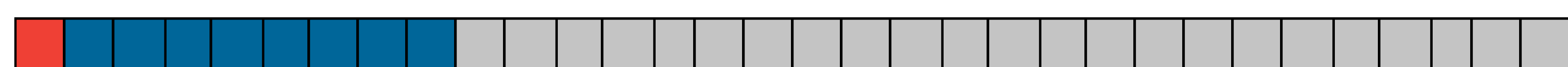


Our Framework

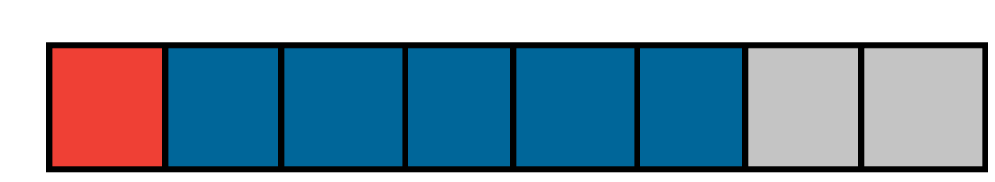
- Facilitates low-precision training research
- Supports various number formats and rounding options
- Implements state-of-the-art low-precision training techniques

Low-Precision Computation

IEEE FP32



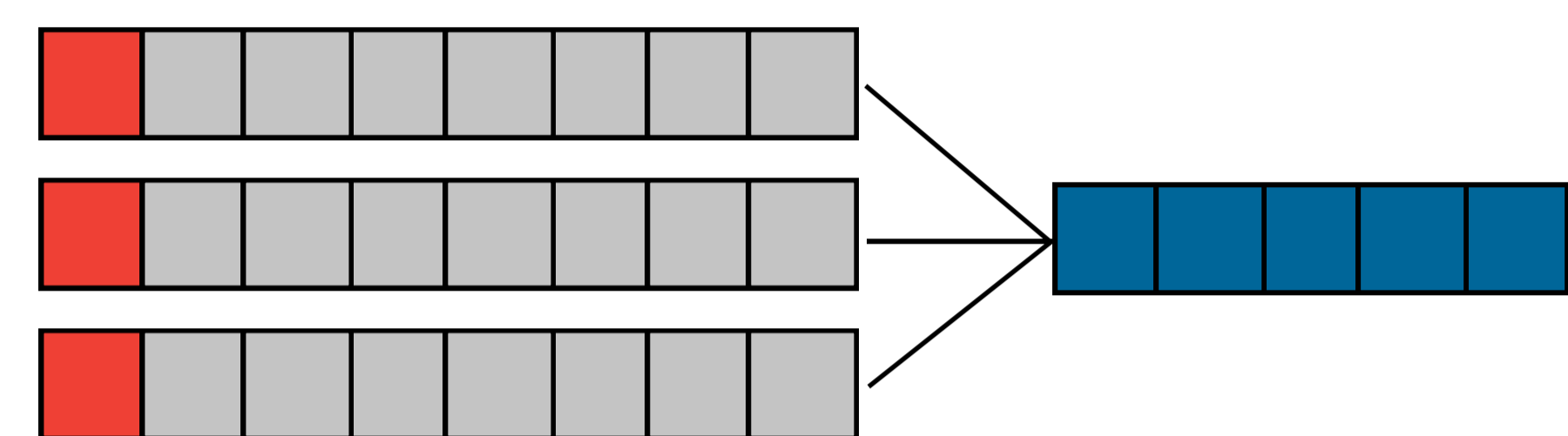
Custom FP8



Fixed 8



Block FP 8

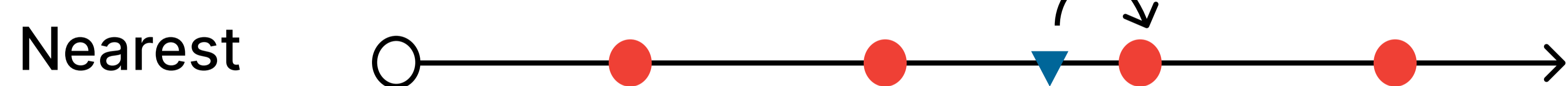


■ Sign ■ Exponent ■ Mantissa

QPyTorch can simulate various number formats:

