

GRAPHCORE

POPART – SUPPORT ONNX ON IPU

Han Zhao 赵晗

Graphcore Tech Team Lead

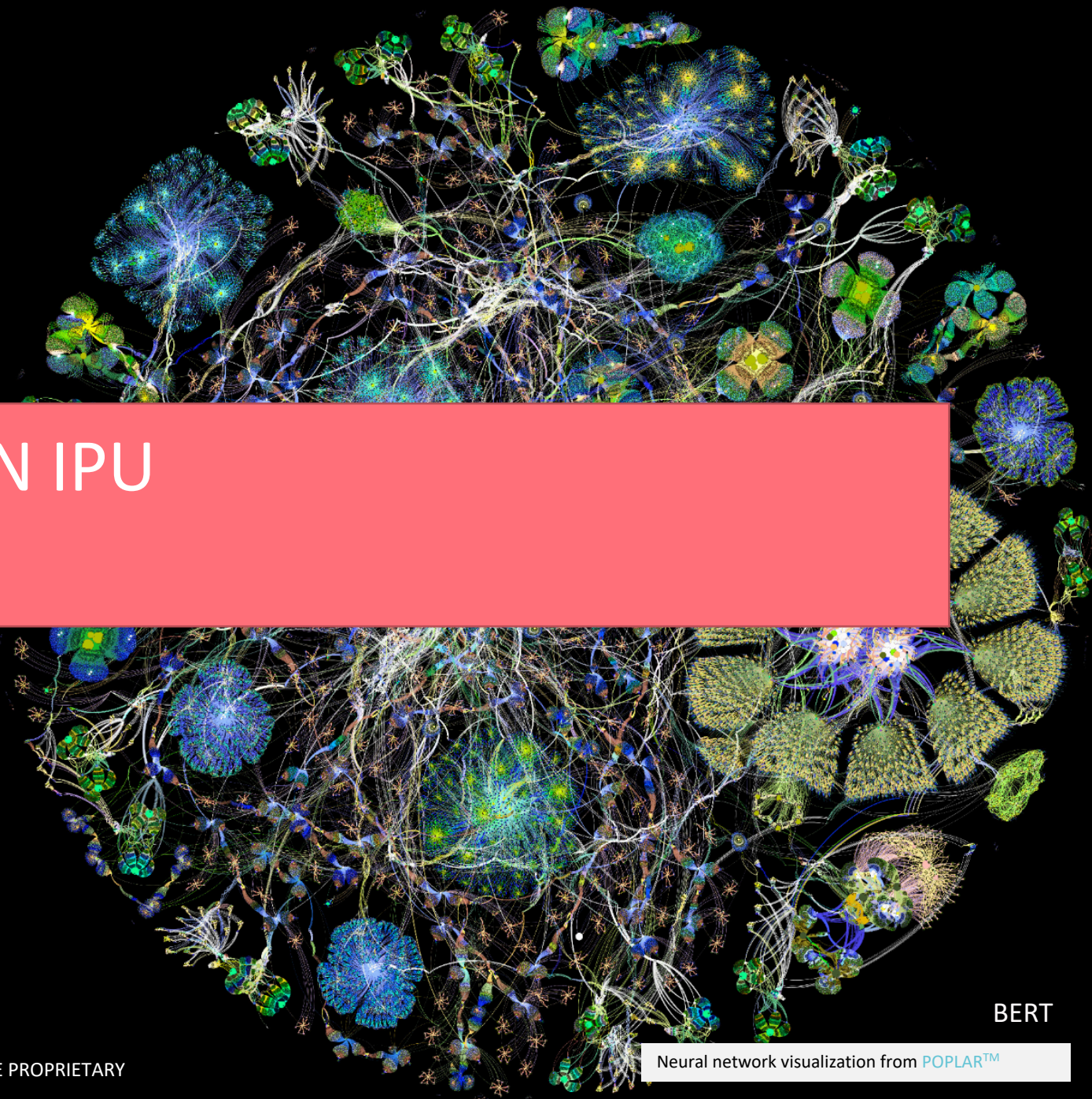
hanzhao@graphcore.ai



GRAPHCORE PROPRIETARY

BERT

Neural network visualization from POPLAR™



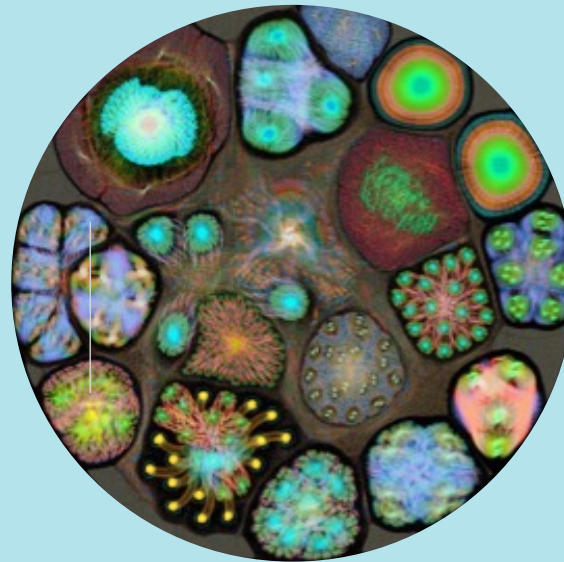
ABOUT GRAPHCORE

Hardware



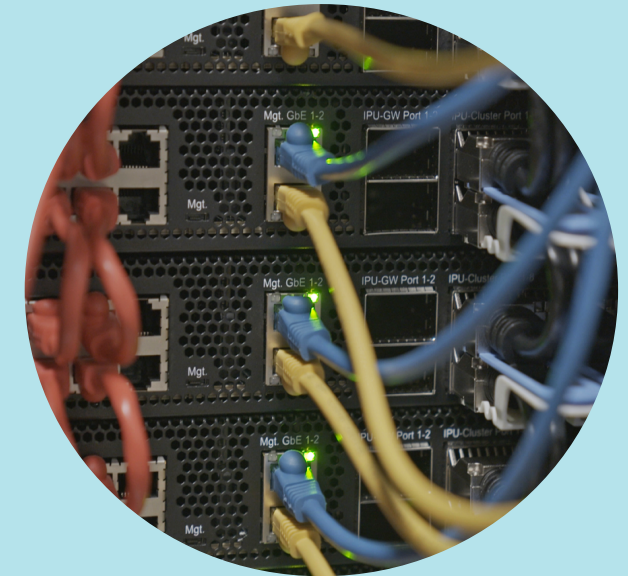
IPU processors
designed for AI

Software



Poplar® software stack &
development tools

Platforms



M2000 and Server
IPU-POD₆₄ scale-out

IPU-Tiles™

