for out-of-the-box training. Unit scaling addresses this by introducing a principled approach to model numerics: seeking unit variance of Show more

8 8 7 1 1 1 1

Join at graphcore.ai/join-community

The background of the image is a dark, solid black. Overlaid on this are numerous small, spherical particles that appear to be made of a granular or textured material. These spheres are illuminated from within, creating a bright glow that varies in intensity. The colors of the spheres are diverse, including shades of yellow, orange, red, green, blue, and purple. Some spheres are larger and more prominent, while others are smaller and scattered throughout the scene. The overall effect is reminiscent of a microscopic view of a complex biological or physical system.

THANK YOU!