

# IPU HANDS-ON SESSION

# GRAPHCORE



Reinforcement Learning Neural network visualization from [POPLAR™](#)



# Day03

## IPU optimization



# OVERVIEW

## GOAL

Start with a configuration where  
Out Of Memory error occurs

Step-by-step optimization

Successfully run the same model  
with same batch size  
with high performance

- File Structure
  - main.py : main training file
  - utils.py : contains functions/classes used in main
  - requirements.txt : required packages  
(pip install -r requirements.txt)

- ResNet101 pipelined over 2 IPUs

```
model = torchvision.models.resnet101()  
model.layer3[8] = poptorch.BeginBlock(model.layer3[8], ipu_id=1)
```

- CIFAR10 dataset

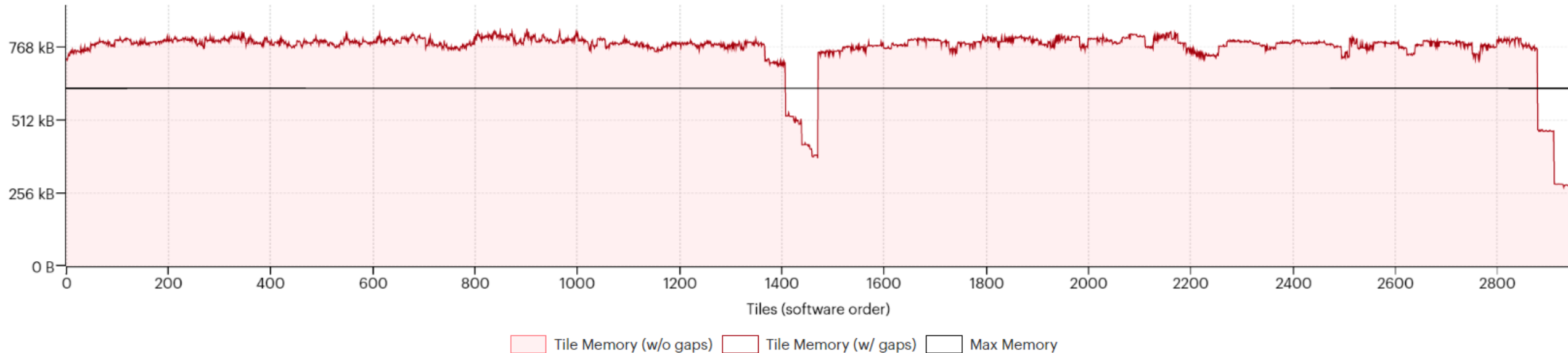
```
train_dataset = torchvision.datasets.CIFAR10(root='../datasets', \br/>| train=True, download=True, transform=transform)
```

# BASE

- Configs where OOM occurs

```
python main.py --batch-size 64 --precision 32.32 --gradient-accumulation 3 --epochs 1
```

```
Namespace(async_data_loader=False, batch_size=64, device_iterations=1,  
| | | | | eight_bit_io=False, epochs=1, gradient_accumulation=3, precision='32.32', replicas=1)
```





# HALF PRECISION

- Graph analyzer provides some tips for reducing memory usage

Summary

Insights

Memory Report

Liveness Report

Program Tree

### Reducing memory usage

You can reduce the memory requirements of your application in the following ways. More information on these can be found in the [Graphcore Memory and Performance Optimisation Guide](#)

These solutions may affect the performance, throughput, convergence and/or training characteristics of your model

Reducing the batch size

Using FP16

Using recomputation

Using recomputation checkpoints

Setting the available memory proportion

Offloading variables

Reusing identical parts of your graph

Disable profiling on the IPU

+ Using different optimizers

Top priority is to use FP16!

## Using FP16 where appropriate

### What is FP16?

Using FP16 instead of FP32 will reduce the memory required for a variable by half and can double the execution speed of operations that use them. For a model you can choose if you want to use FP16 for partials, compute and weights. FP16 may not be suitable in all cases due to a reduced numerical range. Techniques to mitigate this are described here.

### How changing from FP32 to FP16 can fix OOM

Reducing the **partials** format also improves the performance (See Table 5.1 in the [AI-Float™ white paper](#)). So, each time the numerical precision is halved, the execution speed of matmul and convolution operations is doubled.

### Useful links

[Memory and Performance Optimisation on the IPU](#)

# HALF PRECISION

Different precision formats during model execution

- Single precision : FP32 for both the data and the model
- Mixed precision : FP16 for the data and FP32 for the model
- Half precision : FP16 for both the data and the model

2 ways to set half precision

- `model = model.to(torch.float16)`
- `model = model.half()`

```
input_dtype = torch.float16 if args.precision[:2] == '16' else torch.float32
model_dtype = torch.float16 if args.precision[-2:] == '16' else torch.float32

if args.eight_bit_io:
    model = ModelWithNormalization(model, dtype=input_dtype)
model = ModelWithLoss(model, criterion).to(model_dtype)
```

```
class ModelWithNormalization(torch.nn.Module):
    def __init__(self, model: torch.nn.Module, dtype: torch.dtype):
        super().__init__()
        self.model = model
        self.dtype = dtype
        self.normalize = transforms.Normalize(**IMAGENET_STATISTICS)

    def forward(self, img):
        return self.model(self.normalize(img.to(self.dtype)))
```

# HALF PRECISION

Poplar SDK contains many tools to support stable training with half precision, so it's highly recommended to use half precision on the IPU

```
class AdamW(Optimizer, torch.optim.AdamW):
    """ Adam optimizer with true weight decay.

    This optimizer matches PyTorch's implementation
    (torch.optim.AdamW https://pytorch.org/docs/1.10.0/optim.html#torch.optim.AdamW)
    with optional loss scaling.

    AMSGrad is currently not supported."""

    # Variables which don't exist in the parent optimizer class and are
    # global (Cannot be set per group).
    _child_vars = ["loss_scaling"]
    # All the attributes and variables which don't exist in the parent optimizer class.
    _child_only = _child_vars + [
        "bias_correction",
        "accum_type",
        "first_order_momentum_accum_type",
        "second_order_momentum_accum_type",
        "max_grad_norm",
    ]
```

```
if args.precision[-3:] == '.16':
    optimizer = poptorch.optim.AdamW(model.parameters(), lr=0.001,
                                     loss_scaling=1000,
                                     accum_type=torch.float16,
                                     first_order_momentum_accum_type=torch.float16,
                                     second_order_momentum_accum_type=torch.float32)
```

- Poptorch optimizers are implemented by inheriting the original torch.optim
- Provides more flexible control over data types during execution

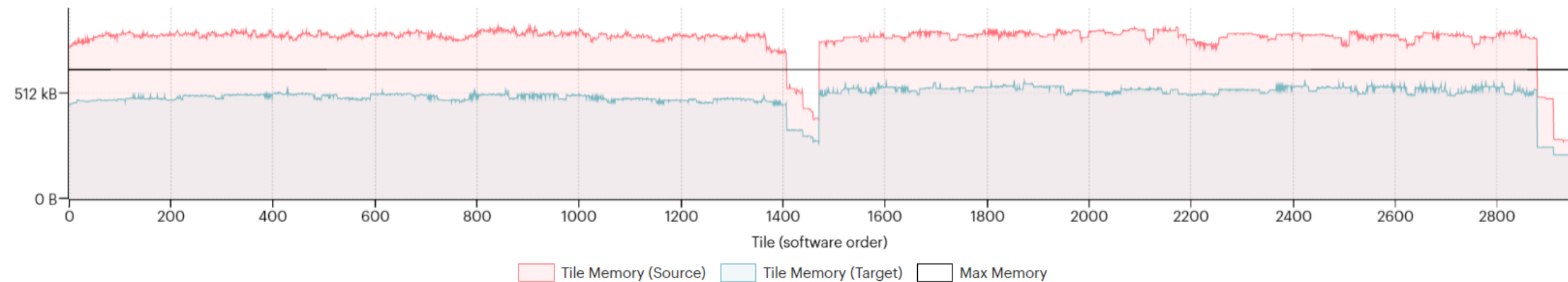
- Factor by which to scale the loss
- Automatic Loss Scaling is supported for SDK>=2.5

