

DEEP DIVE ON GRAPH NEURAL NETWORKS AND LARGE LANGUAGE MODELS

June 12, 2024

Alexander Tsyplikhin

GRAPHCORE



AGENDA

GNNs

- Graphcore IPUs and PyTorch Geometric
- Case study: SchNet for molecular property prediction

LLMs


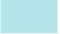
- HuggingFace Optimum
- PopART for GPT-3 175B

Q&A

IPU – Architected For AI

Massive parallelism with ultrafast memory access

Parallelism

Processors 
Memory 

Memory Access

CPU

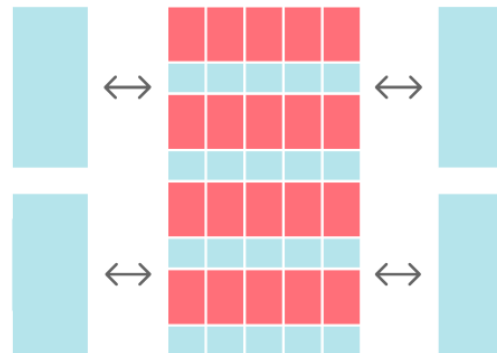
Designed for scalar processes



Off-chip memory

GPU

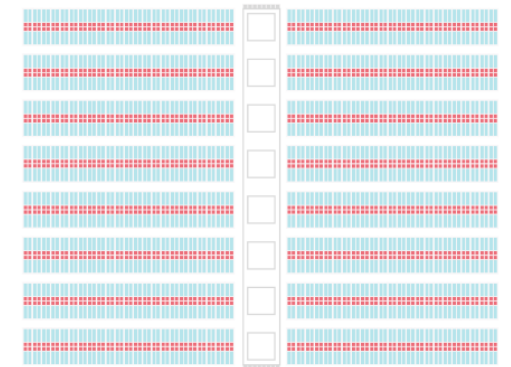
SIMD/SIMT architecture. Designed for large blocks of dense contiguous data



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

IPU

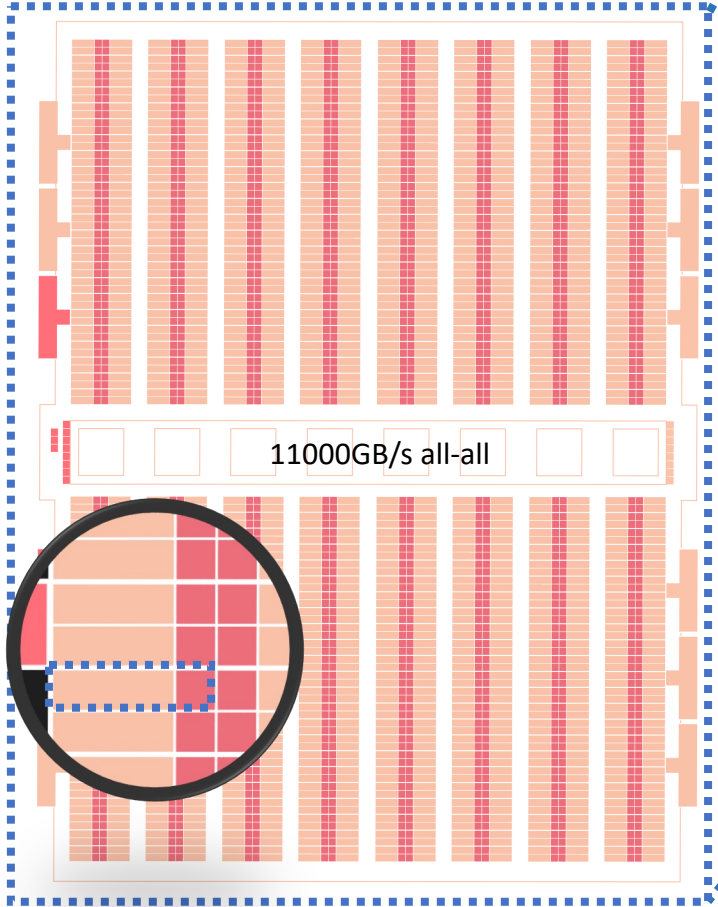
Massively parallel MIMD. Designed for fine-grained, high-performance computing



Model and data tightly coupled, and large locally distributed SRAM

BOW IPU

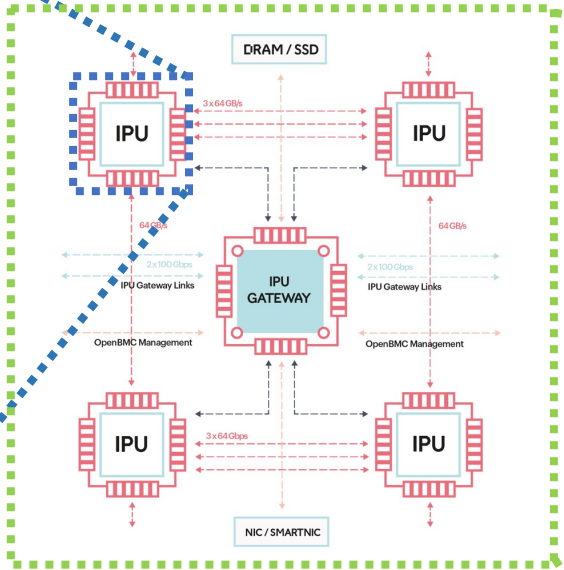
350TFLOPS(F16)
900MB SRAM
1472 IPU-Tiles
8832 independent instruction streams



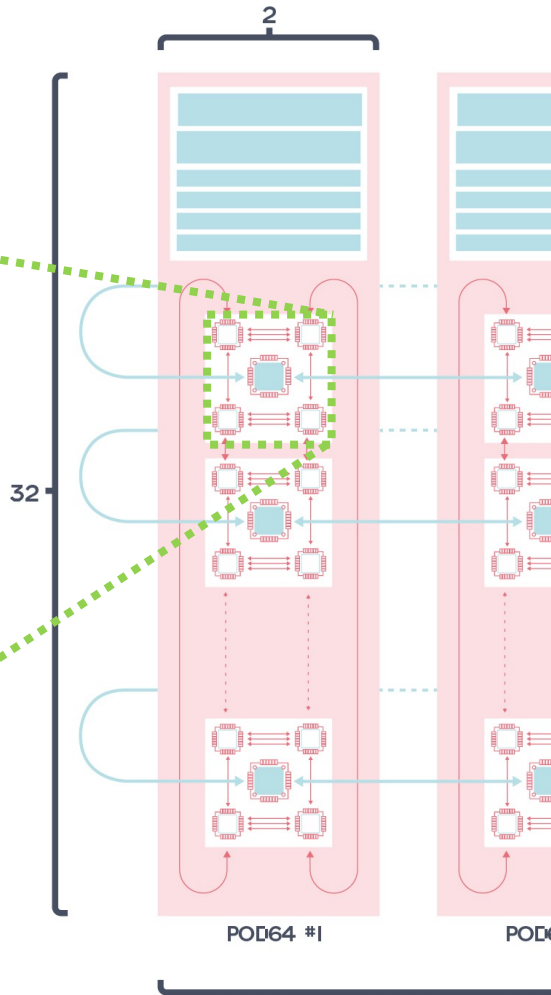
IPU-Tile
236GFLOPS(F16)
640KB SRAM
6 HW Threads
1.83GHz
128 F16Ops/Cycle

BOW-2000 IPU-Machine

4x IPU
256 GB DRAM



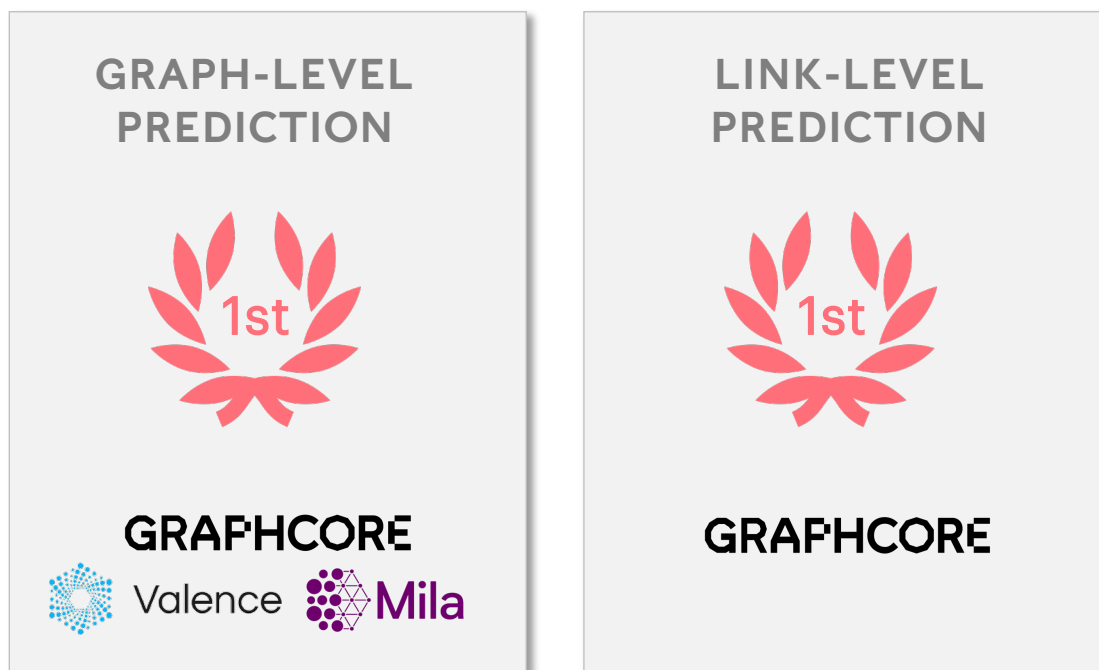
- 64 GB/s IPU-Link
- 100GE Host-Link Network I/F
- 100Gbps IPU-GW Link
- 32 GB/s PCIe G4





@ NeurIPS 2022

GRAPHCORE IPU ACHIEVES DOUBLE FIRST PLACE!



Open Graph Benchmark was established in 2020 with the aim of objectively measuring the performance of different graph models and compute systems

“As I started applying IPUs for molecular property predictions, I was shocked to see the speed improvements over traditional methods.”