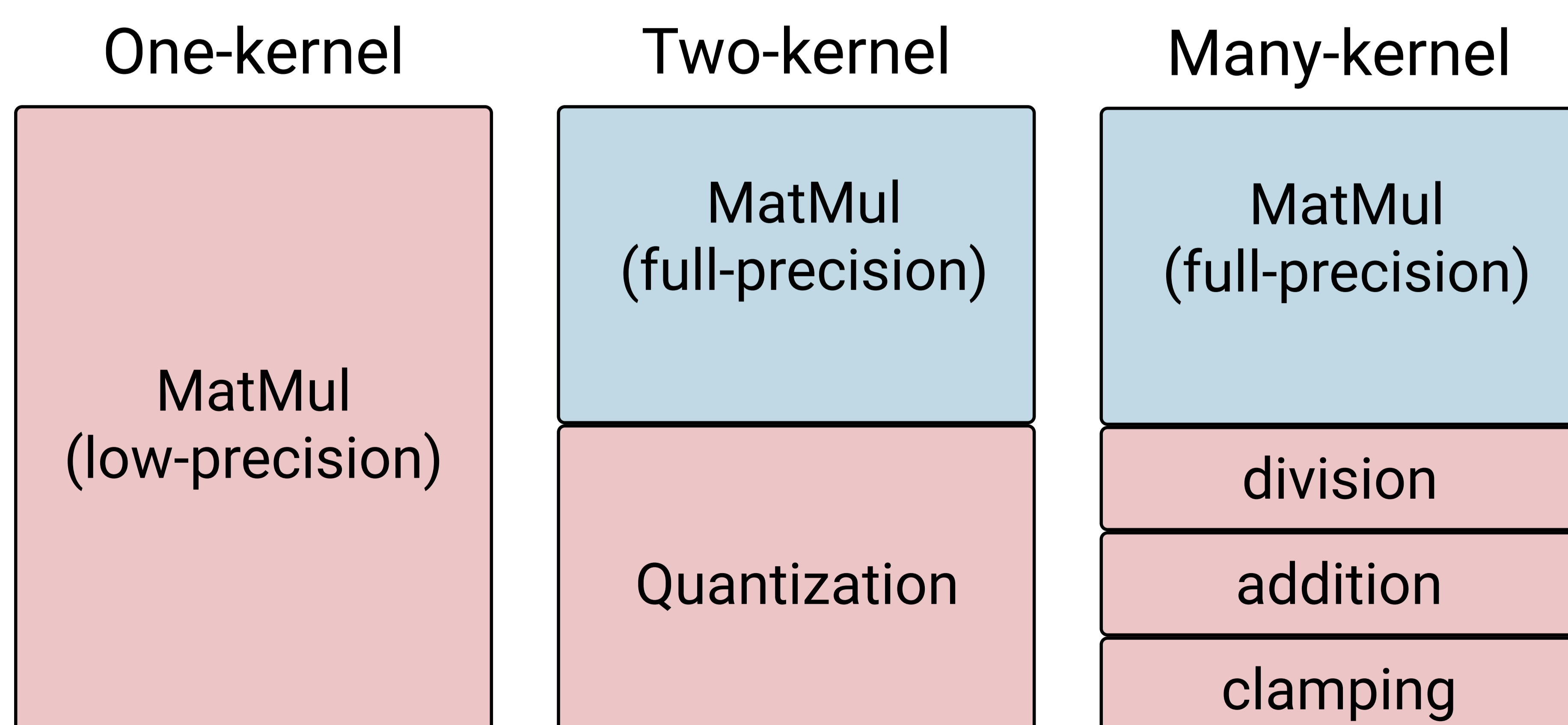
- floating point: lower precision than single precision
- fixed point: arbitrary precision
- block floating point: lower precision than single precision