# HALF PRECISION

- Using half precision

```
python main.py --batch-size 64 --precision 16.16 --gradient-accumulation 3
```

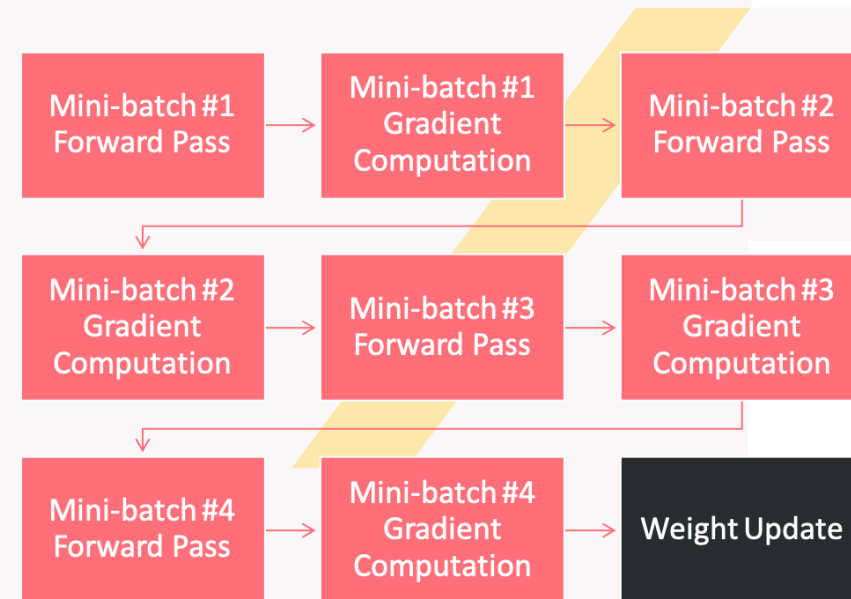
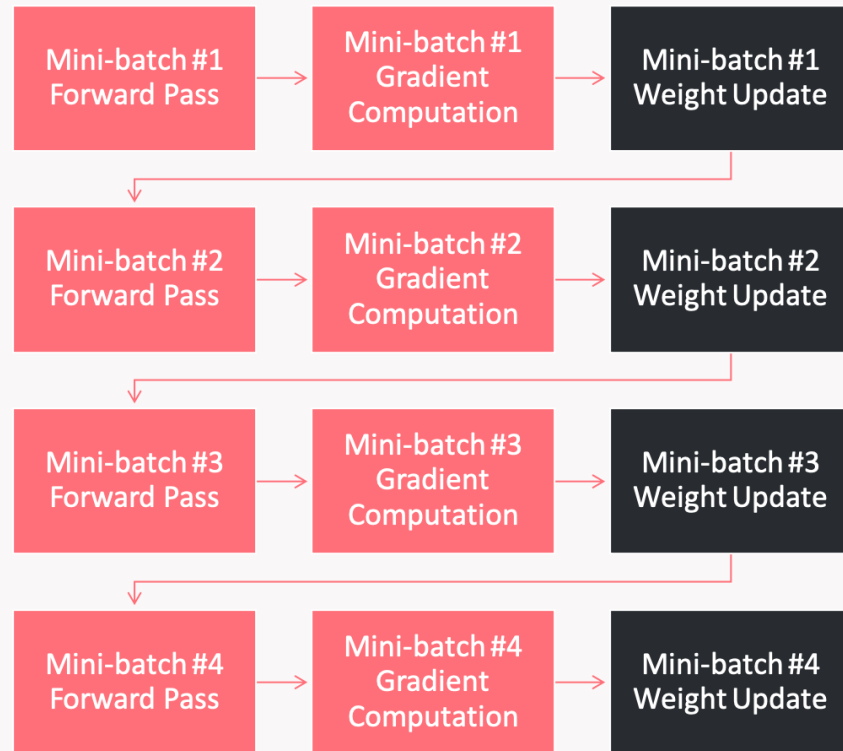
```
Namespace(async_dataloader=False, batch_size=64, device_iterations=1,  
| | | | | eight_bit_io=False, epochs=5, gradient_accumulation=3, precision='16.16', replicas=1)
```





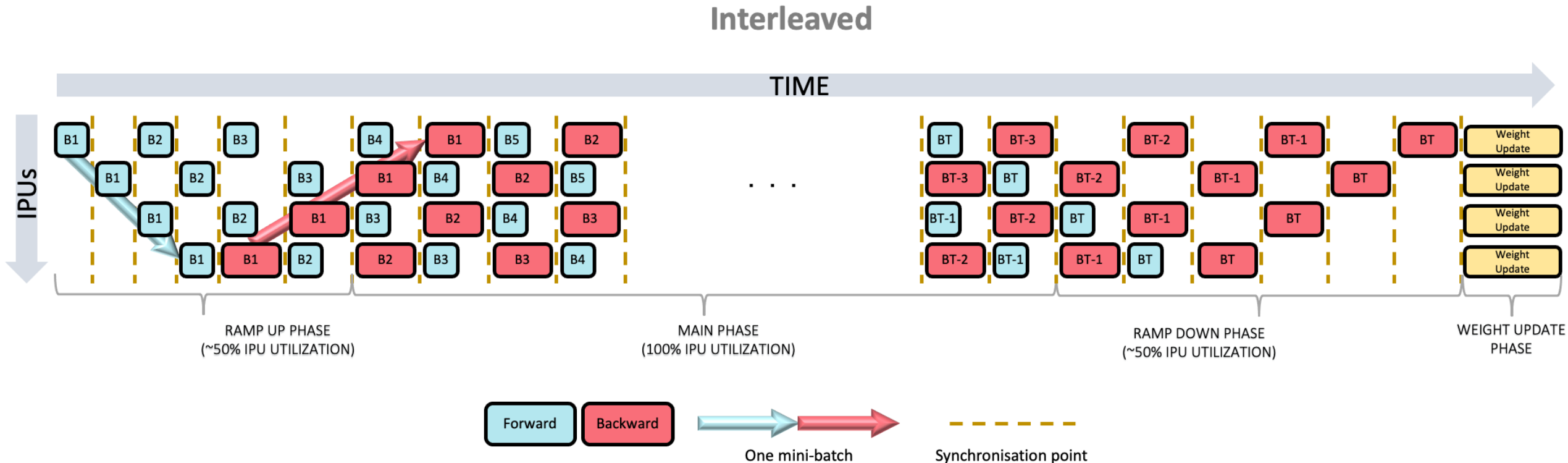
# GRADIENT ACCUMULATION

- More samples per a weight update without increase in memory usage
- Speed up due to less weight updates



# GRADIENT ACCUMULATION

- Maximize the portion of main phase → higher throughput



# GRADIENT ACCUMULATION

- How to set

```
opts = poptorch.Options()  
opts.Training.gradientAccumulation(args.gradient_accumulation)
```

- Command using higher gradient accumulation factor

```
python main.py --batch-size 64 --precision 16.16 --gradient-accumulation 8
```

- Output

```
[Epoch 01]: 100%|██████████| 97/97 [09:20<00:00, 5.78s/it, loss=1.56, acc=39.1]  
[Epoch 01] loss: 1.900, acc: 33.7  
[Epoch 02]: 100%|██████████| 97/97 [00:46<00:00, 2.10it/s, loss=1.01, acc=67.2]  
[Epoch 02] loss: 1.375, acc: 50.0  
[Epoch 03]: 100%|██████████| 97/97 [00:46<00:00, 2.10it/s, loss=1.31, acc=54.7]  
[Epoch 03] loss: 1.177, acc: 58.0  
[Epoch 04]: 100%|██████████| 97/97 [00:46<00:00, 2.11it/s, loss=0.935, acc=67.2]  
[Epoch 04] loss: 1.056, acc: 61.9  
[Epoch 05]: 100%|██████████| 97/97 [00:46<00:00, 2.09it/s, loss=0.905, acc=64.1]  
[Epoch 05] loss: 0.921, acc: 67.3
```

