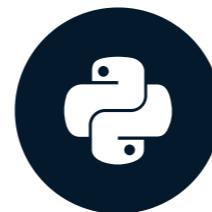


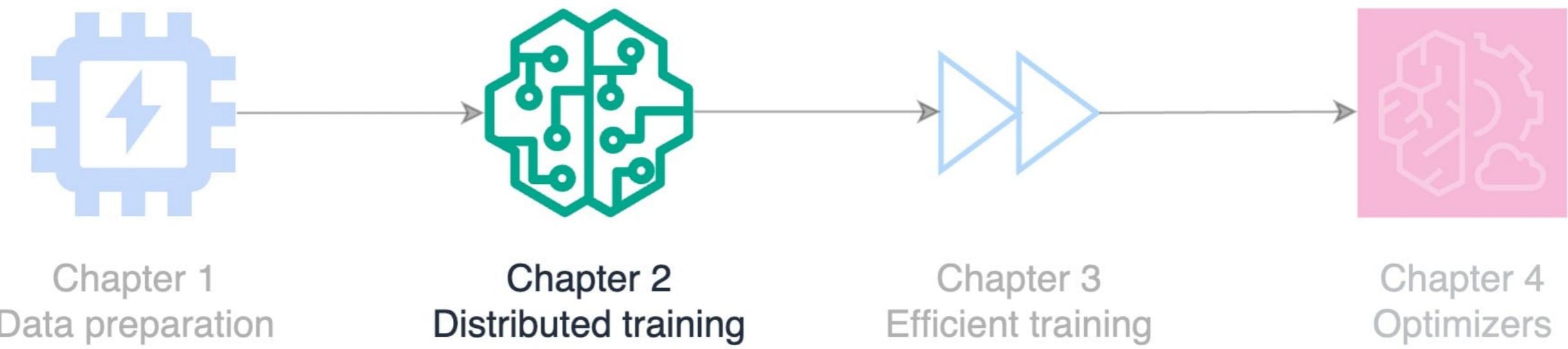
Gradient accumulation

EFFICIENT AI MODEL TRAINING WITH PYTORCH

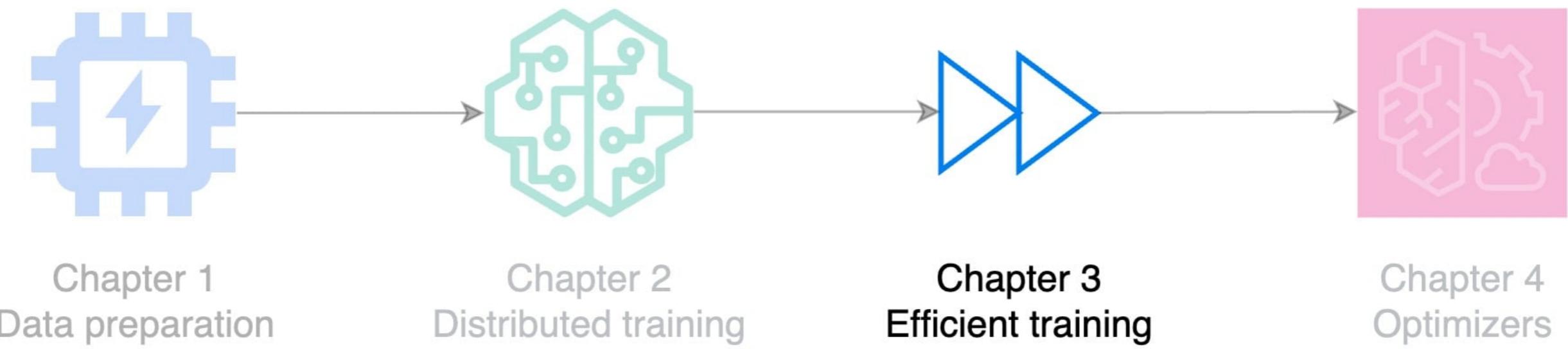


Dennis Lee
Data Engineer

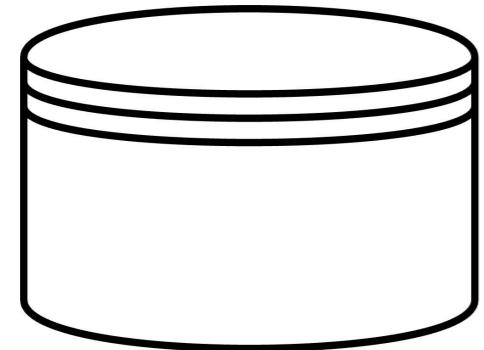
Distributed training



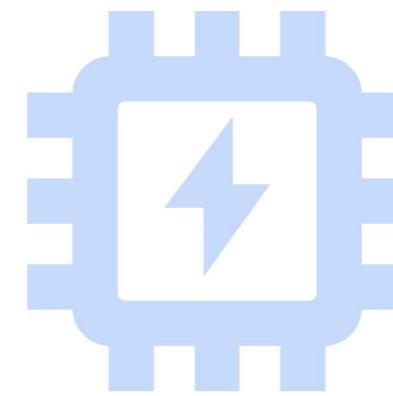
Efficient training



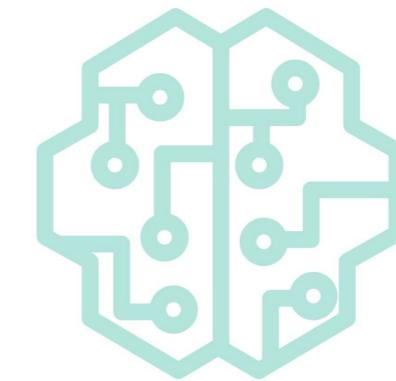
Improving training efficiency