Dominique Beaini, Research Team Lead at Valence Discovery and Associate Professor at Mila



<https://ogb.stanford.edu/neurips2022/results/>

OGB-LSC PCQM4MV2 CHALLENGE

THE IPU ADVANTAGE

- Simulating molecular properties using traditional methods (like DFT – Dense Functional Theory) is a very slow process
- Finding the optimal model & implementation required fast experimentation and innovation to explore combined benefits of GNN approaches with transformer-style attention
- The IPUs unique MIMD architecture and ultra-fast memory bandwidth enables :
 - Flexibility for **innovation**
 - High performance for **speed of experimentation**
- IPUs efficient scaling enabled quick experimentation on small models & efficient tuning on larger ‘production’ models



OGB-LSC 2022
GRAPH-LEVEL
PREDICTION

PCQM4MV2



GRAPHCORE



Valence



Mila

OGB-LSC WIKIKG90MV2 CHALLENGE

THE IPU ADVANTAGE

- Knowledge graph completion challenge using WikiKG90Mv2 dataset, based on the knowledge graph consisting of pages extracted from Wikipedia
- Dataset scale presents a problem for standard techniques
- This is addressed efficiently by exploitation of the IPU systems high capacity streaming memory, supplementing the large and ultra-fast In-Processor memory & inter-processor communication via IPU-Links
- This enabled quick iteration across the hyperparameter space and experimentation with new ideas, training of hundreds of models to convergence, and in the end construction of an ensemble of models for increased predictive power



OGB-LSC 2022
LINK-LEVEL
PREDICTION

WIKIKG90MV2



GRAPHCORE



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

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 IPU: UNLEASH THE POWER OF GRAPH NEURAL NETWORKS



Written By:
Blazej Banaszewski, Adam S
Akash Swamy & Mihai Polce

Apr 05, 2023 | Developer, GNN, Paperspace

GETTING STARTED WITH PYTORCH GEOMETRIC (PYG) ON GRAPHCORE IPU

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

dataset = GNN(dataset_num_features, dataset_num_classes)
model = GCN(dataset_num_features, dataset_num_classes)
optimizer = torch.optim.Adam(model.parameters(), lr=0.001)

print("Training on CPU")
for epoch in range(1, 5):
    optimizer.zero_grad()
```



PYTORCH GEOMETRIC + IPU



- Hardware lends itself to **GNNs** - fast gather scatter operations
- Already possible to run **PyTorch** on **IPUs**
- **PyTorch Geometric** is the PyTorch library to unify deep learning on graph-structured data
- Aim to make it as easy as possible to use **PyTorch Geometric** on **IPUs** and start **accelerating** your **GNNs**

AHEAD OF TIME COMPILATION

What?

- The model is **compiled** into a single compute graph with forward and backward passes.

Why?

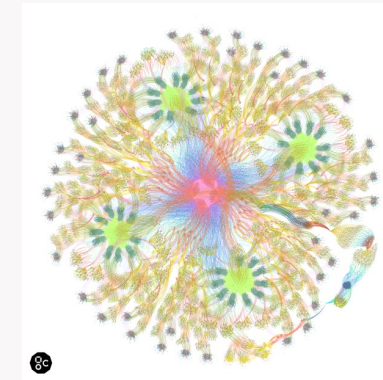
- **Efficient memory & communication**
- Allows **optimisations** to be applied during compilation

What does it mean for you?

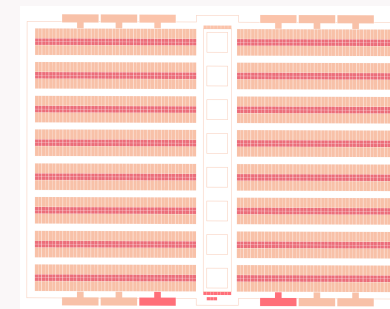
- All tensors in your model must be **fixed size**
- This includes the model **inputs**



PyTorch
model



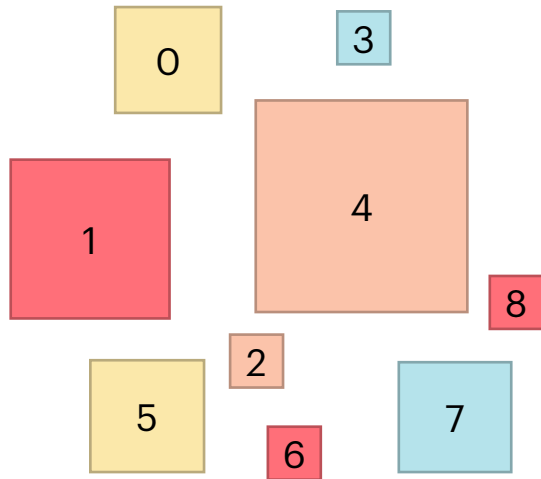
Compiled
graph



IPU

PYG MINI-BATCHING OF SMALL GRAPHS

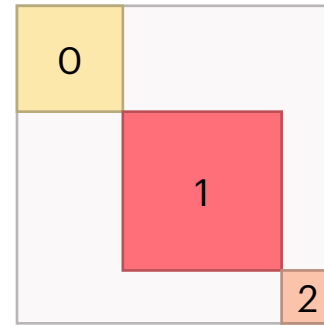
Adjacency of samples in dataset



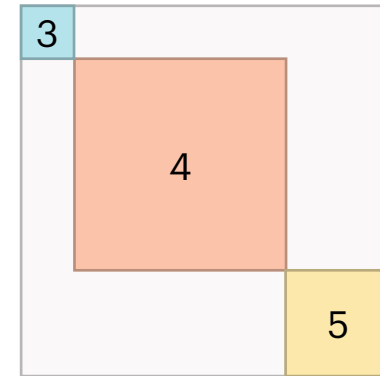
Mini-batching
(batch size 3)



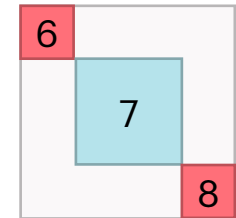
Adjacency of each mini-batch



Mini-batch 1

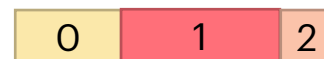


Mini-batch 2



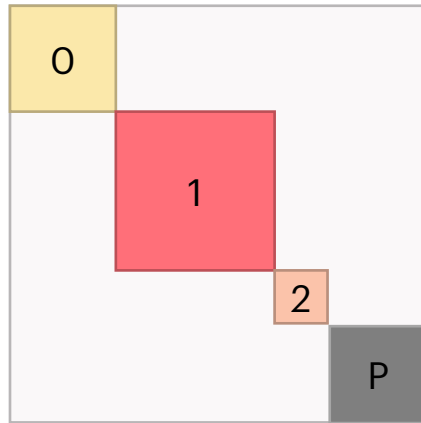
Mini-batch 3

Sparse representation of
each mini-batch

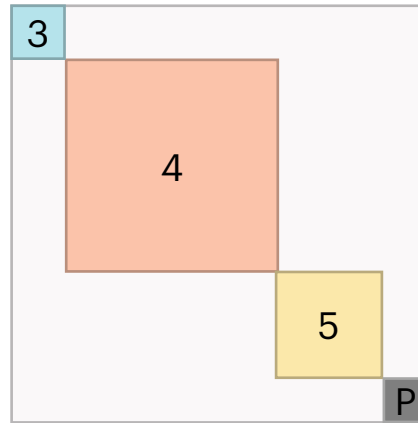


FIXED SIZE MINI BATCHING

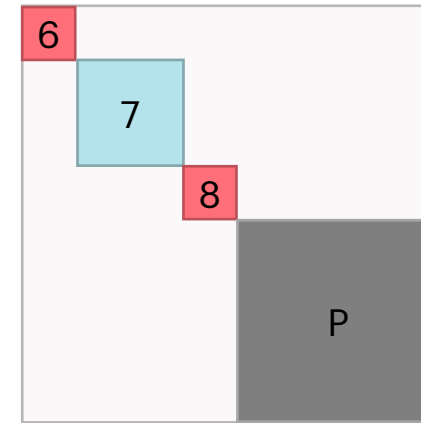
Adjacency
of each
mini-batch



Mini-batch 1



Mini-batch 2



Mini-batch 3

Sparse
representation
of each mini-
batch



- **Message passing** just works!
- Do we have to do any **masking**?

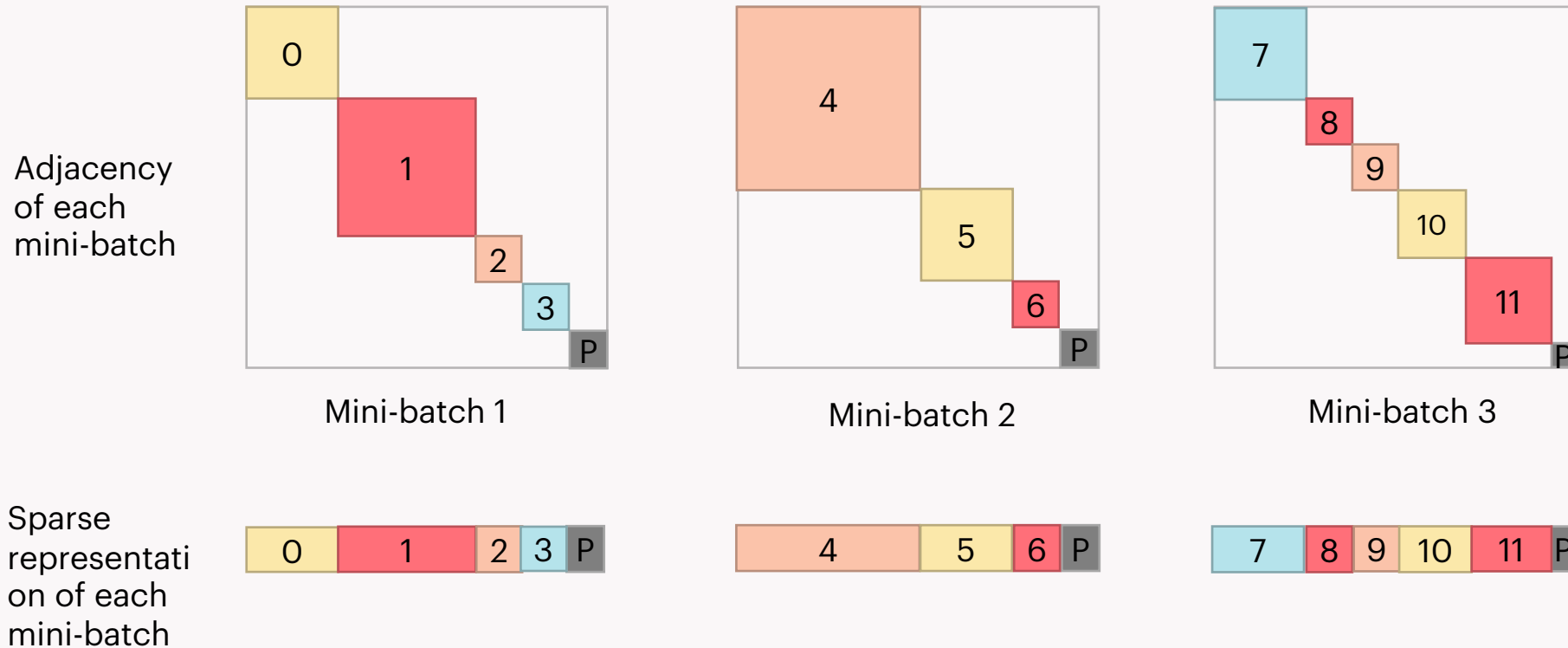
```
1 # Ignore final padded graph  
2 loss = F.mse_loss(x[:-1], y[:-1])
```