# GRADIENT ACCUMULATION

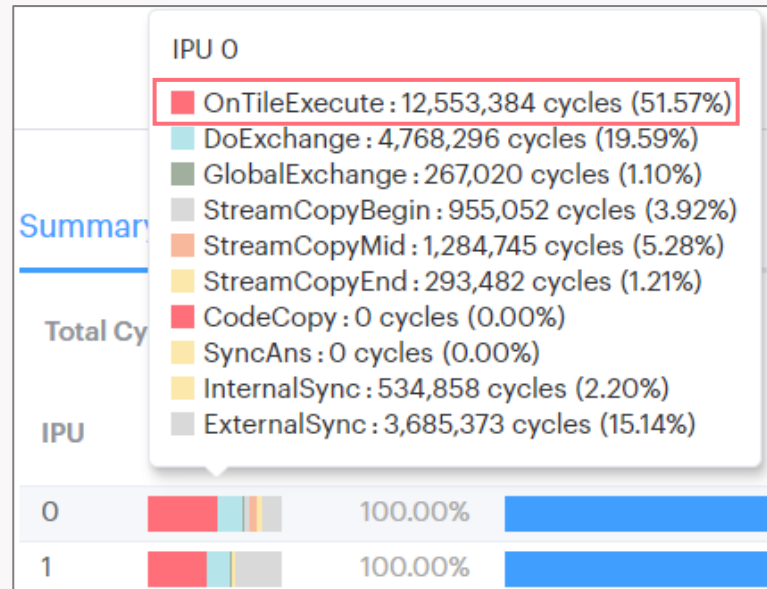
## Execution Trace

Summary

Details

Total Cycles 60,111,417

IPU	Proportions	Tile Balance	SteamCopy (Rx/Tx)	G
0		100.00%	4.5 MB / 51.8 kB	
1		100.00%	768 B / 141 kB	



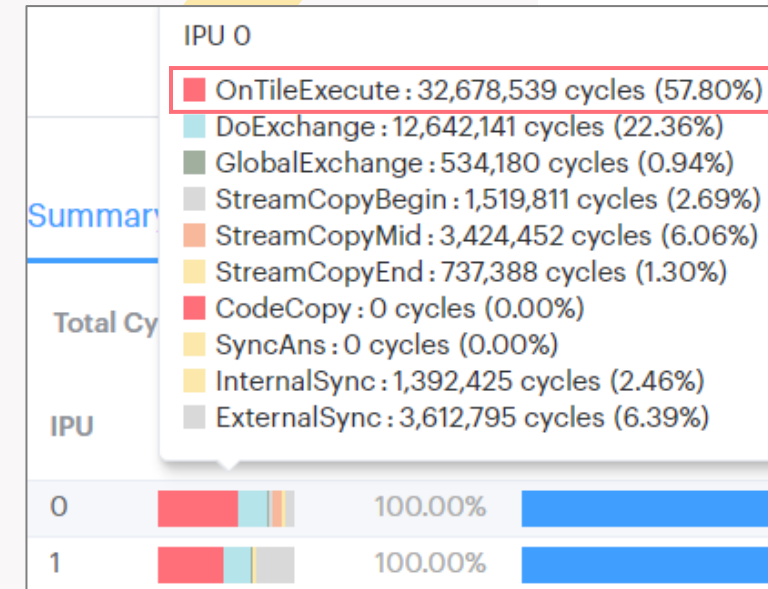
The portion of OnTileExecute became larger

Summary

Details

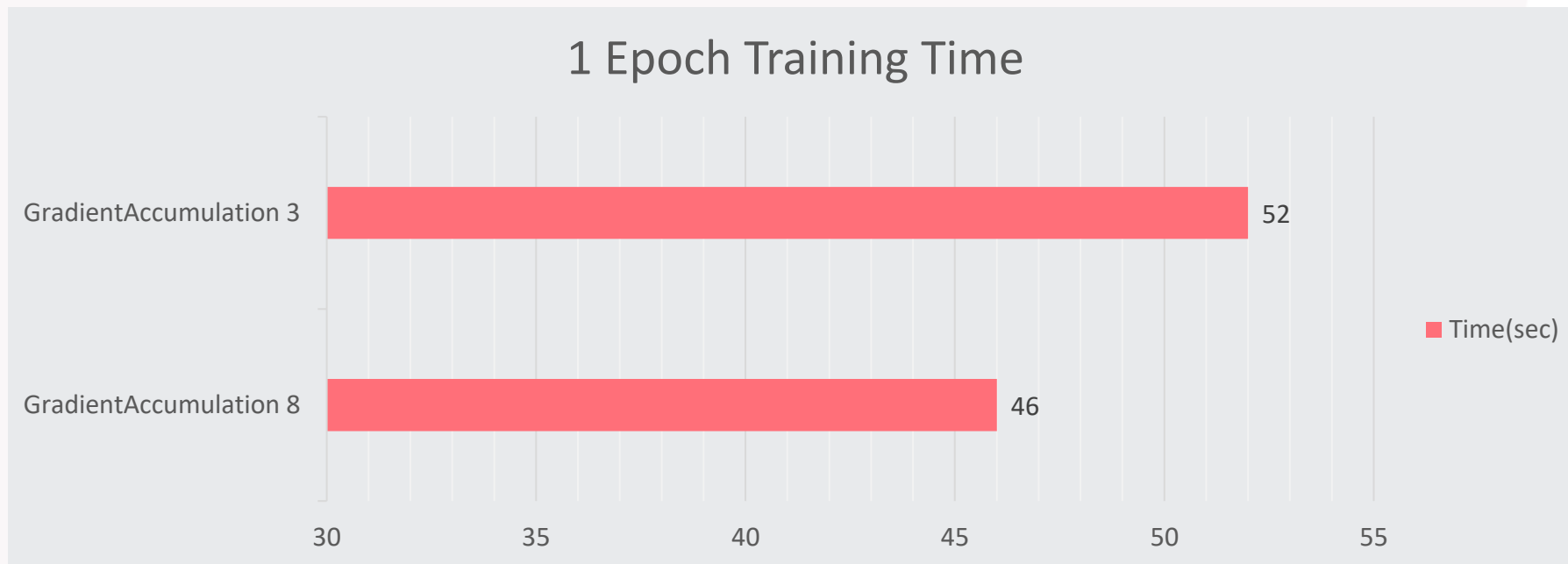
Total Cycles 95,653,315

IPU	Proportions	Tile Balance	SteamCopy (Rx/Tx)	G
0		100.00%	12.0 MB / 138 kB	
1		100.00%	2.0 kB / 141 kB	