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

IPU-Core™

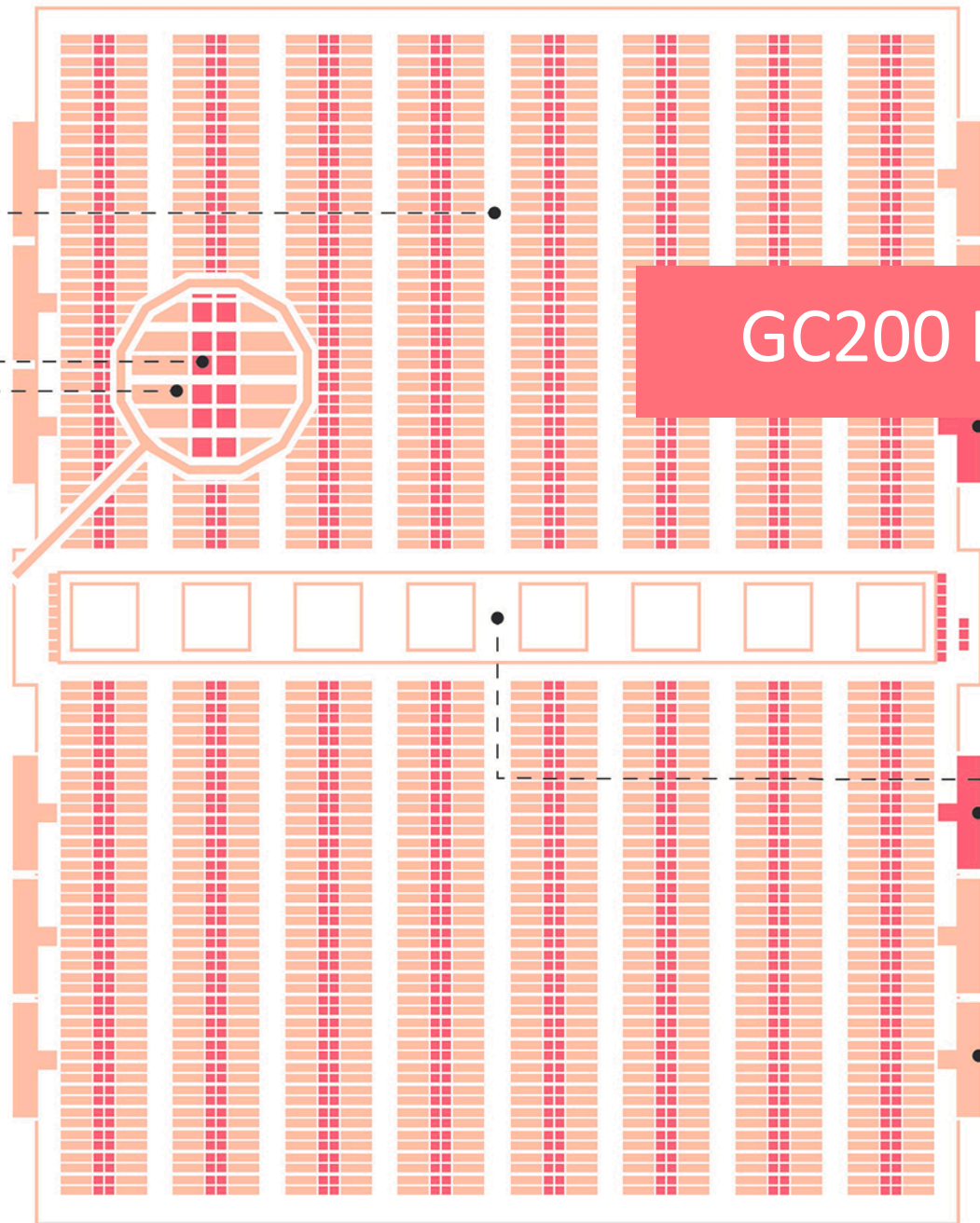
1472 independent IPU-Core™

8832 independent program threads executing in parallel

In-Processor-Memory™

900MB In-Processor-Memory™ per IPU

47.5TB/s memory bandwidth per IPU



GC200 IPU PROCESSOR

IPU-Exchange™

8 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



IPU SCALEOUT: INNOVATIONS & DIFFERENTIATORS

DESIGNED FOR AI

MIMD ARCHITECTURE

ULTRA-FAST
IN-PROCESSOR MEMORY

TRAINING & INFERENCE