```
1 dataloader = FixedSizeDataLoader(dataset,  
2                                     num_nodes=300,  
3                                     num_edges=600,  
4                                     batch_size=10)
```



FIXED SIZE INPUTS WITH PACKING

Stream packing



Global packing

<https://arxiv.org/abs/2209.06354>



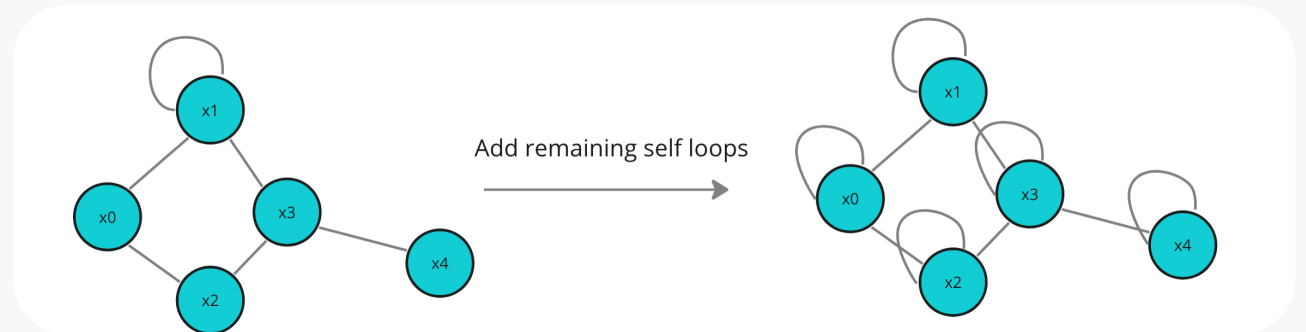
AND OTHER DYNAMIC THINGS

Other operations in your model may be **dynamic** that you wouldn't expect

Adding self-loops

```
1 # Adds remaining self loops
2 conv = GCNConv(10,
3               10,
4               add_self_loops=True)
```

```
1 transform = T.AddSelfLoops()
```

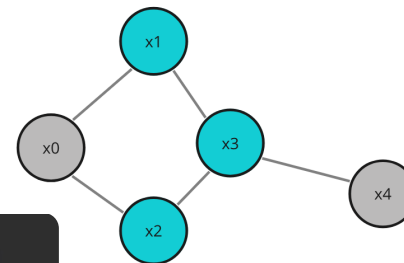


Using masks of different sizes each batch

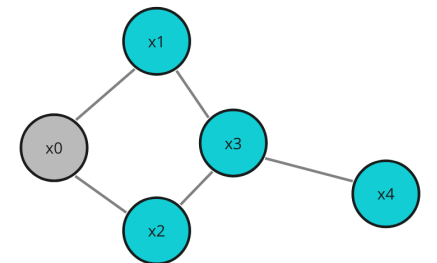
```
1 out = out[batch.train_mask]
```