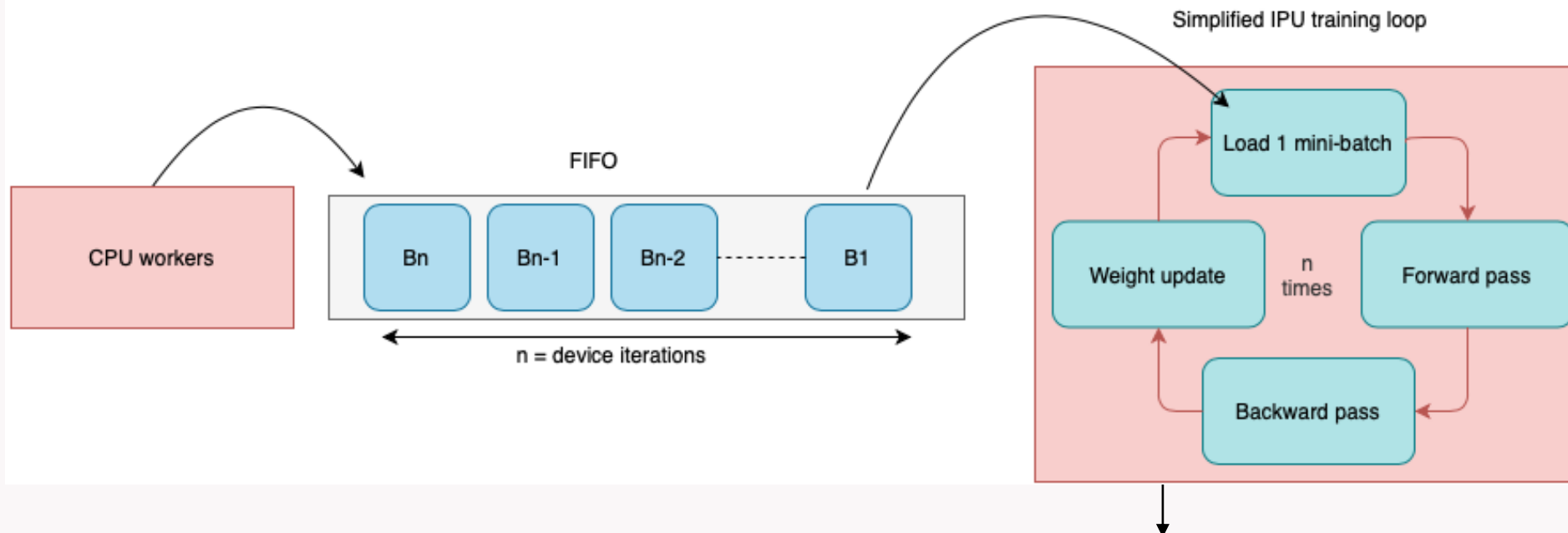
# GRADIENT ACCUMULATION



# DEVICE ITERATION

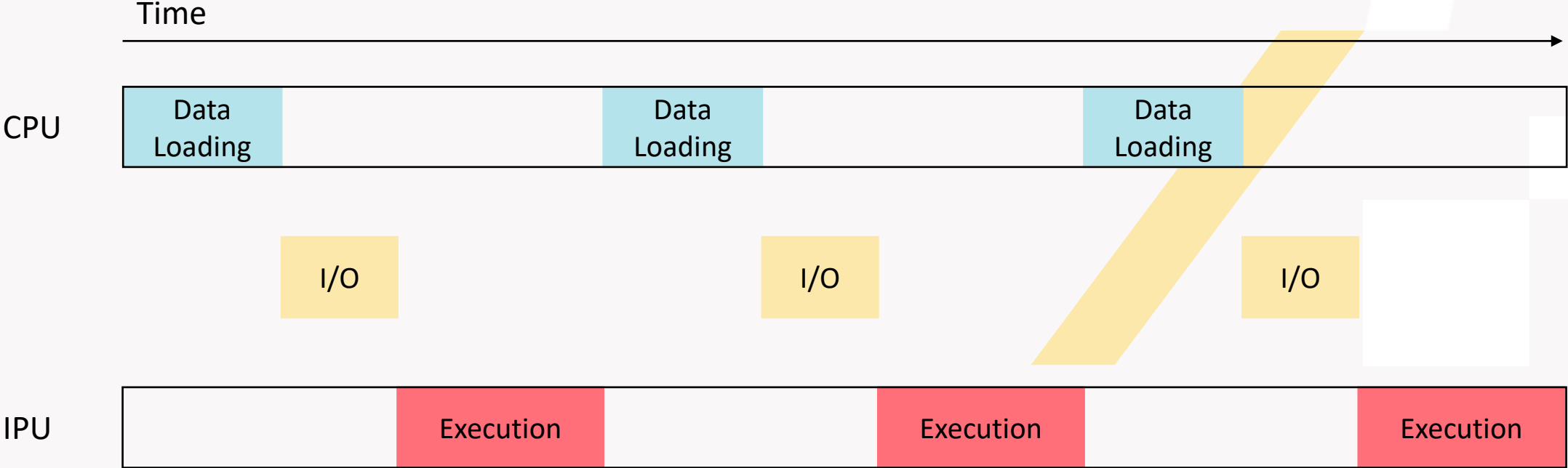
- Iterate training loops over multiples batches from an input queue
- Accelerate training by reducing the number of host-device communications

Basic training: no pipeline, no replication

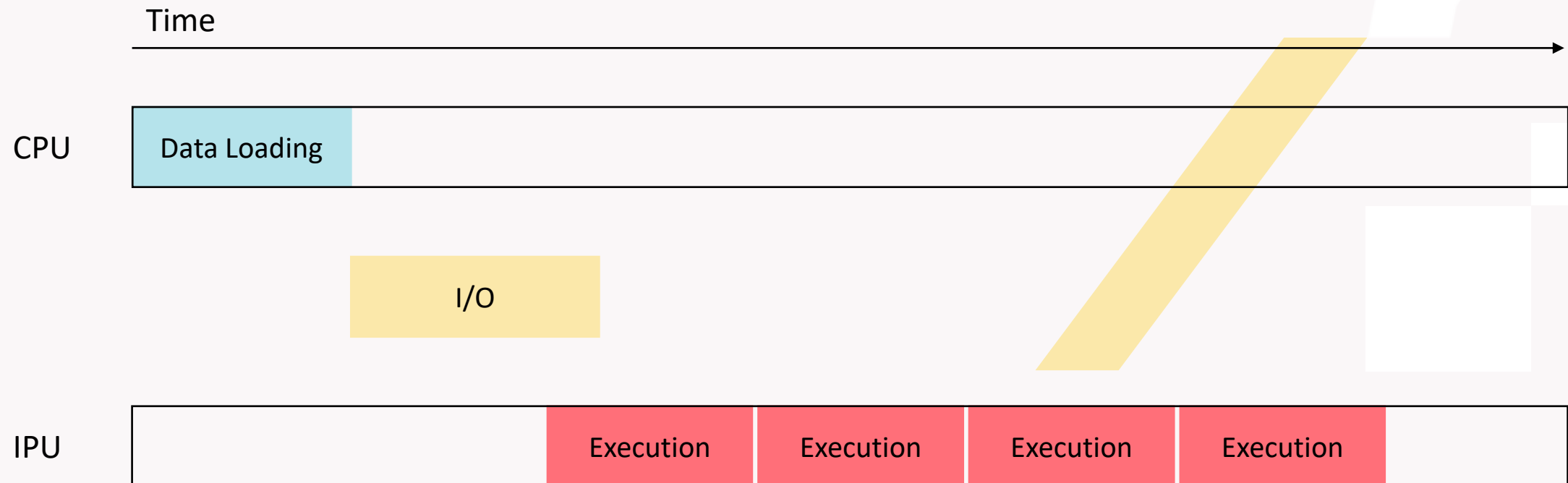


- `output_mode` (*`poptorch.OutputMode`*) –
  - `All`: Return a result for each batch.
  - `Sum`: Return the sum of all the batches.
  - `Final`: Return the last batch.
  - `EveryN`: Return every N batches: N is passed in as `output_return_period`.
  - Default: `All` for inference, `Final` for training.