1U FLEXIBILITY &
SCALABILITY

HIGH CAPACITY
EXCHANGE MEMORY

IPU-FABRIC BUILT-IN
FOR SCALEOUT

HIGHEST PERFORMANCE
PLATFORM IN 5U FORMFACTOR

FULLY RECONFIGURABLE
FOR SCALEOUT

DISAGGREGATED HOST SERVER
FOR OPTIMAL TCO

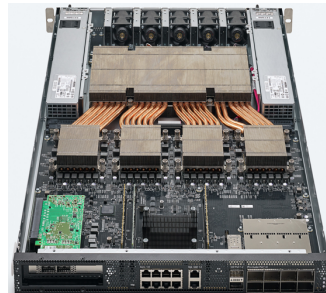
BUILDING BLOCK FOR SCALEOUT

IPU-FABRIC
COMMUNICATION FOR
EFFICIENT
SCALEOUT
PERFORMANCE



x4

GC200



x4

IPU-M2000

4 IPU
1 PetaFlop



IPU-POD16 Direct Attach

x4

IPU-POD₁₆

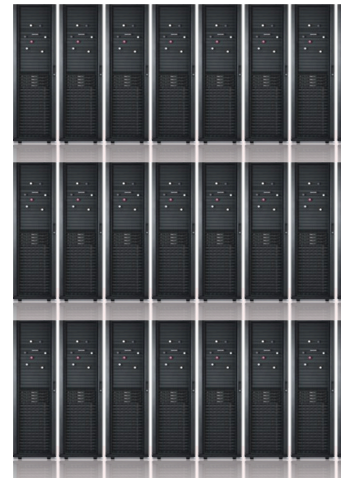
16 IPU