```
1 out = torch.where(batch.train_mask, out, -100)
```

Batch 1

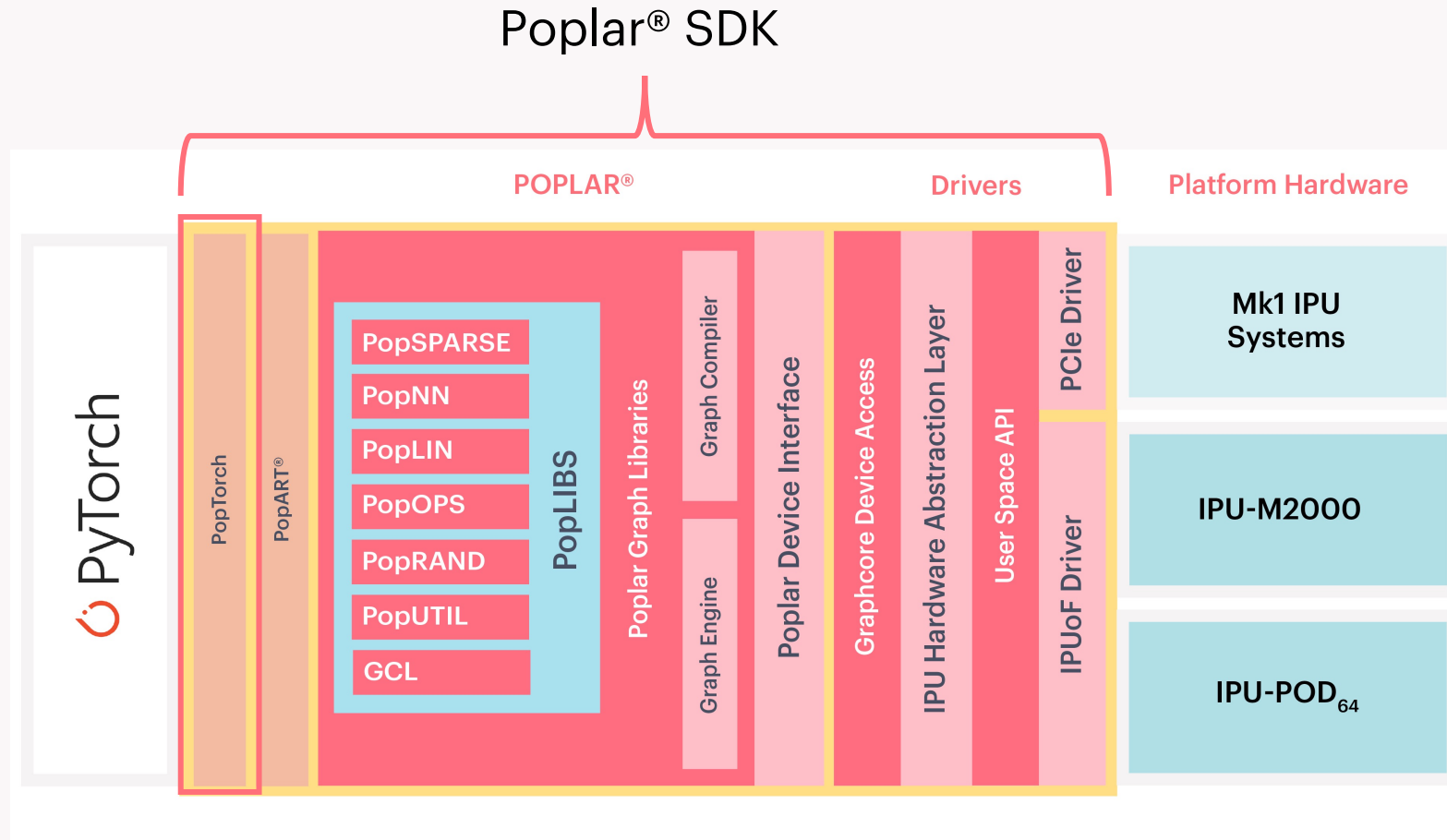


Batch 2



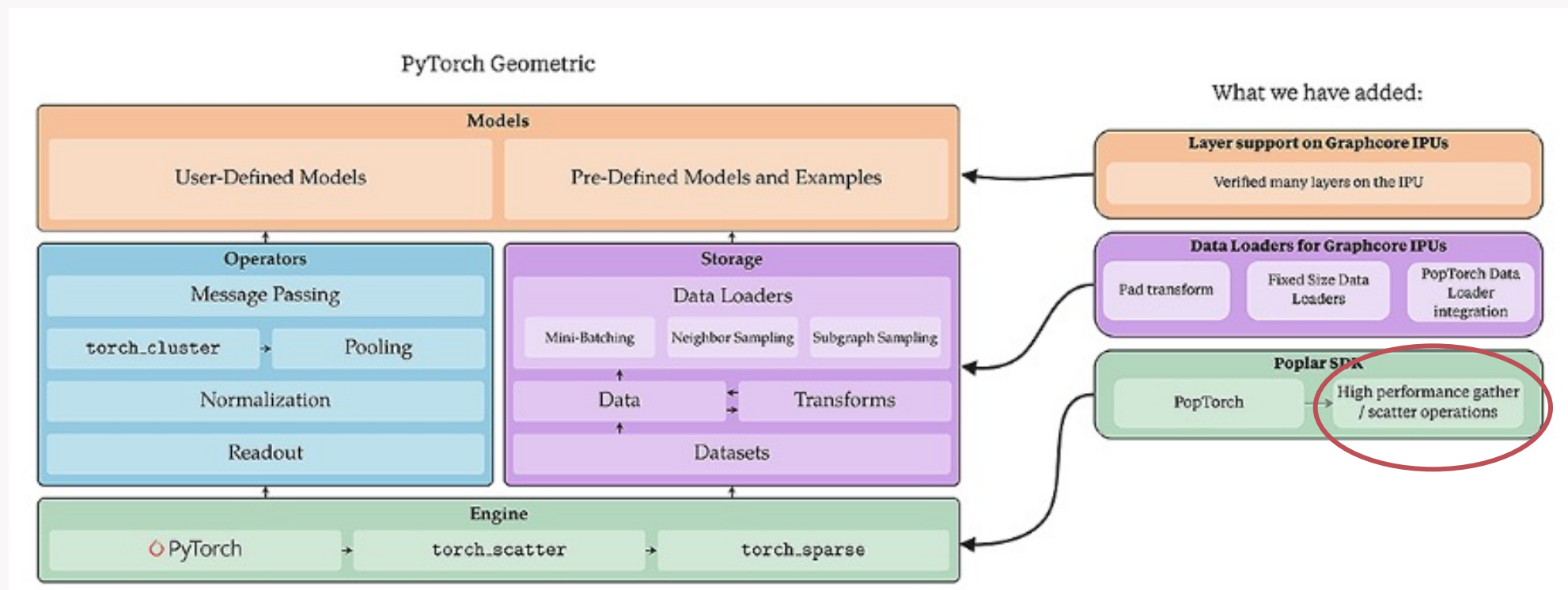
RUNNING PYG ON IPU: POPTORCH

PopTorch compiles PyTorch models into Poplar executables

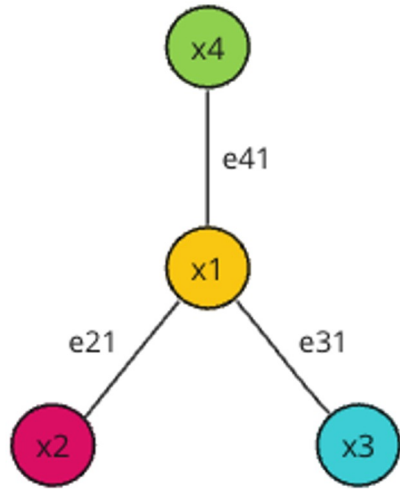


RUNNING PYG ON IPUS: POPTORCH GEOMETRIC

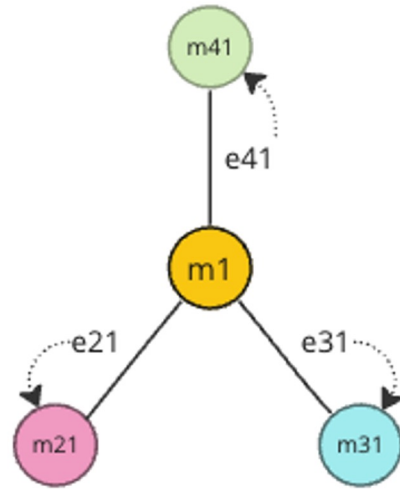
PopTorch Geometric enables GNN models to be run on Graphcore IPUs



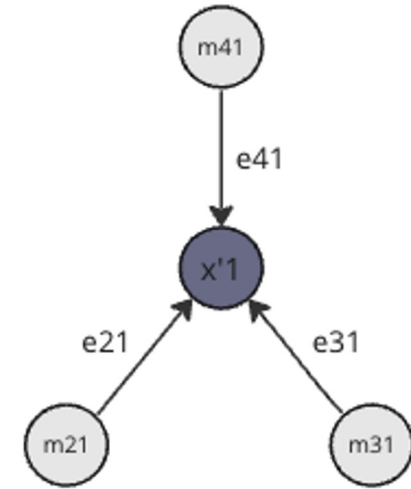
BENCHMARKING MESSAGE PASSING AS GATHER/SCATTER OPERATIONS



Original graph
x1 target node



Gather
Messages are collected from neighboring nodes



Scatter(-add)
Messages are aggregated along outbound edges

$$\vec{x}'_i = \gamma \left(\vec{x}_i, \underbrace{\text{aggregate}_{j \in \mathcal{N}(i)} [\phi(\vec{x}_i, \vec{x}_j, \vec{e}_{j,i})]}_{\text{Gather}} \right)$$

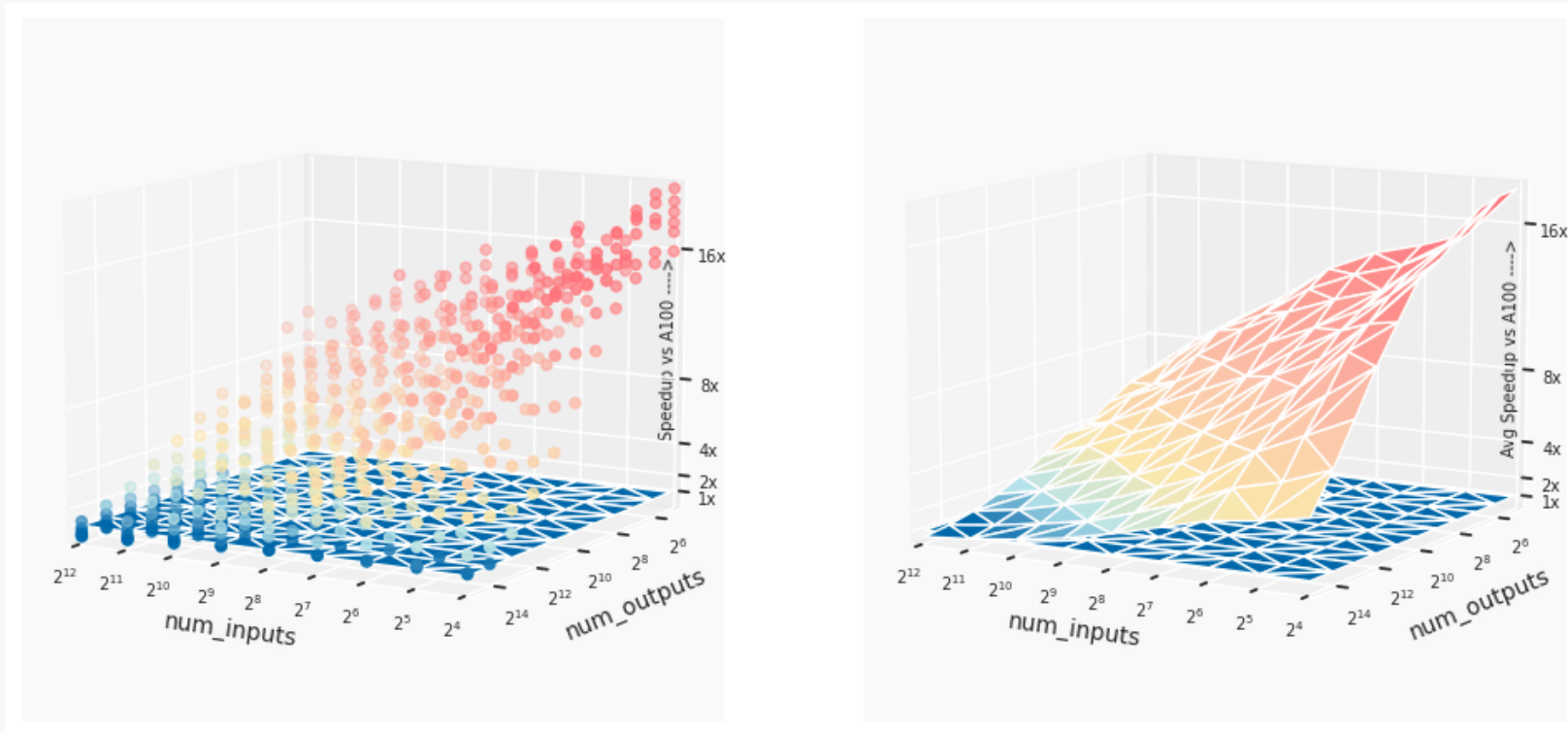
Update
Message

Scatter-reduce
Gather



HIGH PERFORMANCE SCATTER-ADD ON IPU

For small scatter input size, IPU achieves >16x speedups vs GPU

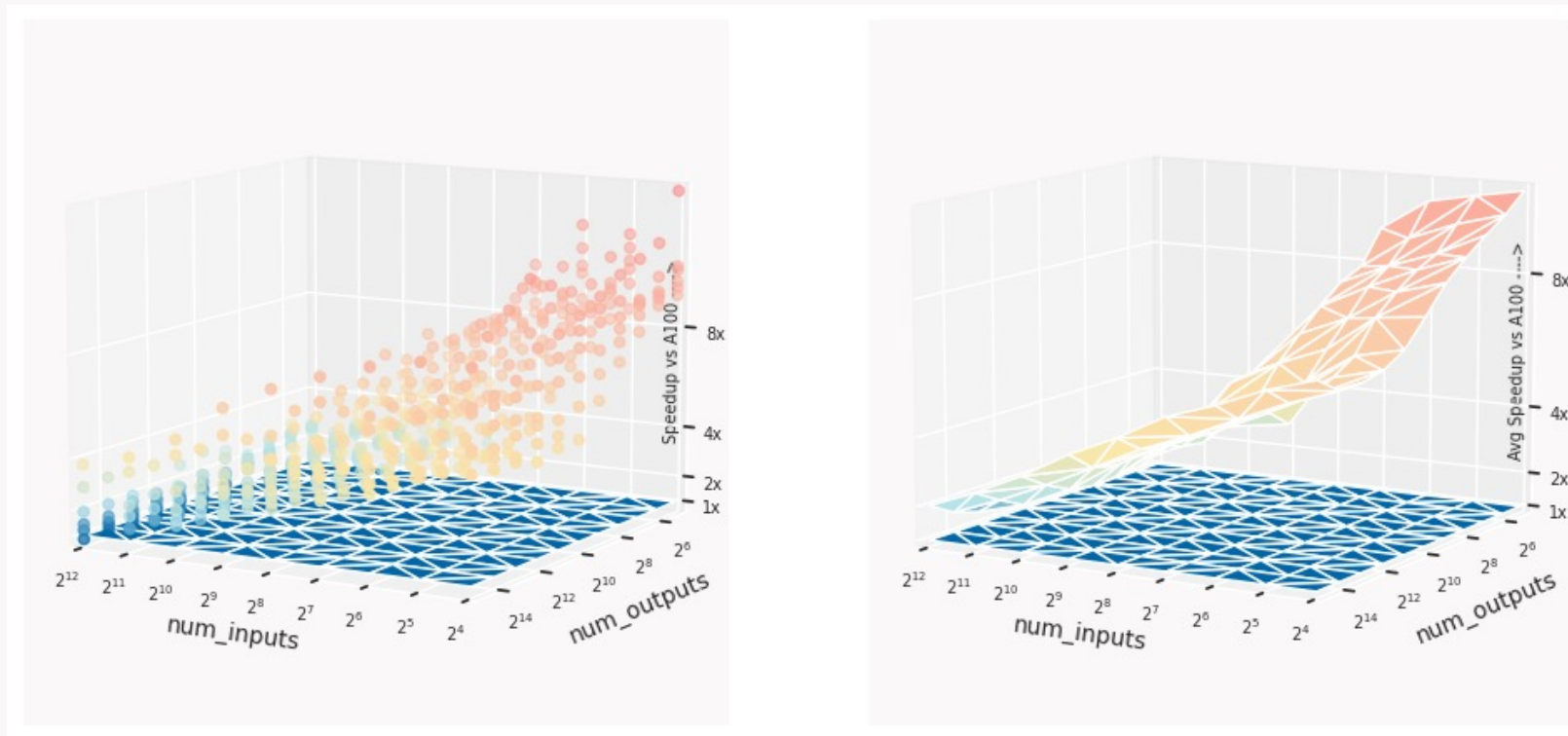


Graphcore BOW-M2000 vs NVIDIA A100 (1x, blue plane)



HIGH PERFORMANCE GATHER ON IPU

For small gather input size, IPU achieves >8x speedups vs GPU



Graphcore BOW-M2000 vs NVIDIA A100 (1x, blue plane)

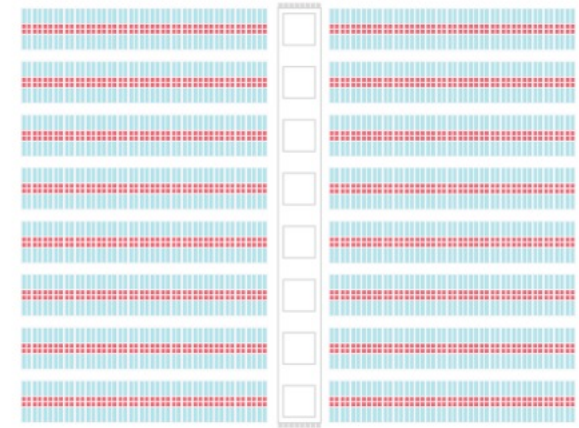
HIGH PERFORMANCE GATHER-SCATTER OPS ON IPUS

Why faster on IPUs?

- Large, high bandwidth on-chip SRAM.
- Support for fine-grained parallelism.
- Fast all-to-all communication links.

IPU

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



Model and data tightly coupled, and large locally distributed SRAM

CASE STUDY: SCHNET FOR MOLECULAR PROPERTY PREDICTION

Training on Graphcore IPUs with PyG

Use the QM9 dataset from MoleculeNet to train the SchNet model to predict a graph-level property, the HOMO-LUMO energy gap

Copyright (c) 2023 Graphcore Ltd. All rights reserved.

Molecular property prediction on IPU using SchNet - Training with PyTorch Geometric

This notebook demonstrates training a [SchNet graph neural network](#) with PyTorch Geometric on the Graphcore IPU. We will use the QM9 dataset from the [MoleculeNet: A Benchmark for Molecular Machine Learning](#) paper and train the SchNet model to predict the HOMO-LUMO energy gap.

Domain	Tasks	Model	Datasets	Workflow	Number of IPUs	Execution time
GNNs	Graph Classification	SchNet	QM9	Training, evaluation	recommended: 16 (min: 4)	~4 minutes

This notebook assumes some familiarity with PopTorch as well as PyTorch Geometric (PyG). For additional resources please consult:

- [PopTorch Documentation](#)
- [PopTorch Examples and Tutorials](#)
- [PyTorch Geometric](#)
- [PopTorch Geometric Documentation](#)

Running on Paperspace

The Paperspace environment lets you run this notebook with no set up. To improve your experience we preload datasets and pre-install packages, this can take a few minutes, if you experience errors immediately after starting a session please try restarting the kernel before contacting support. If a problem persists or you want to give us feedback on the content of this notebook, please reach out to through our community of developers using our [slack channel](#) or raise a [GitHub issue](#).



CASE STUDY: SCHNET FOR MOLECULAR PROPERTY PREDICTION

Notebook walkthrough

QM9 dataset

Molecular properties of interest to train SchNet are:

- `z` atomic number for each atom in the molecule
- `pos` contains the 3D structure of the molecule
- `y` contains the 19 regression targets: we slice it `y[:, 4]` where the HOMO-LUMO gap is stored