**Memory
Efficiency**

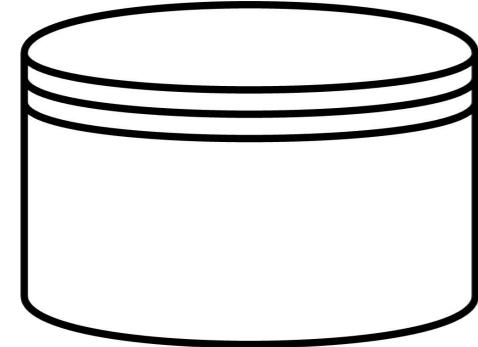


**Communication
Efficiency**

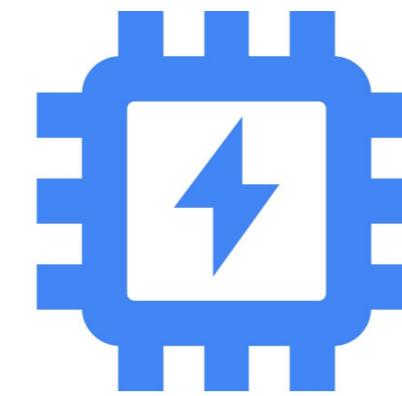


**Computational
Efficiency**

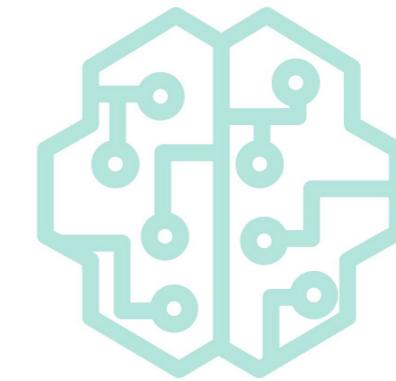
Improving training efficiency



Memory
Efficiency

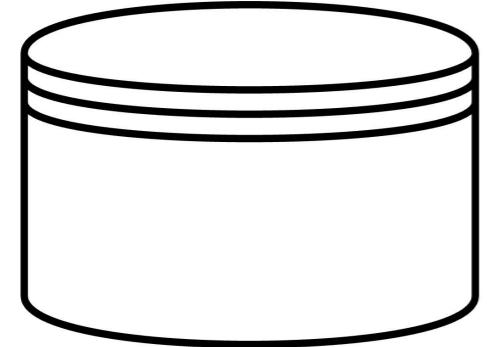


Communication
Efficiency

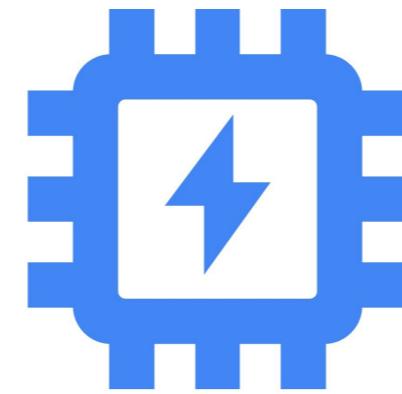


Computational
Efficiency