# DEVICE ITERATION



# DEVICE ITERATION





# DEVICE ITERATION

- How to set

```
opts = poptorch.Options()
opts.deviceIterations(args.device_iterations)
```

- Command using device iteration

```
python main.py --batch-size 64 --precision 16.16 --gradient-accumulation 8 --device-iteration 4
```

- Output

```
[Epoch 01]: 100%|██████████| 24/24 [09:11<00:00, 22.97s/it, loss=1.51, acc=45.3]
[Epoch 01] loss: 1.759, acc: 36.5
[Epoch 02]: 100%|██████████| 24/24 [00:45<00:00, 1.88s/it, loss=1.11, acc=57.8]
[Epoch 02] loss: 1.395, acc: 48.3
[Epoch 03]: 100%|██████████| 24/24 [00:45<00:00, 1.88s/it, loss=0.841, acc=75]
[Epoch 03] loss: 1.251, acc: 55.1
[Epoch 04]: 100%|██████████| 24/24 [00:45<00:00, 1.88s/it, loss=1.37, acc=54.7]
[Epoch 04] loss: 1.086, acc: 61.8
[Epoch 05]: 100%|██████████| 24/24 [00:45<00:00, 1.88s/it, loss=0.744, acc=73.4]
[Epoch 05] loss: 0.863, acc: 68.9
```

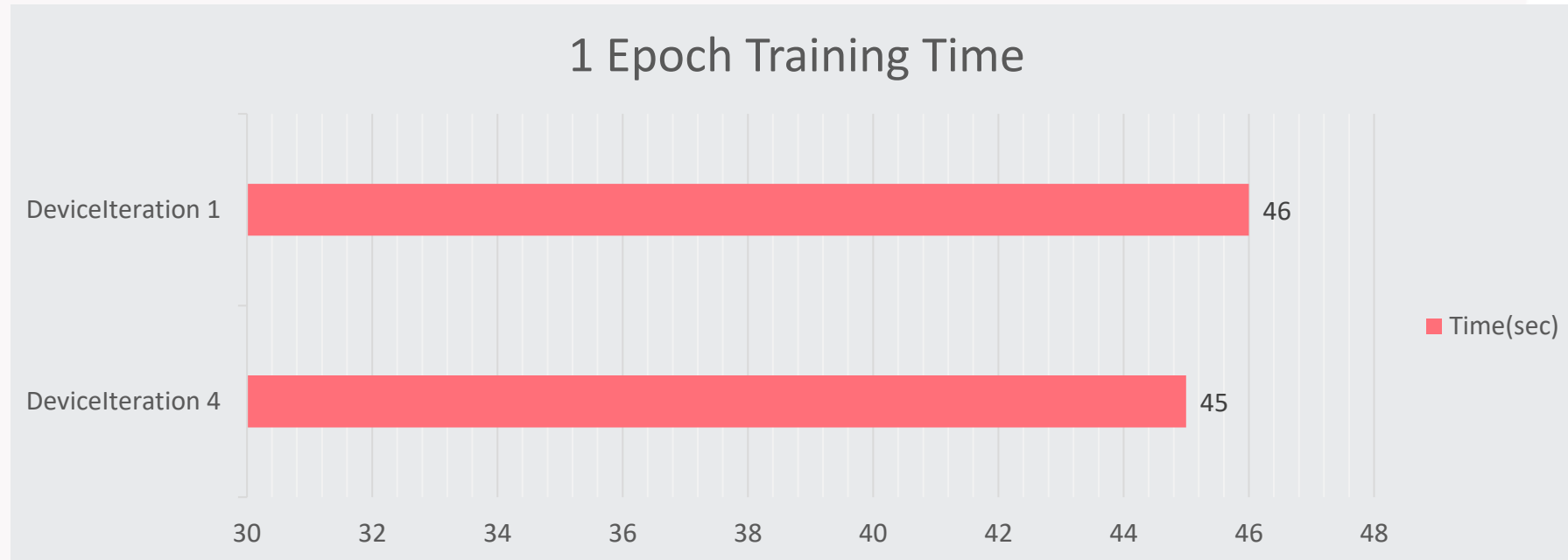
# DEVICE ITERATION

Execution Trace



Iterating 4 training loops on device

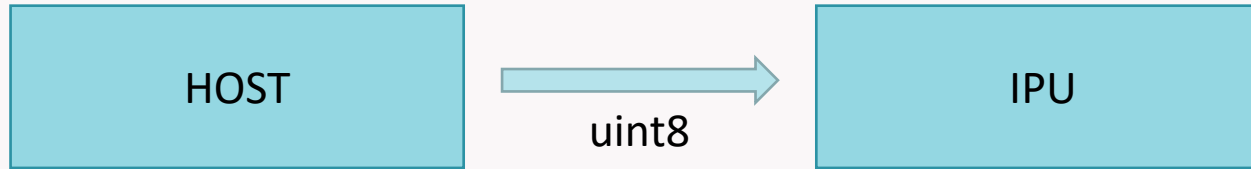
# DEVICE ITERATION



Host overhead might be occurred while processing the larger data

# 8-BIT I/O

- Casting the input data to uint8 can reduce I/O overhead



HOST side

```
def get_transform(precision='16.16', eight_bit_io=False):
    preprocssing_steps = [
        transforms.RandomHorizontalFlip(),
        transforms.Resize((64, 64)),
        transforms.PILToTensor()
    ]

    if not eight_bit_io:
        dtype = torch.float16 if precision[:2] == '16' else torch.float32
        preprocssing_steps.extend([
            transforms.ConvertImageDtype(dtype),
            transforms.Normalize(**IMAGENET_STATISTICS)
        ])

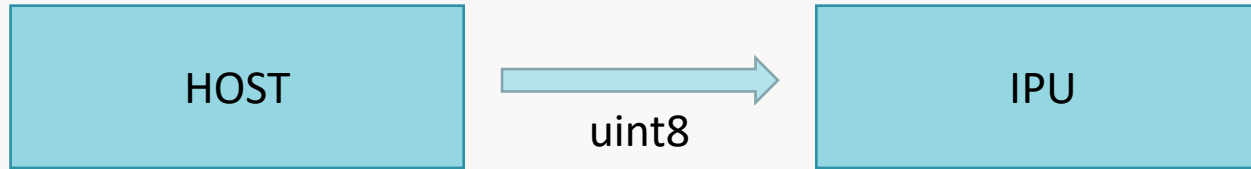
    return transforms.Compose(preprocssing_steps)
```

In usual case (not using 8-bit I/O),  
type casting and normalization conducted on CPU



# 8-BIT I/O

- Casting the input data to uint8 can reduce I/O overhead



IPU side

```
if args.eight_bit_io:
    model = ModelWithNormalization(model, dtype=input_dtype)

class ModelWithNormalization(torch.nn.Module):
    def __init__(self, model: torch.nn.Module, dtype: torch.dtype):
        super().__init__()
        self.model = model
        self.dtype = dtype
        self.normalize = transforms.Normalize(**IMAGENET_STATISTICS)

    def forward(self, img):
        # One-liner in Poplar SDK >= 2.6:
        # return self.model(self.normalize(img.to(self.dtype)))
        if self.dtype == torch.float32:
            img = img.float()
        elif self.dtype == torch.float16:
            img = img.half()
        else:
            raise Exception(f'Normalization dtype must be one of torch.float16 or torch.float32,'
                             f' but {self.dtype} is given')
        return self.model(self.normalize(img))
```

When using 8-bit I/O,  
type casting and normalization conducted on IPU

# 8-BIT I/O

- Command using 8-bit I/O

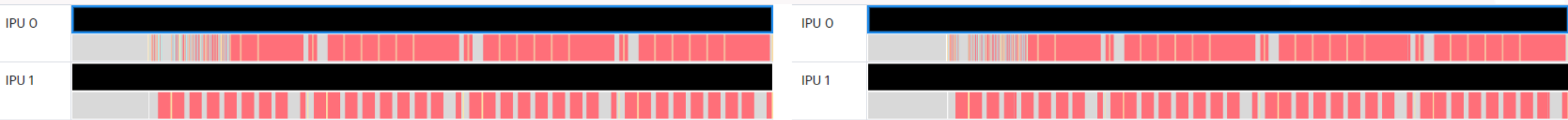
```
python main.py --batch-size 64 --precision 16.16 --gradient-accumulation 8 \
    --device-iteration 4 --eight-bit-io
```