```
Run 6
1 datum = dataset[123244]
2 datum, datum.z, datum.pos, datum.y[:, 4]

(Data(x=[13, 11], edge_index=[2, 28], edge_attr=[28, 4], y=[1, 19], pos=[13, 3], idx=[1], name='gdb_125563', z=[13]),
tensor([[6, 6, 7, 7, 6, 7, 7, 7, 1, 1, 1, 1]),
tensor([[ -2.7500e-02,  1.4963e+00,  5.2800e-02],
        [-9.1000e-03,  1.2800e-02,  2.6000e-03],
        [-4.2200e-02, -7.5060e-01, -1.0686e+00],
        [-9.3000e-03, -2.1018e+00, -7.2150e-01],
        [ 4.3600e-02, -2.0859e+00,  5.8330e-01],
        [ 1.0010e-01, -2.8658e+00,  1.7010e+00],
        [ 1.3480e-01, -2.0775e+00,  2.8101e+00],
        [ 1.0380e-01, -8.4570e-01,  2.4693e+00],
        [ 4.7300e-02, -8.2400e-01,  1.1011e+00],
        [ 8.7460e-01,  1.8923e+00,  5.3110e-01],
        [-8.0800e-02,  1.8742e+00, -9.6890e-01],
        [-8.9120e-01,  1.8675e+00,  6.1440e-01],
        [ 1.1880e-01, -3.8660e+00,  1.7999e+00]]),
tensor([[5.2708]]))
```



CASE STUDY: SCHNET FOR MOLECULAR PROPERTY PREDICTION

Notebook walkthrough

Data loading and mini-batching

```
Run | 12  
  
1 loader = DataLoader(dataset, batch_size=4)  
2  
3 it = iter(loader)  
4 next(it), next(it)  
  
(DataBatch(y=[4], pos=[16, 3], z=[16], batch=[16], ptr=[5]),  
DataBatch(y=[4], pos=[21, 3], z=[21], batch=[21], ptr=[5]))
```

AOT compilation requirement on IPU
The mini-batches will need to be adapted to be fixed size

Padding individual dataset samples

```
Run | 17  
  
1 data = Batch.from_data_list([dataset[0]])  
2 pad_transform = Pad(32, node_pad_value=AttrNamePadding({"z": 0, "pos": 0, "batch": 1}))  
3 padded_batch = pad_transform(data)  
4 padded_batch  
  
DataBatch(y=[1], pos=[32, 3], z=[32], batch=[32], ptr=[2], num_nodes=32)
```



CASE STUDY: SCHNET FOR MOLECULAR PROPERTY PREDICTION

Notebook walkthrough

Efficient data loading: padding the mini-batch

```
Run | 23  
1 batch_size = 8  
2  
3 dataloader = CustomFixedSizeDataLoader(  
4     dataset, batch_size=batch_size, num_nodes=32 * (batch_size - 1)  
5 )
```

```
# slice off the padding molecule and calculate the mse loss  
prediction = prediction[0:-1]  
target = target[0:-1]  
loss = F.mse_loss(prediction, target)  
return prediction, loss
```

The mini-batches have now the same sizes



```
Run | 24  
1 dataloader_iter = iter(dataloader)  
2 first_batch = next(dataloader_iter)  
3 second_batch = next(dataloader_iter)  
4 print(first_batch)  
5 print(second_batch)
```

```
DataBatch(y=[8], pos=[224, 3], batch=[224], ptr=[9], z=[224], num_nodes=224, num_edges=0)  
DataBatch(y=[8], pos=[224, 3], batch=[224], ptr=[9], z=[224], num_nodes=224, num_edges=0)
```



CASE STUDY: SCHNET FOR MOLECULAR PROPERTY PREDICTION

Notebook walkthrough

Train SchNet on IPU

Select your hyperparameters and PopTorch options:

```
Run | 27
1 replication_factor = int(num_ipus)
2 device_iterations = 32
3 gradient_accumulation = max(1, 16 // replication_factor)
4 learning_rate = 1e-4
5 num_epochs = 5
```

```
Run | 28
1 options = poptorch.Options()
2 options.enableExecutableCaching(executable_cache_dir)
3 options.outputMode(poptorch.OutputMode.All)
4 options.deviceIterations(device_iterations)
5 options.replicationFactor(replication_factor)
6 options.Training.gradientAccumulation(gradient_accumulation)
```

Recreate the data loader to pass it the selected hyperparameters and options, define the model and compile it on IPU:

```
Run | 32
1 torch.manual_seed(0)
2 knn_graph = KNNInteractionGraph(cutoff=cutoff, k=28)
3 model = SchNet(cutoff=cutoff, interaction_graph=knn_graph)
4 model.train()
5 model = TrainingModule(
6     model, batch_size=batch_size, replace_softplus=additional_optimizations
7 )
8 optimizer = poptorch.optim.AdamW(model.parameters(), lr=learning_rate)
9 training_model = poptorch.trainingModel(model, options, optimizer)
10
11 data = next(iter(train_loader))
12 training_model.compile(data.z, data.pos, data.batch, data.y)
```

Graph compilation: 100% | ██████████ | 100/100 [00:05<00:00]



CASE STUDY: SCHNET FOR MOLECULAR PROPERTY PREDICTION

Notebook walkthrough

Train SchNet on IPU

Define the training loop and finally plot the mean of the loss

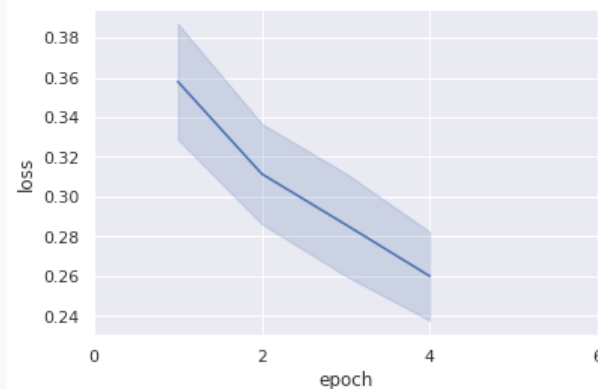
Run 33

```
1 train = []
2
3 for epoch in range(num_epochs):
4     bar = tqdm(train_loader)
5     for i, data in enumerate(bar):
6         _, mini_batch_loss = training_model(data.z, data.pos, data.batch, data.y)
7         loss = float(mini_batch_loss.mean())
8         train.append({"epoch": epoch, "step": i, "loss": loss})
9         bar.set_description(f"Epoch {epoch} loss: {loss:0.6f}")
```

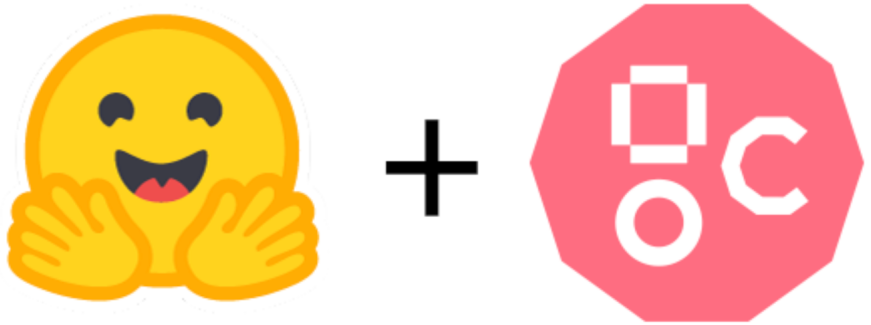
```
Epoch 0 loss: 0.375165: 100%|██████████| 30/30 [00:24<00:00, 1.21it/s]
Epoch 1 loss: 0.292912: 100%|██████████| 30/30 [00:24<00:00, 1.21it/s]
Epoch 2 loss: 0.270268: 100%|██████████| 30/30 [00:25<00:00, 1.20it/s]
Epoch 3 loss: 0.255817: 100%|██████████| 30/30 [00:24<00:00, 1.21it/s]
Epoch 4 loss: 0.233512: 100%|██████████| 30/30 [00:24<00:00, 1.21it/s]
```

Run 35