Quantization and Rounding



QPyTorch can simulate various rounding strategies

Efficient Simulation

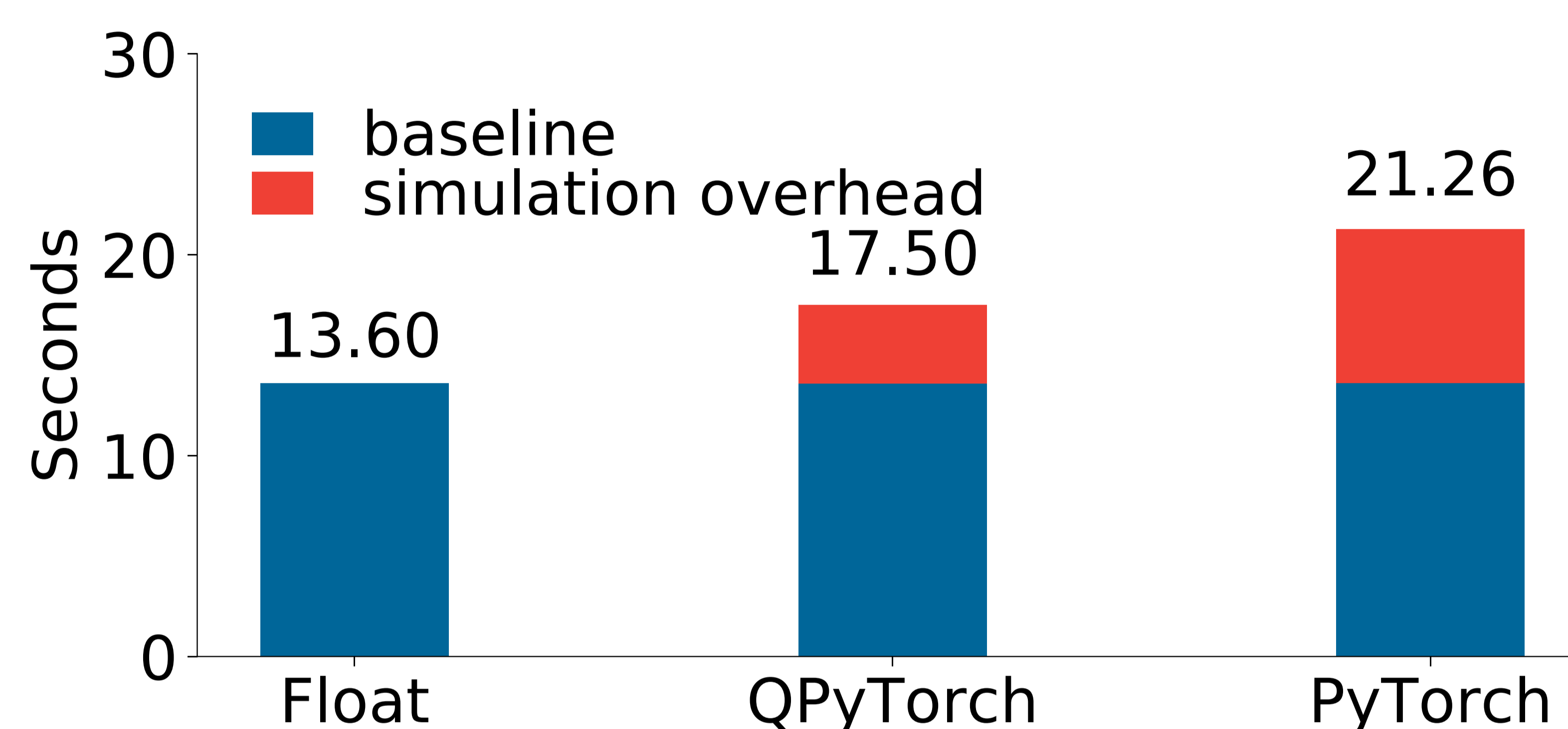
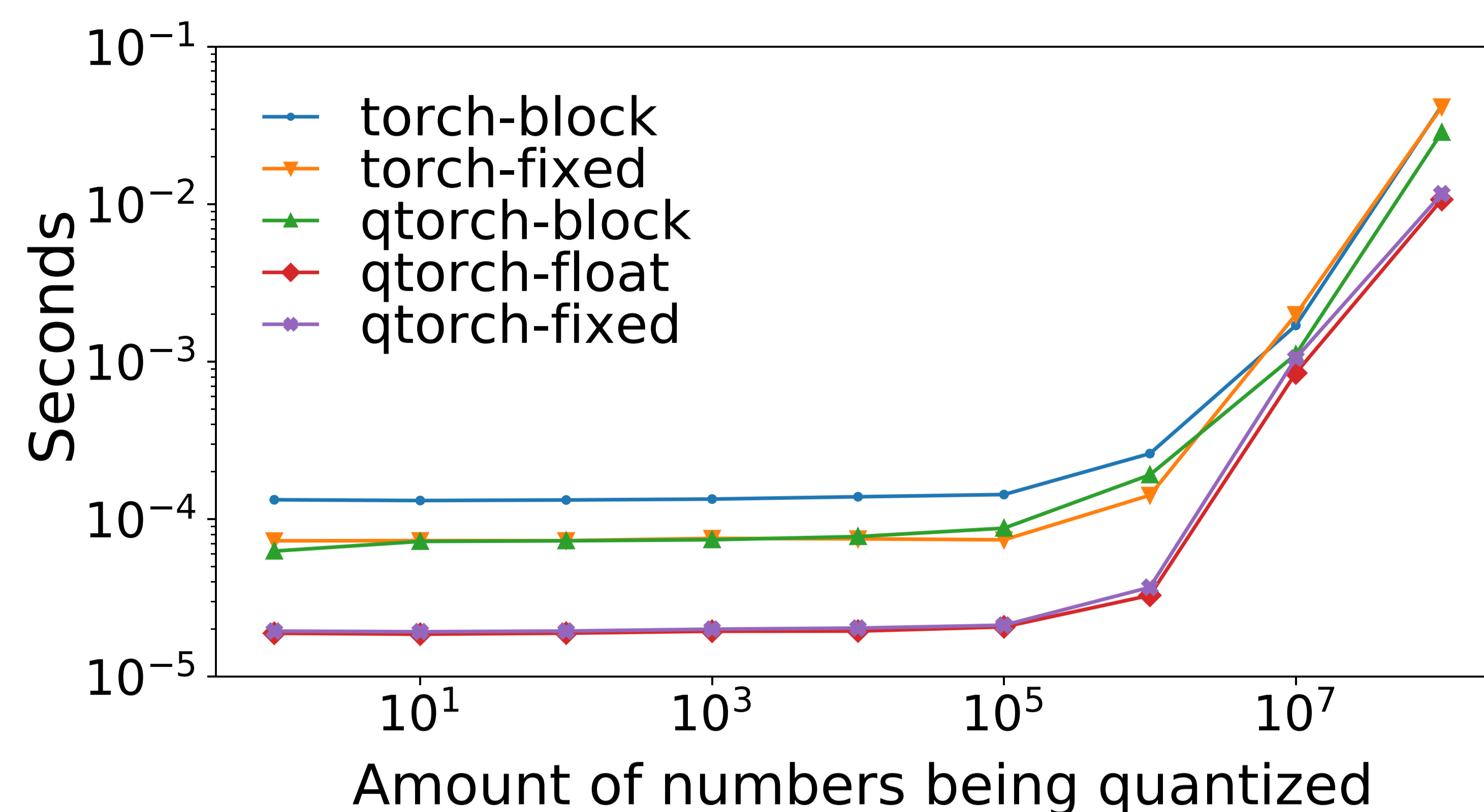