- Output

```
[Epoch 01]: 100%|██████████| 24/24 [08:35<00:00, 21.48s/it, loss=1.65, acc=34.4]
[Epoch 01] loss: 1.848, acc: 32.6
[Epoch 02]: 100%|██████████| 24/24 [00:13<00:00, 1.77it/s, loss=1.33, acc=51.6]
[Epoch 02] loss: 1.334, acc: 51.8
[Epoch 03]: 100%|██████████| 24/24 [00:13<00:00, 1.77it/s, loss=1.44, acc=50]
[Epoch 03] loss: 1.205, acc: 56.3
[Epoch 04]: 100%|██████████| 24/24 [00:13<00:00, 1.77it/s, loss=0.639, acc=81.2]
[Epoch 04] loss: 1.004, acc: 63.4
[Epoch 05]: 100%|██████████| 24/24 [00:13<00:00, 1.77it/s, loss=0.903, acc=71.9]
[Epoch 05] loss: 0.926, acc: 67.2
```

# 8-BIT I/O

## Execution Trace



Summary

Details

Total Cycles 283,021,420

IPU	Proportions	Tile Balance	SteamCopy (Rx/Tx)	G
0		100.00%	48.0 MB / 552 kB	
1		100.00%	0 kB / 142 kB	

Summary

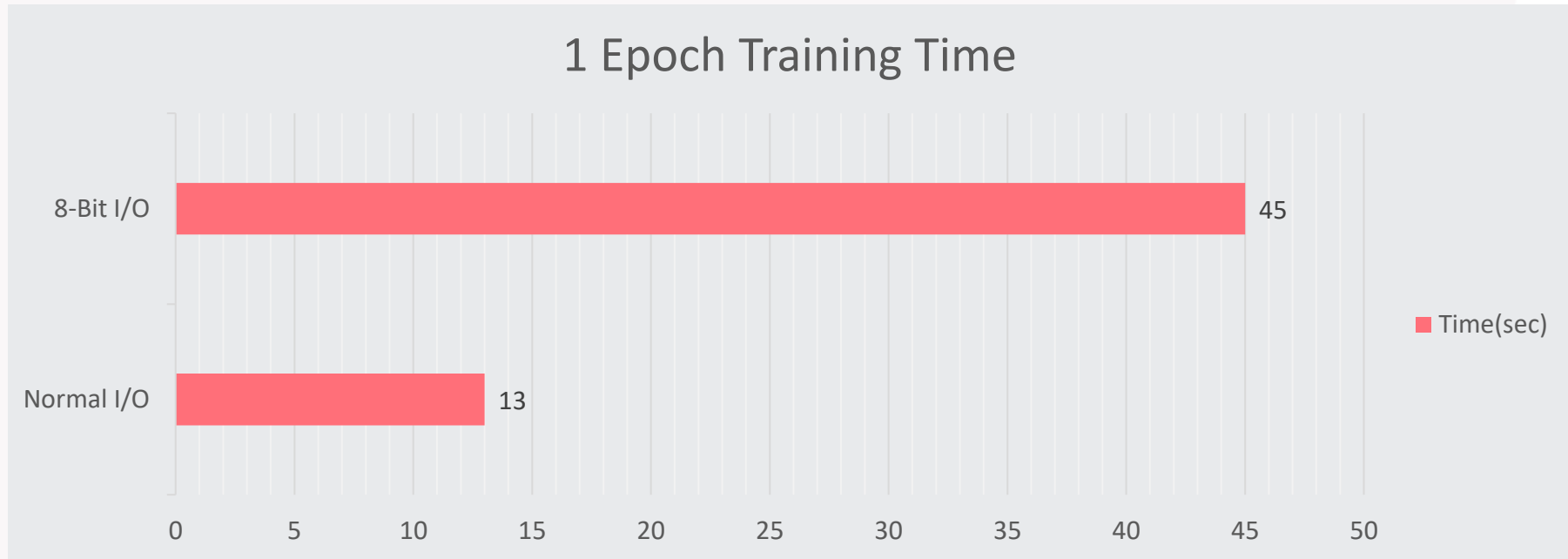
Details

Total Cycles 278,393,266

IPU	Proportions	Tile Balance	SteamCopy (Rx/Tx)
0		100.00%	24.0 MB / 551 kB
1		100.00%	0 kB / 142 kB

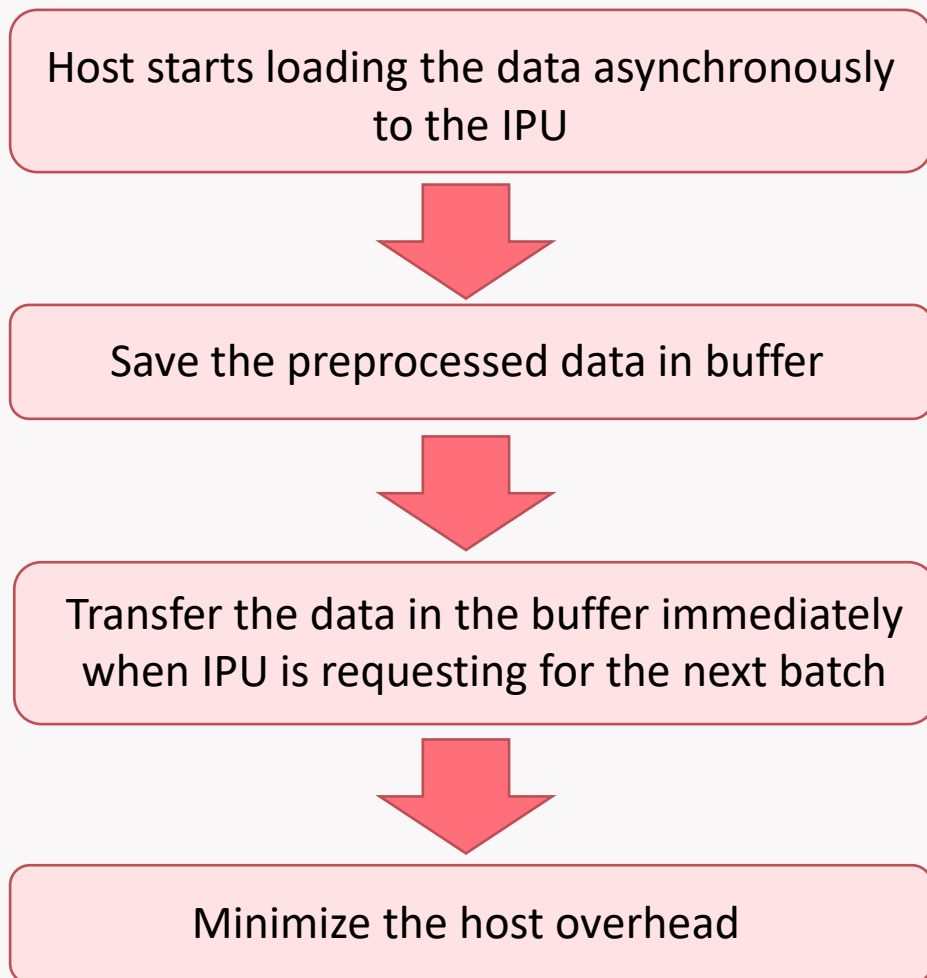
Size of the received data reduced to half