```
1 df = pd.DataFrame(train)
2 g = sns.lineplot(data=df[df.epoch > 0], x="epoch", y="loss", errorbar="sd")
3 g.set_xticks(range(0, num_epochs + 2, 2))
4 g.figure.show()
```



LARGE LANGUAGE MODELS



Optimum Graphcore

👏 Optimum Graphcore is the interface between the 🧠 Transformers library and [Graphcore IPU](#)s. 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](#):

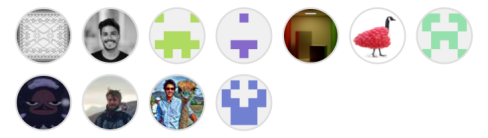
IPUs are the processors that power Graphcore's IPU-POD datacenter compute systems. This new type of processor is designed to support the very specific computational requirements of AI and machine learning. Characteristics such as fine-grained parallelism, low precision arithmetic, and the ability to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

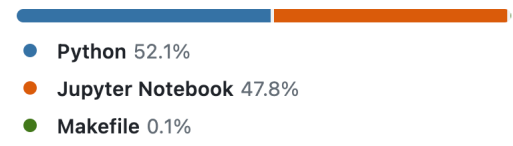


Contributors 36



[+ 25 contributors](#)

Languages



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,
```



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 Answ
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.



Code

main

Go to file

- > .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
 - > stable_diffusion
 - > text_embeddings_models
 - > wav2vec2
 - 📄 README.md
 - 📄 audio_classification.ipynb
 - 📄 deberta-blog-notebook.ipynb
 - 📄 external_model.ipynb
 - 📄 flan_t5_inference.ipynb
 - 📄 image_classification.ipynb

optimum-graphcore / notebooks /

↑ Top

Introduction to Optimum Graphcore	Introduce Optimum-Graphcore with a BERT fine-tuning example.	
Sentiment analysis with pipelines	Use the sentiment-analysis pipeline to quickly evaluate pre-trained models on the IPU.	
Real Time Name Entity Recognition on the IPU	Use Gradio and pipelines to prototype a web application doing fast token classification.	
Train an external model	Show how to train an external model that is not supported by Optimum or Transformers.	
Train your language model	Show how to train a model for causal or masked language modelling from scratch.	
How to fine-tune a model on text classification	Show how to preprocess the data and fine-tune a pretrained model on any GLUE task.	
How to fine-tune a model on language modeling	Show how to preprocess the data and fine-tune a pretrained model on a causal or masked LM task.	Coming soon on Gradient
How to fine-tune a model on token classification	Show how to preprocess the data and fine-tune a pretrained model on a token classification task (NER, PoS).	
How to fine-tune a model on question answering	Show how to preprocess the data and fine-tune a pretrained model on SQUAD.	
How to fine-tune a model on multiple choice	Show how to preprocess the data and fine-tune a pretrained model on SWAG.	
How to fine-tune a model on translation	Show how to preprocess the data and fine-tune a pretrained model on WMT.	
How to fine-tune a model on summarization	Show how to preprocess the data and fine-tune a pretrained model on XSUM.	
How to fine-tune a model on audio classification	Show how to preprocess the data and fine-tune a pretrained Speech model on Keyword Spotting	Coming soon on Gradient
How to fine-tune a model on image classification	Show how to preprocess the data and fine-tune a pretrained model on image classification.	
wav2vec 2.0 Fine-Tuning on IPU	How to fine-tune a pre-trained wav2vec 2.0 model with PyTorch on the Graphcore IPU-POD16 system.	
wav2vec 2.0 Inference on IPU	How to run inference on the wav2vec 2.0 model with PyTorch on the Graphcore IPU-POD16 system.	
Stable Diffusion Text-to-Image generation	Run a Stable Diffusion (Conditional UNet) pipeline on the text-to-image generation task.	



Code

master + 🔍

Go to file 🔍

- > ai_for_simulation
- > finance
- > gnn
- > multimodal
- ▼ nlp
 - > bert
 - > bloom/popxl
 - > dolly/popxl
 - > gpt2/pytorch
 - > gpt3_175B/popxl**
 - > gpt3_2.7B/popxl
 - > gpt_j/popxl
 - > t5/popxl
- > preview
- > probability
- > speech
- > tutorials
- > utils
- > vision

.git-blame-ignore-revs

README.md



GPT-3 training on IPU's using PopXL

Framework	Domain	Model	Datasets	Tasks	Training	Inference	Reference
PopXL	NLP	GPT-3	Wikipedia	Next sentence prediction, Question/Answering	✅ Min. 256 IPU's (POD256) required	❌	Language Models are Few-Shot Learners

This README describes how to run [GPT-3](#) models for NLP pre-training on Graphcore IPU's using the PopXL library. A combination of phased execution, tensor model parallelism, data parallelism, and remote tensor sharding are utilised to train the models.

This application shows how to run larger models on IPU. The techniques to do this mean that performance is lower than for models that fit in IPU memory. Large model training or fine-tuning requires a big Pod installation. The minimum to run pre-training with this model is a Pod256. PopXL is an experimental framework and may be subject to change in future releases.

Instructions summary

1. Install and enable the Poplar SDK (see [Poplar SDK setup](#))
2. Install the system and Python requirements (see [Environment setup](#))
3. Download the WIKI-103 dataset (See [Dataset setup](#))

Poplar SDK setup

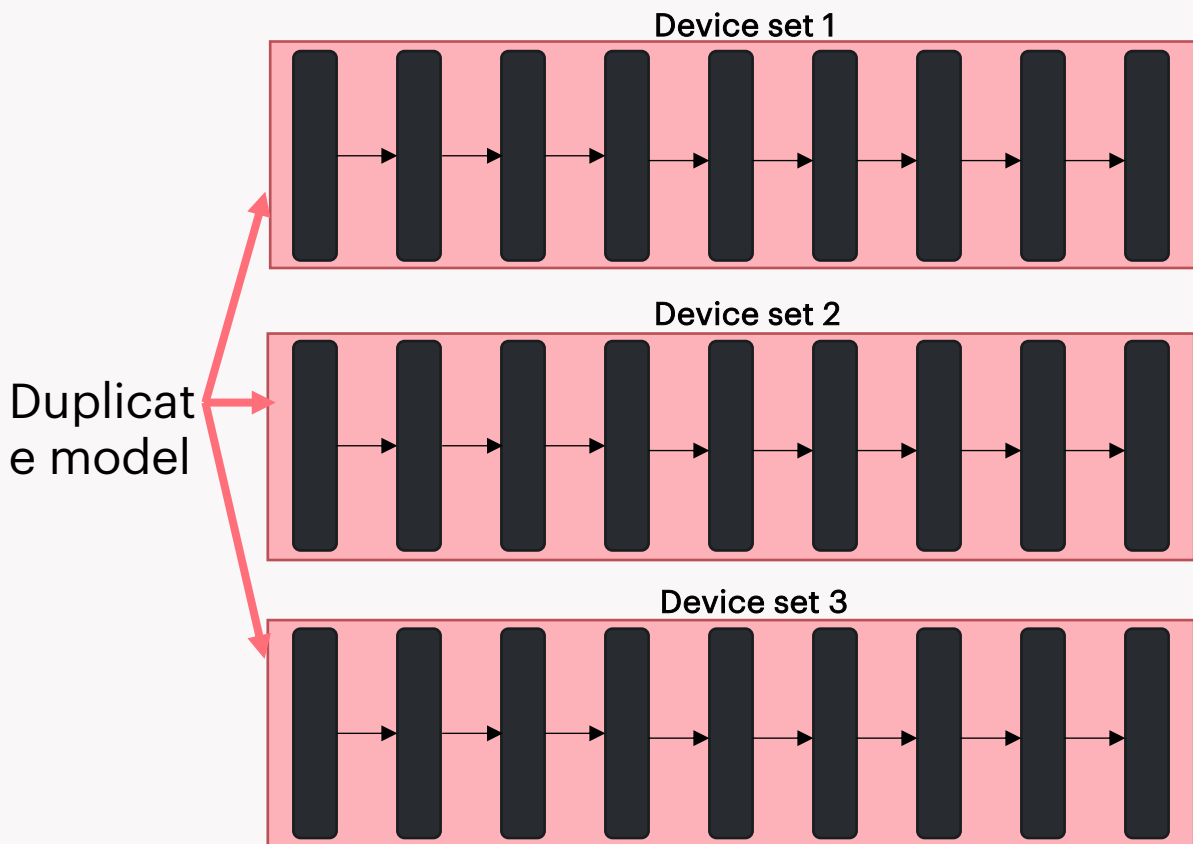
To check if your Poplar SDK has already been enabled, run:

```
echo $POPLAR_SDK_ENABLED
```

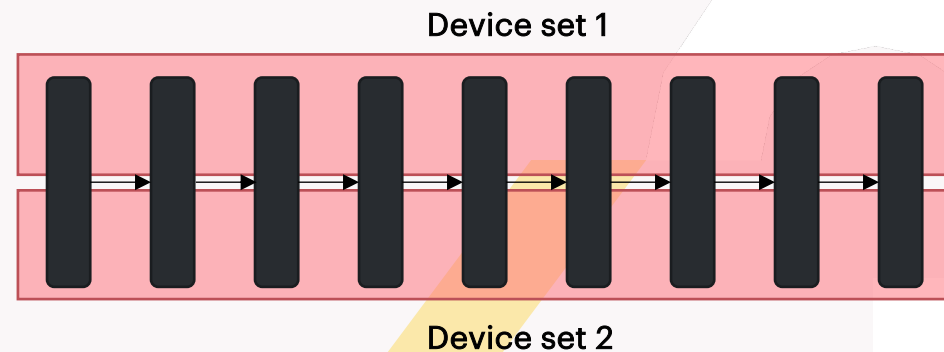


Modes of Execution

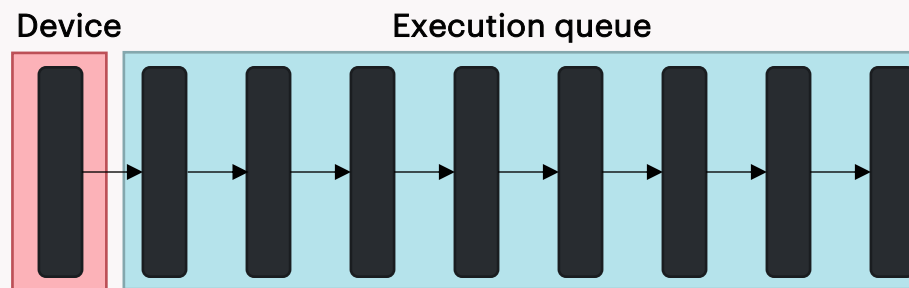
Data Parallelism (DP)