4 PetaFlop
4x IPU-M2000



x1024

IPU-POD₆₄

64 IPU
16 PetaFlop
16x IPU-M2000



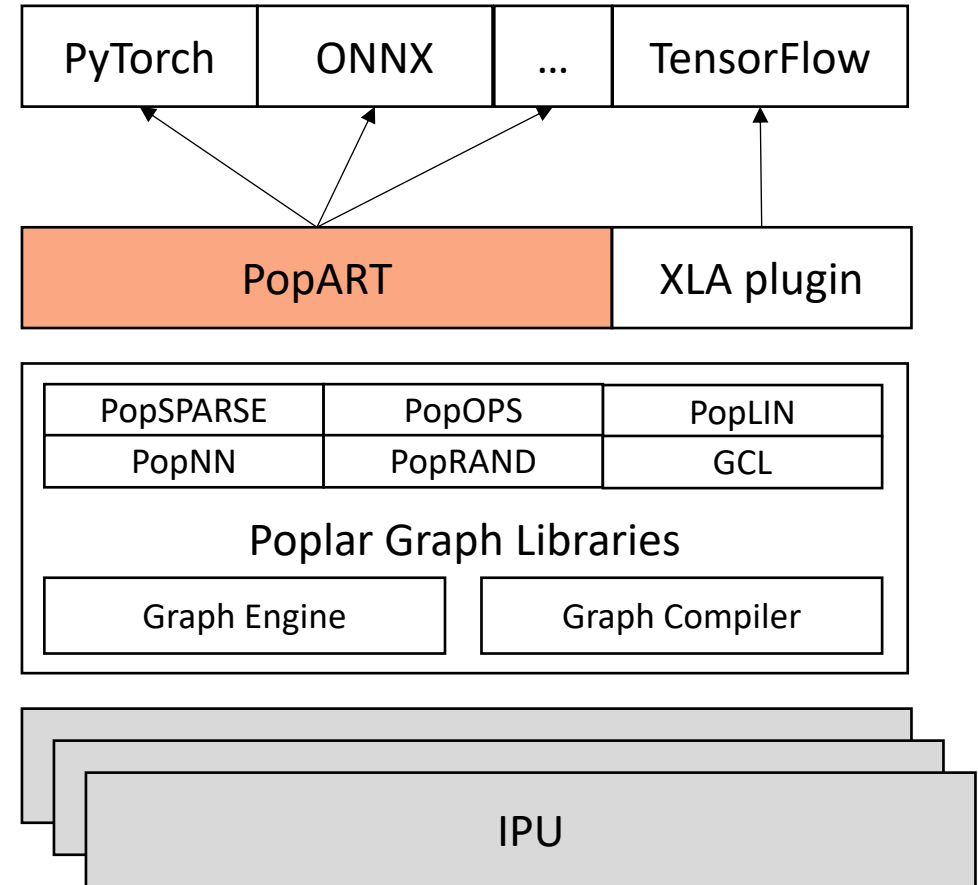
IPU-POD_{64k}

64k IPU
16 ExaFlop
<7.4 PB



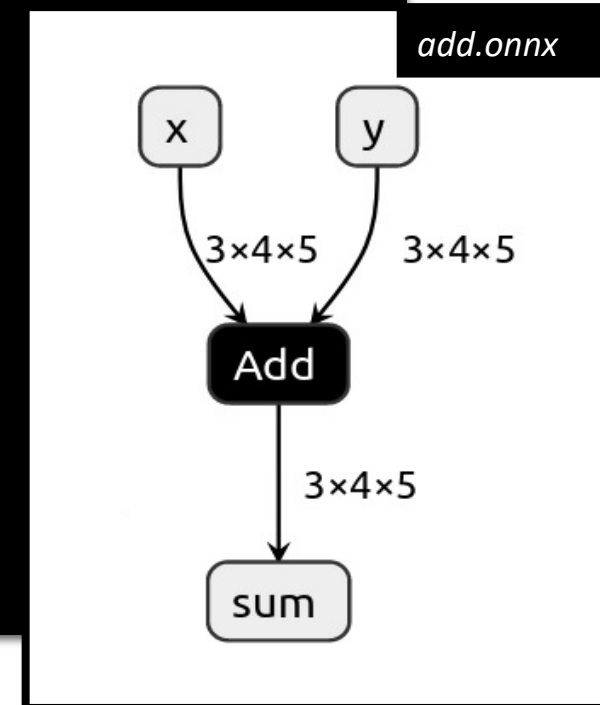
PopART

- PopART : Poplar Advanced Runtime
- PopART is a runtime which can execute ONNX models on IPU(s)
- Implements ONNX IR version 4 , Opset Versions 6 to 11
- Supports inference & training(*)
- Provides a C++ & Python API
- Provides a builder API to construct ONNX models programmatically