# 8-BIT I/O





# ASYNCHRONOUS DATALOADER

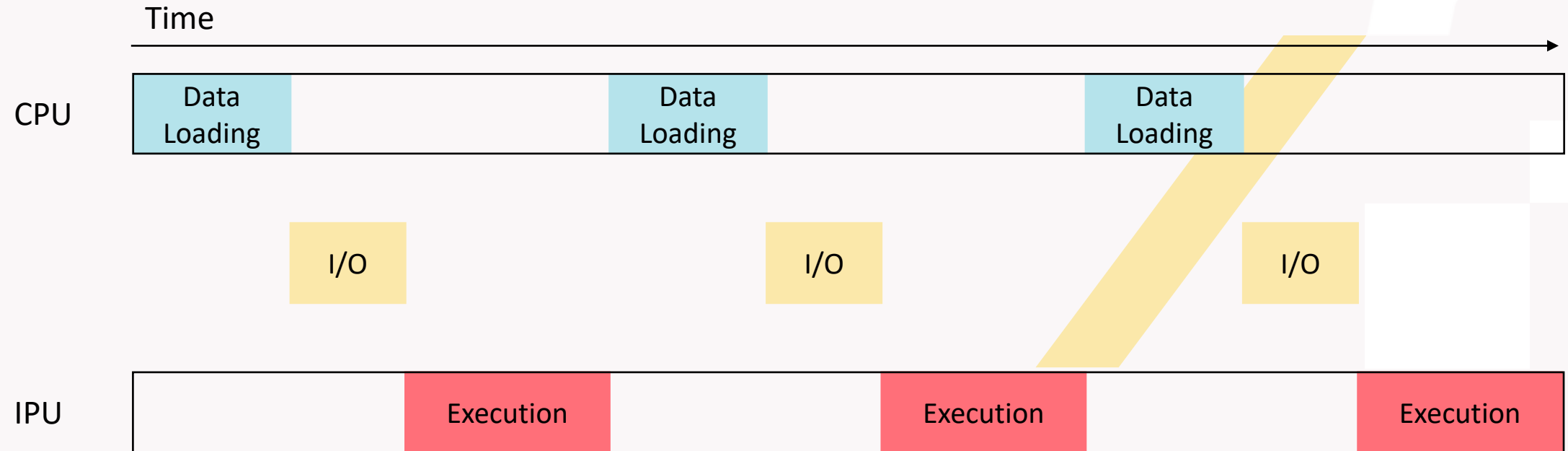


## ▪ How to set

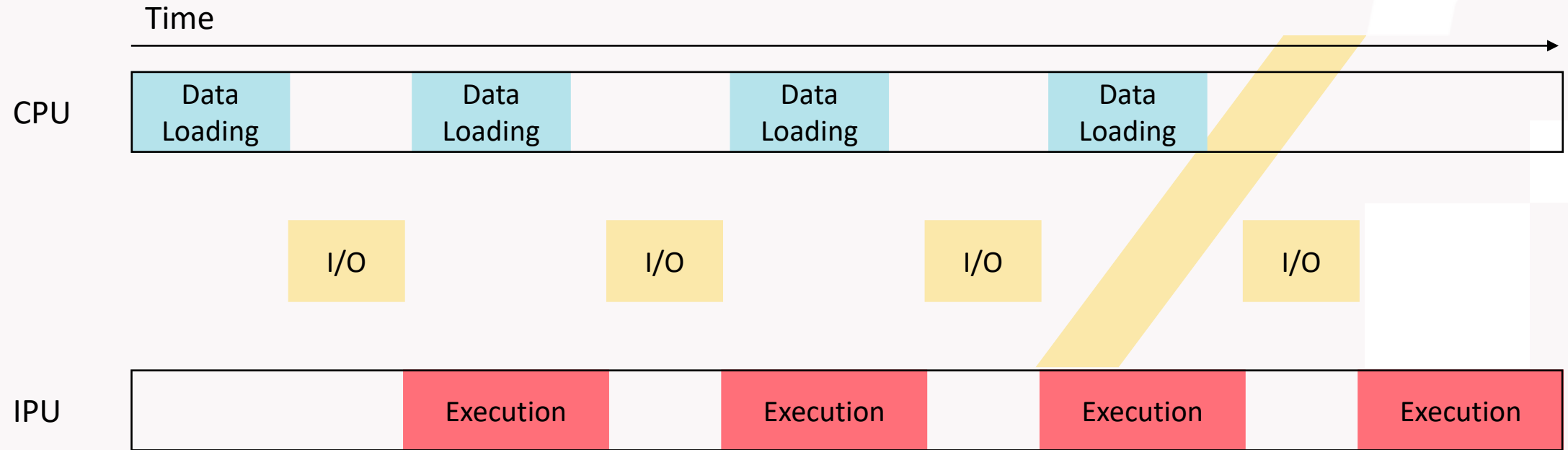
```
mode = poptorch.DataLoaderMode.Async if args.async_dataloader \
| else poptorch.DataLoaderMode.Sync

train_dataloader = poptorch.DataLoader(opts,
|                                     train_dataset,
|                                     batch_size=args.batch_size,
|                                     shuffle=True,
|                                     mode=mode)
```

# ASYNCHRONOUS DATALOADER



# ASYNCHRONOUS DATALOADER



# ASYNCHRONOUS DATALOADER

- Command using asynchronous dataloader

```
python main.py --batch-size 64 --precision 16.16 --gradient-accumulation 8 \
    --device-iteration 4 --eight-bit-io --async-dataloader
```

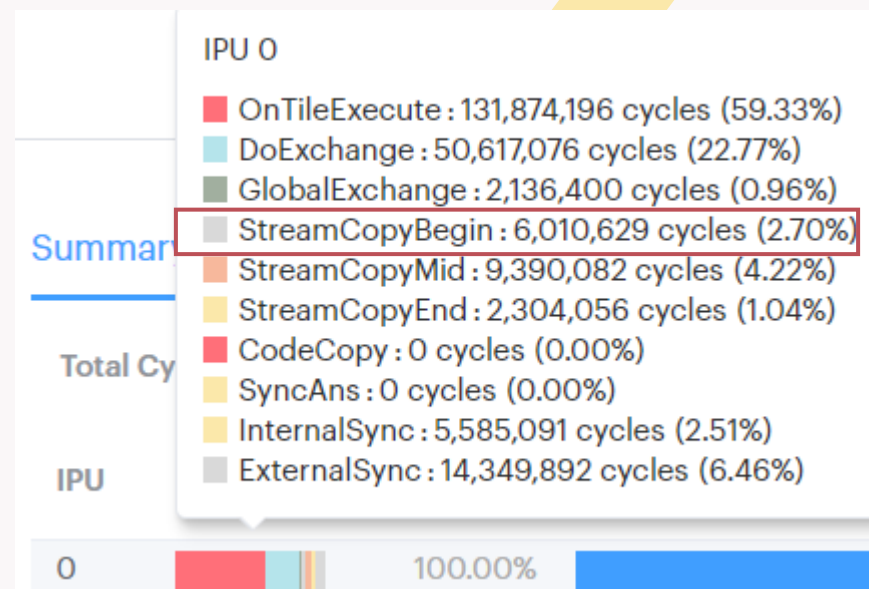
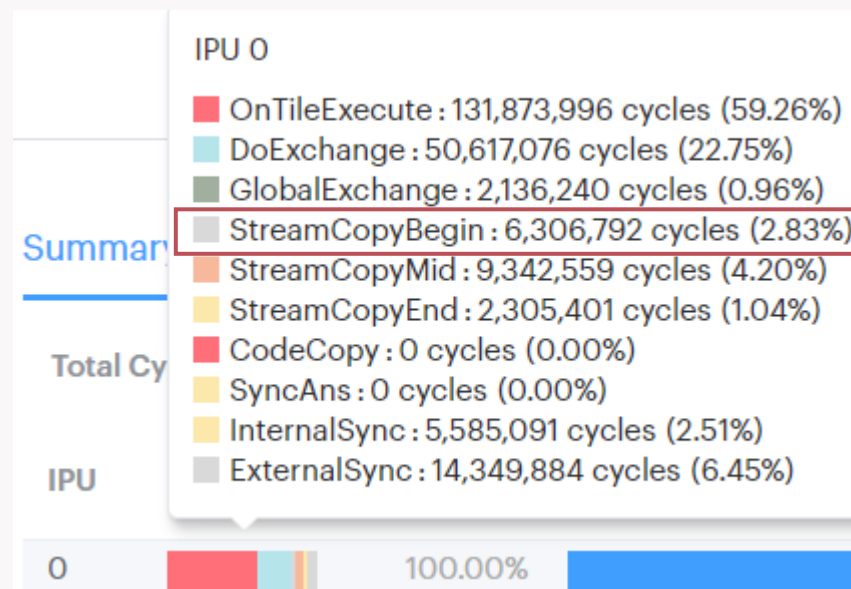
- Output

```
[Epoch 01]: 100%|██████████| 24/24 [08:32<00:00, 21.35s/it, loss=1.27, acc=50]
[Epoch 01] loss: 1.758, acc: 34.2
[Epoch 02]: 100%|██████████| 24/24 [00:08<00:00, 2.75it/s, loss=1.4, acc=50]
[Epoch 02] loss: 1.352, acc: 49.4
[Epoch 03]: 100%|██████████| 24/24 [00:08<00:00, 2.71it/s, loss=1.03, acc=65.6]
[Epoch 03] loss: 1.200, acc: 56.9
[Epoch 04]: 100%|██████████| 24/24 [00:08<00:00, 2.78it/s, loss=0.954, acc=65.6]
[Epoch 04] loss: 1.039, acc: 64.2
[Epoch 05]: 100%|██████████| 24/24 [00:08<00:00, 2.76it/s, loss=1.1, acc=62.5]
[Epoch 05] loss: 1.003, acc: 65.1
```