Tensor Parallelism (TP)



Phased Execution (PE)



Matmul TP

- Consider sharding a matmul in two ways:

$$f(X) = XA$$

$$XA = \begin{pmatrix} X_0 & X_1 \\ X_2 & X_3 \end{pmatrix} \begin{pmatrix} A_0 & A_1 \\ A_2 & A_3 \end{pmatrix}$$

$$XA = \begin{pmatrix} X_0A_0 + X_1A_2 & X_0A_1 + X_1A_3 \\ X_2A_0 + X_3A_2 & X_2A_1 + X_3A_3 \end{pmatrix}$$

$$A_A \rightarrow \{m, k_A\}$$

$$A_B \rightarrow \{m, k_B\}$$

$$\begin{pmatrix} X_0 & X_1 \\ X_2 & X_3 \end{pmatrix} \begin{pmatrix} A_0 & A_1 \\ A_2 & A_3 \end{pmatrix} = \left(\begin{pmatrix} X_0 & X_1 \\ X_2 & X_3 \end{pmatrix} \begin{pmatrix} A_0 \\ A_2 \end{pmatrix} \quad \begin{pmatrix} X_0 & X_1 \\ X_2 & X_3 \end{pmatrix} \begin{pmatrix} A_1 \\ A_3 \end{pmatrix} \right)$$

$$= \underbrace{(XA_A \quad XA_B)}_{\text{Concatenation}}$$

Concatenation

Column-wise sharding:

$$(XA, XA) := \text{AllGather}(XA_A, XA_B)$$

$$X \rightarrow \{n, m\}$$

$$A \rightarrow \{m, k\}$$

$$X_A \rightarrow \{n, m_A\}$$

$$X_B \rightarrow \{n, m_B\}$$

$$A_A \rightarrow \{m_A, k\}$$

$$A_B \rightarrow \{m_B, k\}$$

$$\begin{pmatrix} X_0 & X_1 \\ X_2 & X_3 \end{pmatrix} \begin{pmatrix} A_0 & A_1 \\ A_2 & A_3 \end{pmatrix} = \begin{pmatrix} X_0 \\ X_2 \end{pmatrix} \begin{pmatrix} A_0 & A_1 \end{pmatrix} + \begin{pmatrix} X_1 \\ X_3 \end{pmatrix} \begin{pmatrix} A_2 & A_3 \end{pmatrix}$$

$$= \underbrace{X_A A_A + X_B A_B}_{\text{Summation}}$$

Summation

Row-wise sharding:

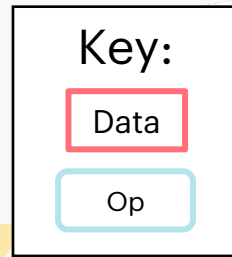
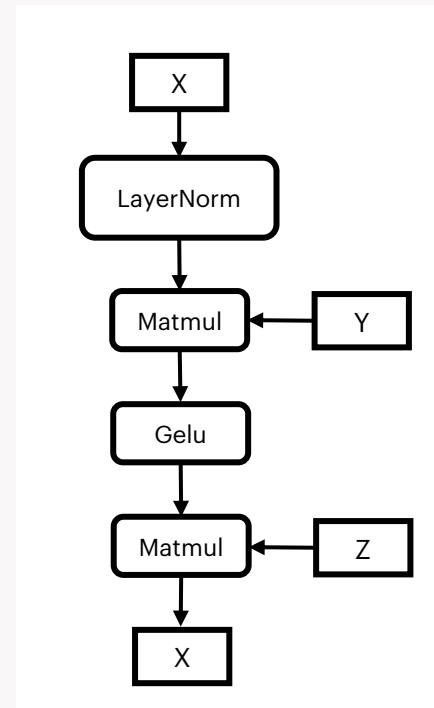
$$(XA, XA) := \text{AllReduce}(X_A A_A, X_B A_B)$$



Feed-Forward Layer: No Parallelism

Shapes

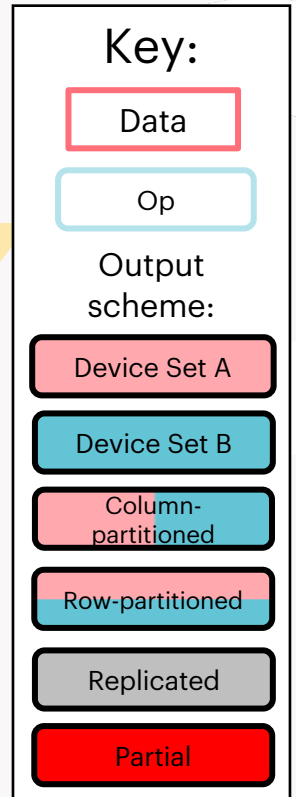
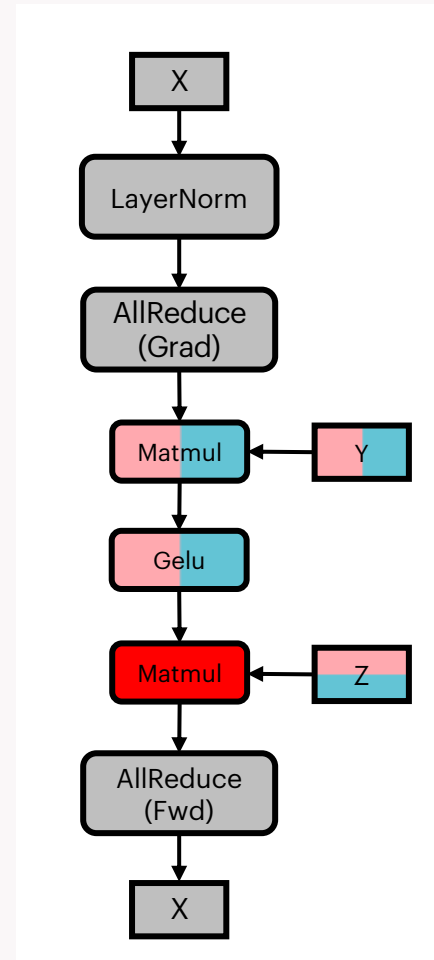
Data/ Op output	Standard
X	[s, h]
Y	[h, 4h]
X := X @ Y	[s, 4h]
Z	[4h, h]
X := X @ Z	[s, h]



Feed-Forward Layer: 1D Tensor Parallelism

Shapes

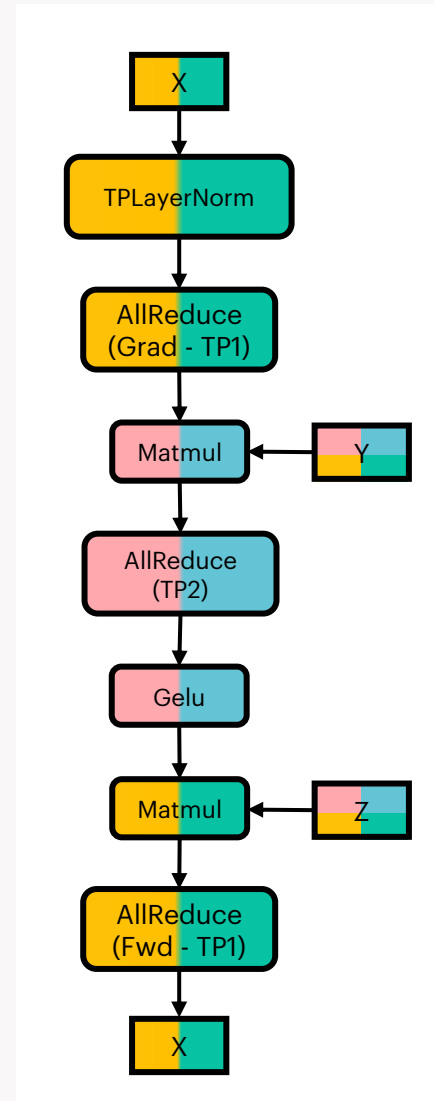
Data/ Op output	Standard	1DTP
X	[s, h]	[s, h]
Y	[h, 4h]	[h, 4h/tp1]
$X := X @ Y$	[s, 4h]	[s, 4h/tp1]
Z	[4h, h]	[4h/tp1, h]
$X := X @ Z$	[s, h]	[s, h]



Feed-Forward Layer: 2D Tensor Parallelism

Shapes

Data/ Op output	Standard	1DTP	2DTP
X	[s, h]	[s, h]	[s, h/tp2]
Y	[h, 4h]	[h, 4h/tp1]	[h/tp2, 4h/tp1]
$X := X @ Y$	[s, 4h]	[s, 4h/tp1]	[s, 4h/tp1]
Z	[4h, h]	[4h/tp1, h]	[4h/tp1, h/tp2]
$X := X @ Z$	[s, h]	[s, h]	[s, h/tp2]



Key:

- Data
- Op

Output scheme:

- TP1 Set A
- TP1 Set B
- TP2 Set A
- TP2 Set B
- Column-partitioned
- Row-partitioned
- Multi-partitioned
- Replicated
- Partial