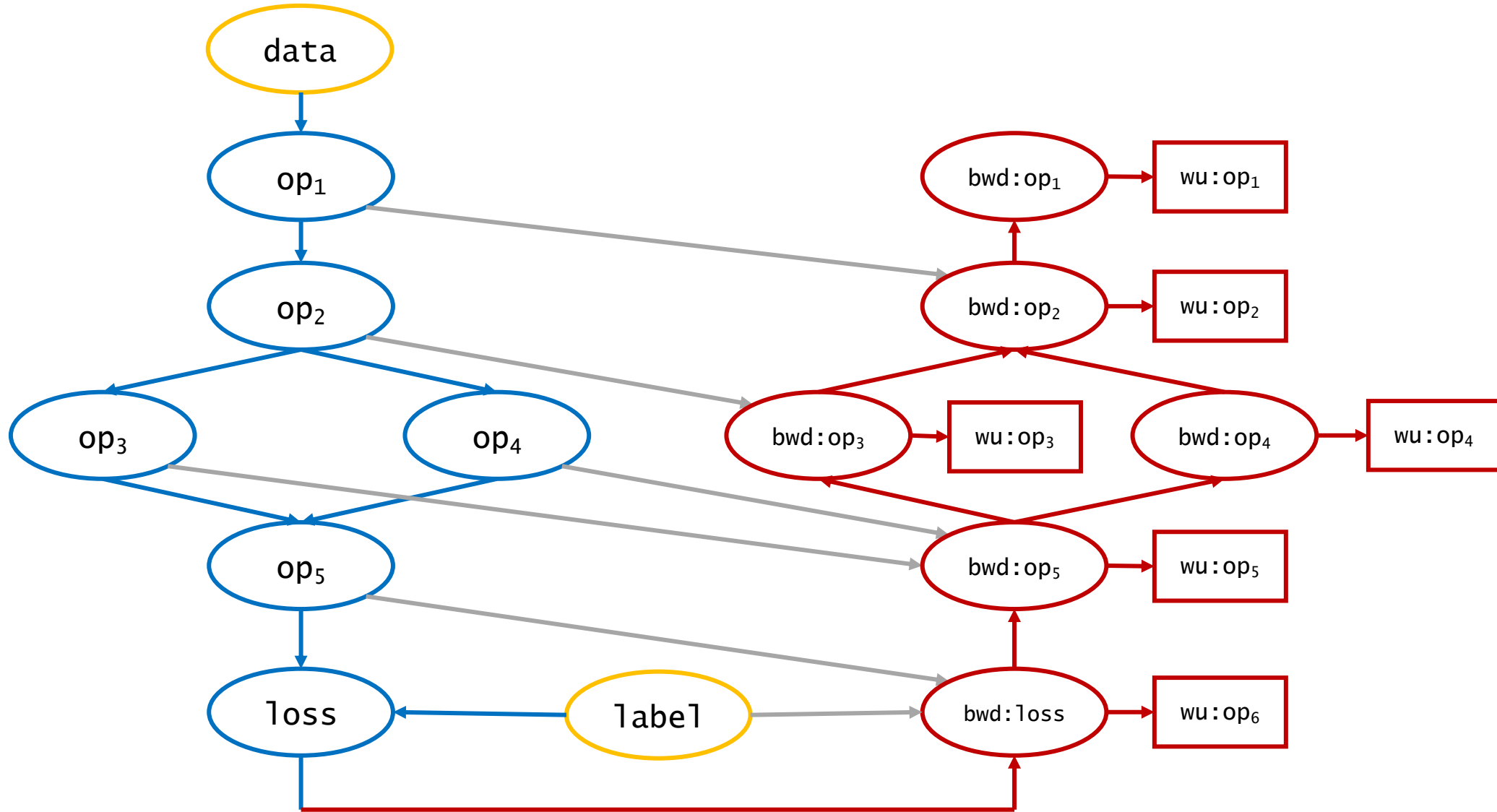


PopART INFERENCE example

```
01 import numpy as np
02 from popart import *
03
04
05 # Create inference session for "add.onnx"
06 session = InferenceSession(fnModel="add.onnx",
07                             dataFeed=DataFlow(batchesPerStep=1, {"sum": AnchorReturnTypes("FINAL")}),
08                             deviceInfo=DeviceManager().acquireAvailableDevice(numIpus=1))
09
10 # Create buffers to receive the output data from the execution of the model
11 anchors = session.initAnchorArrays()
12
13 # Compile graph - where the magic happens
14 session.prepareDevice()
15
16 # Create the input data
17 x = np.random.rand(3, 4, 5).astype(np.float32)
18 y = np.random.rand(3, 4, 5).astype(np.float32)
19
20 # Run the inference session. The PyStepIO defines the inputs and outputs for the run
21 session.run(PyStepIO({"x": x "y": y}, anchors))
22
23 # Print the result
24 print(anchors["sum"])
25
```