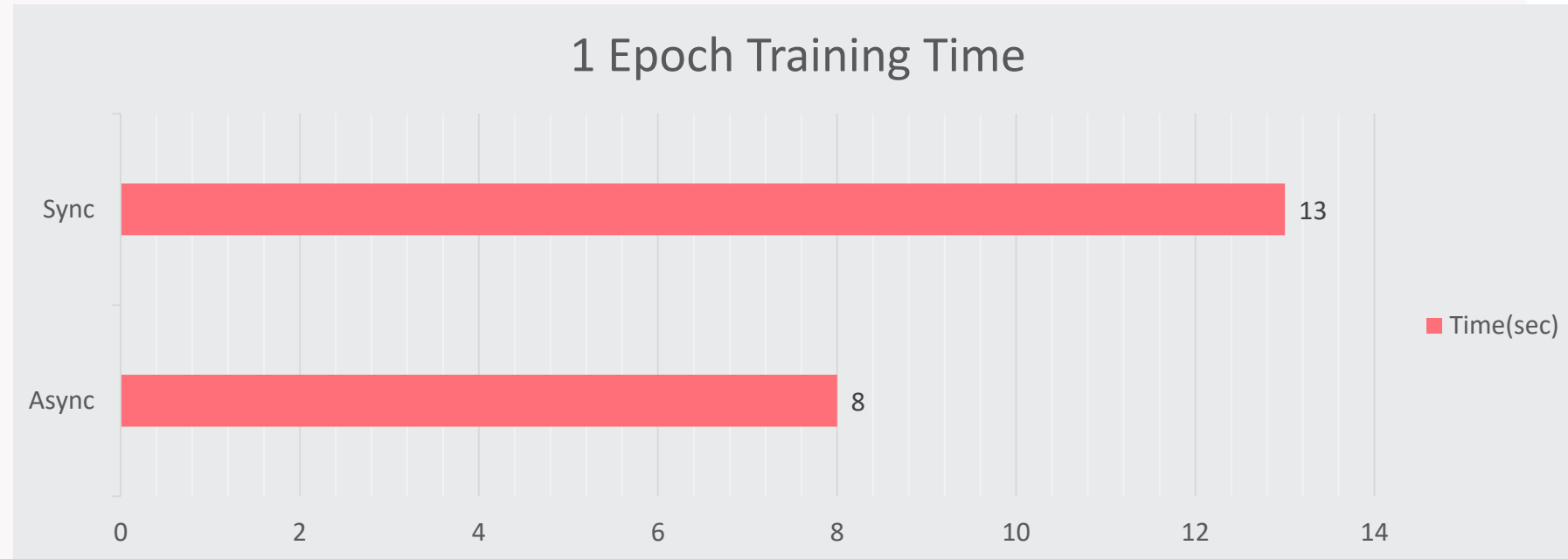
# ASYNCHRONOUS DATALOADER

Execution Trace

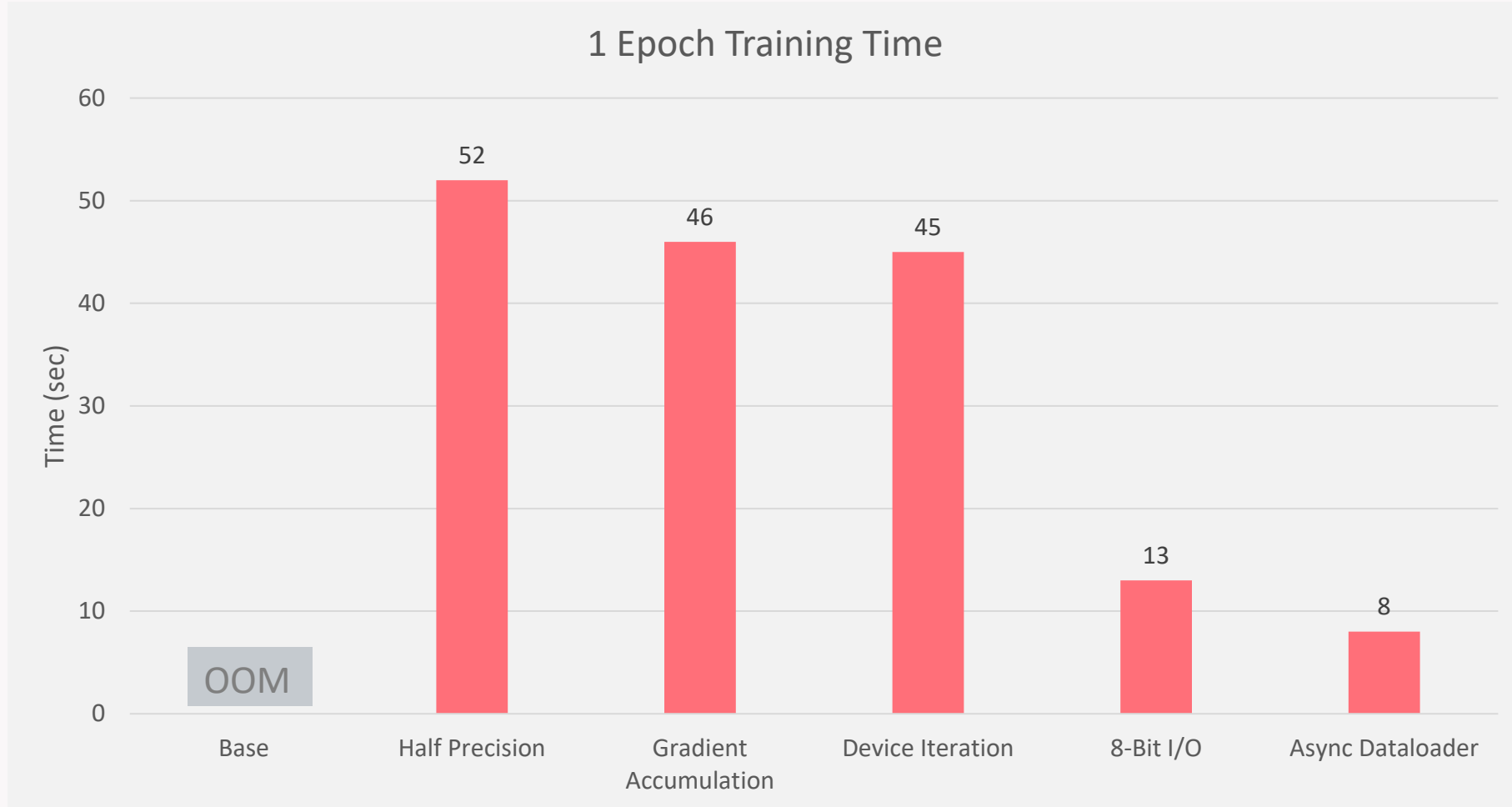
Time spent for waiting CPU to prepare the data is reduced



# ASYNCHRONOUS DATALOADER



# SUMMARY



Performance gain from each techniques might differ for different applications





# THANK YOU