Speed 😞
Usage 😞

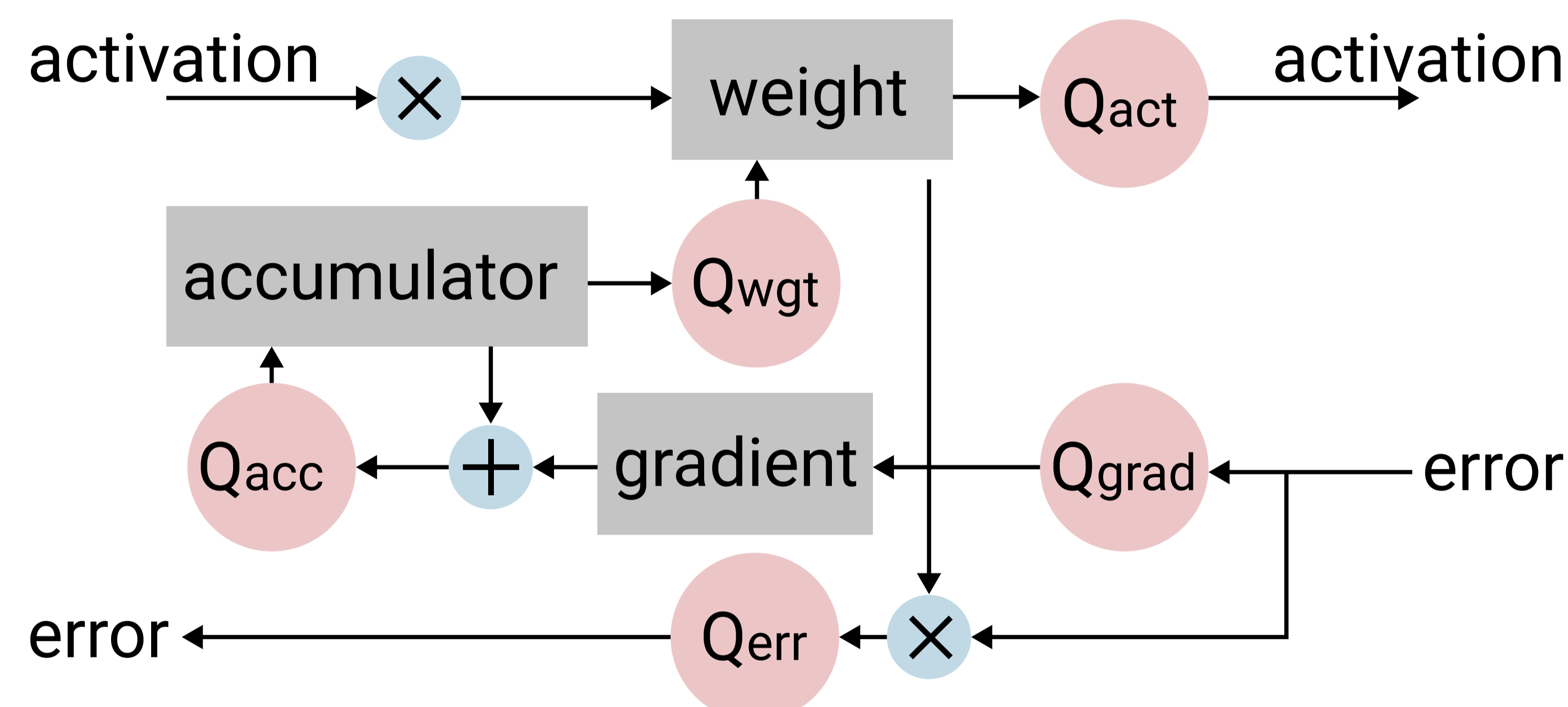
Speed 😊
Usage 😊

Speed 😞
Usage 😊

Experiment



Front-end Interface



Code Example

```
prec = FloatingPoint(bit=8)
acc_prec = FloatingPoint(bit=16)

model = net()
model = lower(model, layer_type=["conv"],
              forward_number=prec,
              backward_number=prec)

optim = SGD(model.parameters(), lr=0.01)
optim = OptimLP(optim, weight_quant=prec,
               grad_quant=prec,
               acc_quant=acc_prec)
```

Installation

`pip install qtorch`