PopART TRAINING example

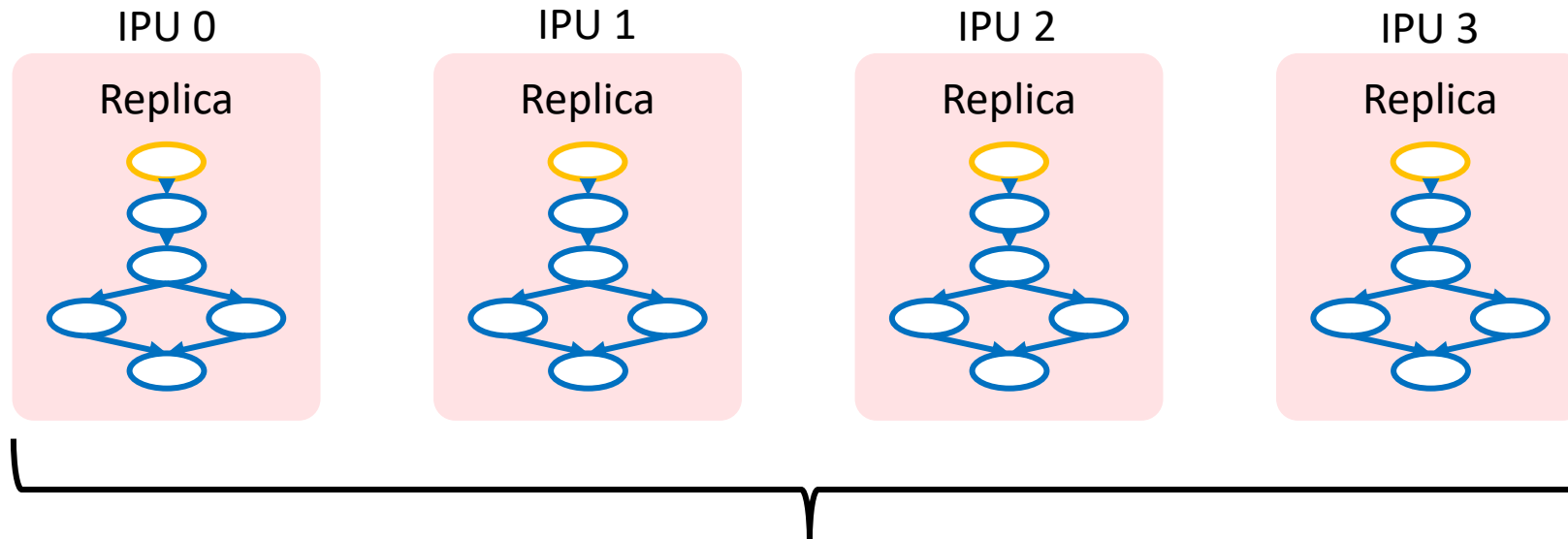


PopART TRAINING example

```
01 import numpy as np
02 from popart import *
03
04
05 builder = Builder("model.onnx")
06
07 # Add a softmax layer to the output of the model & add a label input
08 probs = builder.aiOnnx.softmax(["output"])
09 label = builder.addInputTensor(TensorInfo("INT32", [1]), "label")
10
11 # Create training session (batchesPerStep = 32)
12 session = TrainingSession(fnModel=builder.getModelProto(),
13                           dataFeed=DataFlow(batchesPerStep=32, {"loss": AnchorReturnType("ALL")}),
14                           losses=[NllLoss("output", label, "loss")],
15                           optimizer=SGD(learning_rate=0.001),
16                           deviceInfo=DeviceManager().acquireAvailableDevice(numIpus=1))
17
18 # Create buffers to receive results from the execution
19 anchors = session.initAnchorArrays()