Model-Device Mapping

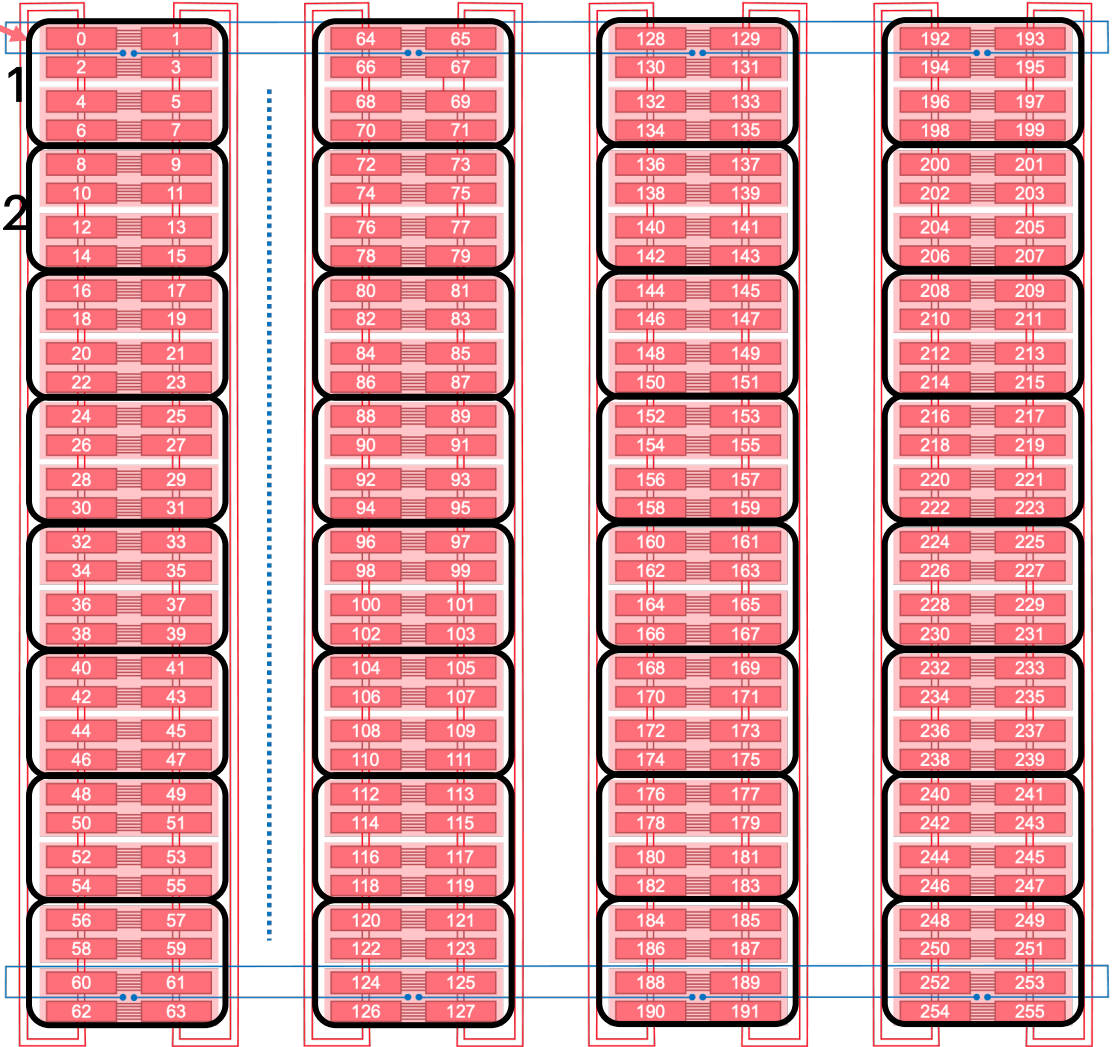
TP2 Groups (size 8)

Pod 256

Group 1

Group 2

...



Rack 0

Rack 1

Rack 2

Rack 4

— IPU Link 15.75GB/s raw

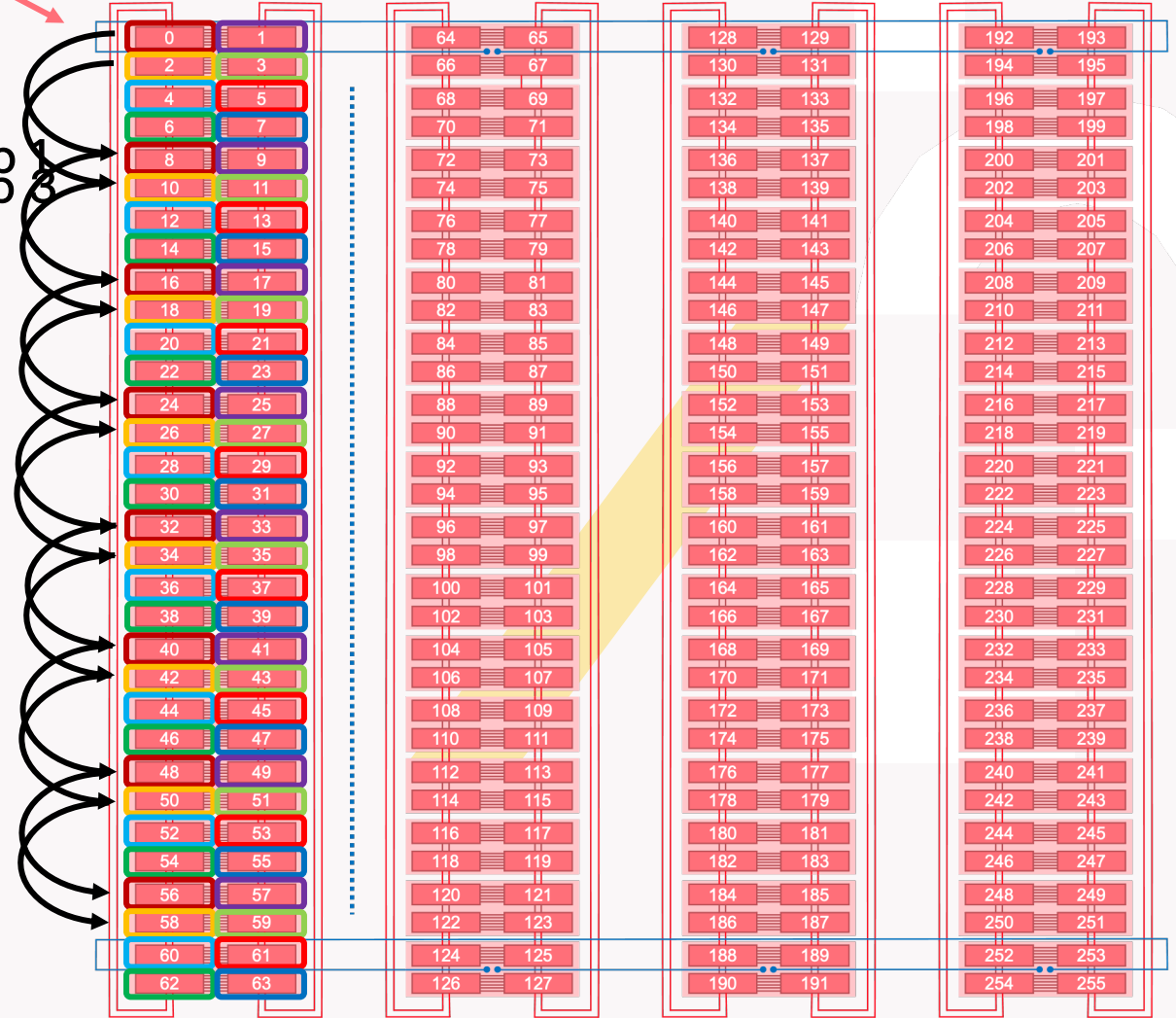
— Gateway Link 12.5GB/s raw

Model-Device Mapping

TP1 Groups (size 8, stride 8)

Pod 256

Group Group 8



Rack 0

Rack 1

Rack 2

Rack 4

— IPU Link 15.75GB/s raw

— Gateway Link 12.5GB/s raw

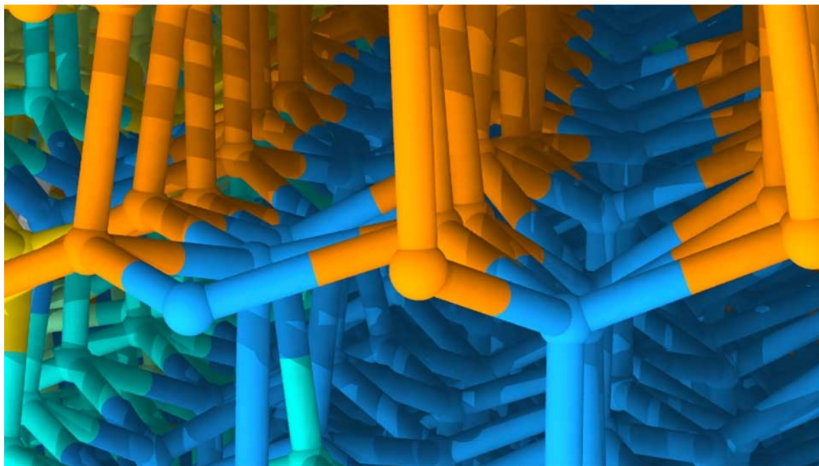


APPLY AND JOIN TODAY



Director's Discretionary Allocation Program

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

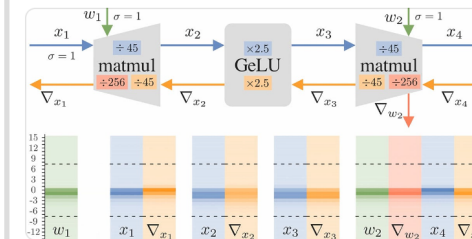
- You can find our blog post here: <https://www.graphcore.ai/posts/simple-fp16-and-fp8-training>
- Paperspace notebook (try it yourself!): <https://ipu.dev/qXfm2a>
- Arxiv paper: <https://arxiv.org/abs/2303.11257>

(& 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



Join at graphcore.ai/join-community

TUESDAY, 11 JUNE



- 1:00 PM** → 1:15 PM **Introduction** ⌚ 15m
- 1:15 PM** → 1:45 PM **Graphcore BowPod64 Hardware** ⌚ 30m
- 1:45 PM** → 2:30 PM **Software Stack: TensorFlow, PyTorch, and Poplar** ⌚ 45m
- 2:30 PM** → 2:45 PM **Break** ⌚ 15m
- 2:45 PM** → 3:15 PM **Porting applications with Poplar** ⌚ 30m
- 3:15 PM** → 4:00 PM **How to use Bow Pod64@ ALCF** ⌚ 45m

WEDNESDAY, 12 JUNE



- 1:00 PM** → 1:45 PM **Deep Dive on Graph neural networks and Large Language Models** ⌚ 45m
- 1:45 PM** → 2:15 PM **Profiling with PopVision** ⌚ 30m
- 2:15 PM** → 2:30 PM **Break** ⌚ 15m
- 2:30 PM** → 3:15 PM **Hands-on session** ⌚ 45m
- 3:15 PM** → 4:00 PM **Best Practices, Q&A** ⌚ 45m



THANK YOU!
Q&A