20
21 # Compile graph
22 session.prepareDevice()
23
24 # Write the weights & optimizer to the device
25 session.weightsFromHost()
26 session.optimizerFromHost()
27
28 # Run the training session
29 for i in range(8):
30     << Read the data (image_data [32, 4, 224, 224] & label_data [32, 4, 1]) >>
36     session.run(PyStepIO({"input": image_data, "label": label_data }, anchors))
37
38 # Write the trained model to file
39 session.modelToHost("trained_model.onnx")
```


DATA PARALLELISM: REPLICATION

```
options = SessionOptions()  
options.enableReplicatedGraphs = True  
options.replicationFactor = 4  
session = TrainingSession(fnModel=builder.getModelProto(), userOptions=options, ...)
```



Data parallel scale-out – each replica has a copy of the model and different data

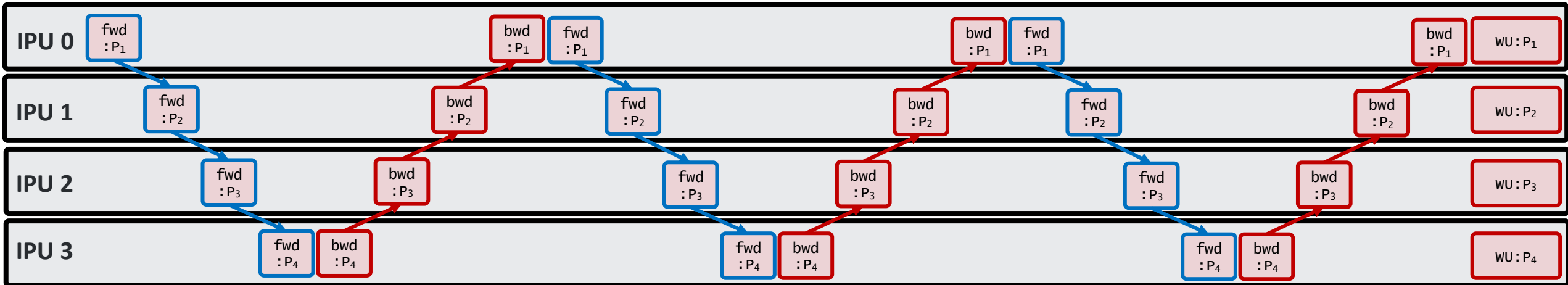
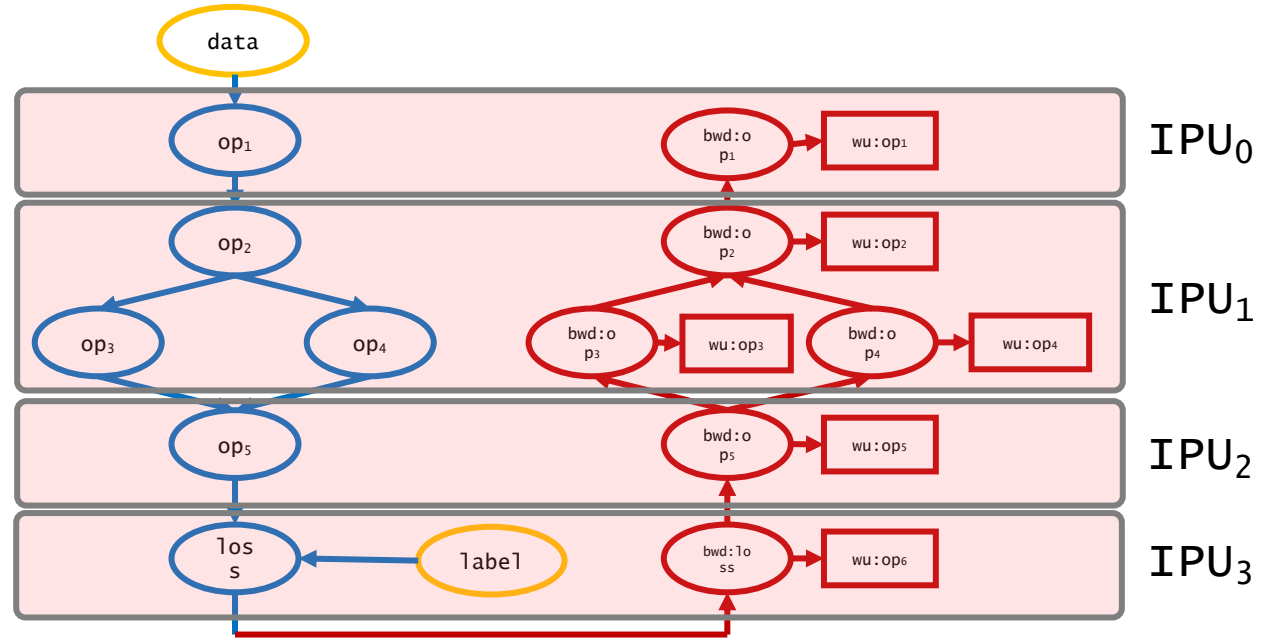
MODEL PARALLELISM: SHARDING LAYERS ACROSS IPUS

Auto sharding:

```
options.virtualGraphMode = VirtualGraphMode::Auto
```

Manual sharding:

```
options.virtualGraphMode = VirtualGraphMode::Manual
with builder.virtualGraph(1):
  o = builder.aiOnnx.add([o1, o2])
```



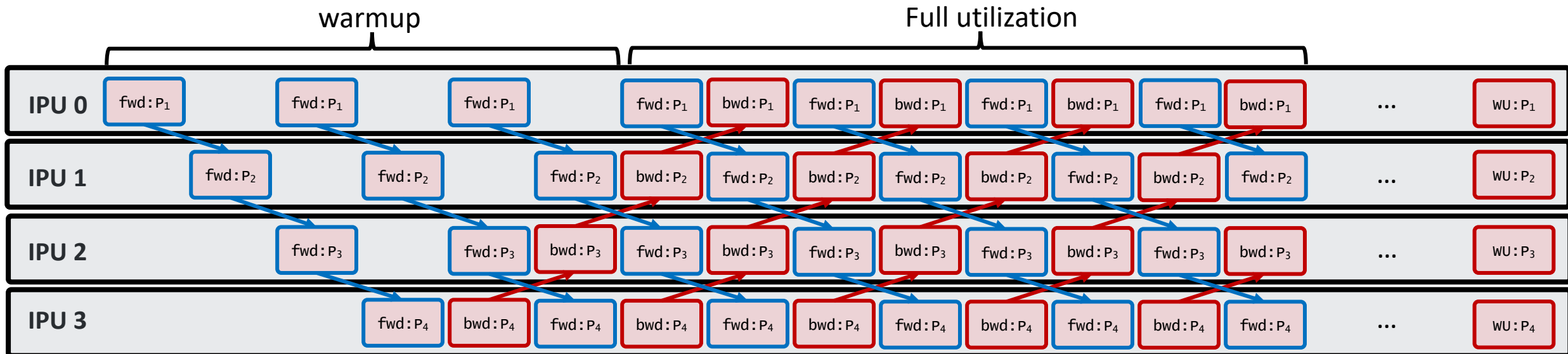
MODEL PARALLELISM: PIPELINING

Enable pipelining:

```
options.virtualGraphMode = VirtualGraphMode::Manual
options.enablePipelining = True

with builder.virtualGraph(0):
    o3 = builder.aiOnnx.add([o1, o2])

with builder.virtualGraph(1):
    o5 = builder.aiOnnx.add([o3, o4])
```



time →

THANK YOU

www.graphcore.ai

www.graphcore.cn

hanzhao@graphcore.ai / info_china@graphcore.ai