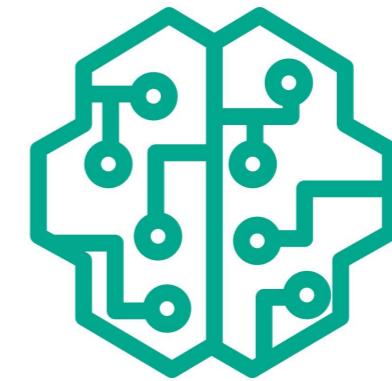
Improving training efficiency



Memory
Efficiency

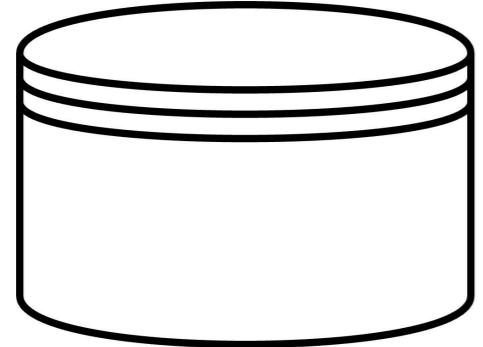


Communication
Efficiency

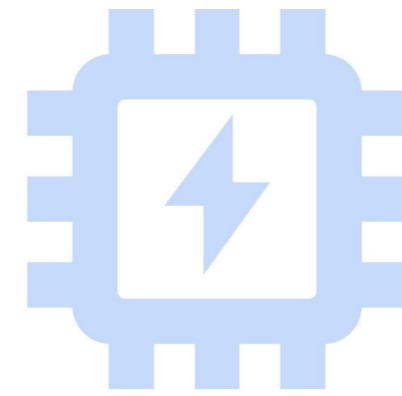


Computational
Efficiency

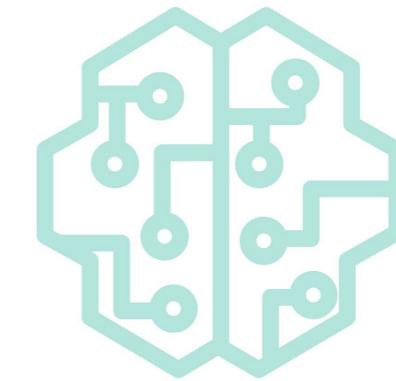
Gradient accumulation improves memory efficiency



Memory
Efficiency



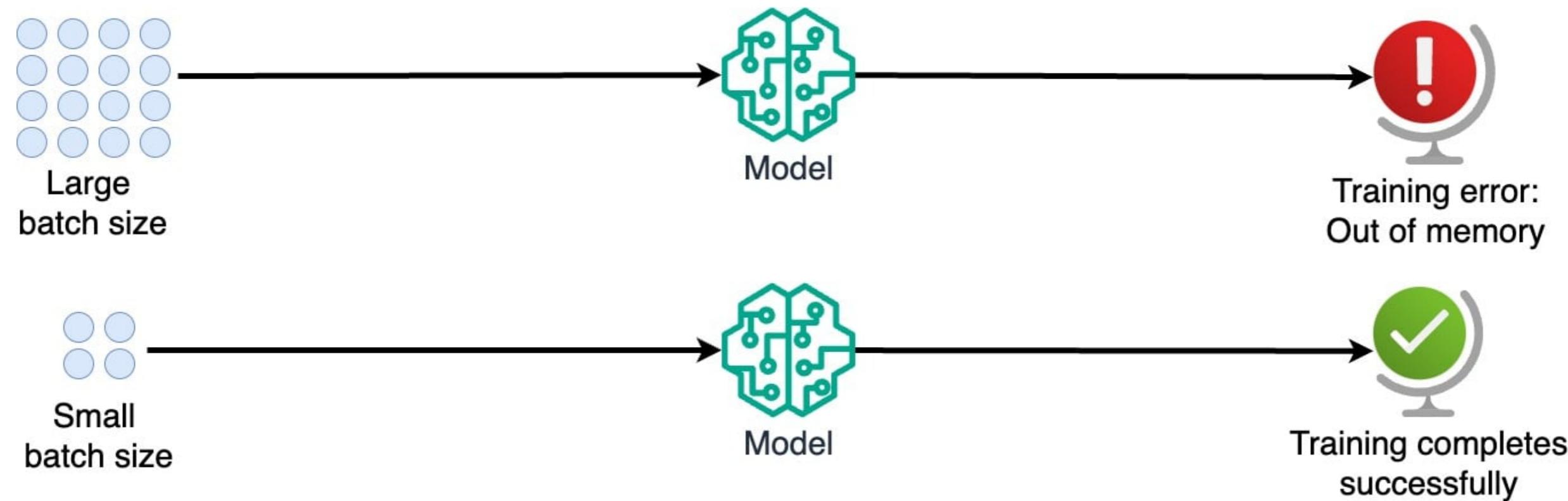
Communication
Efficiency



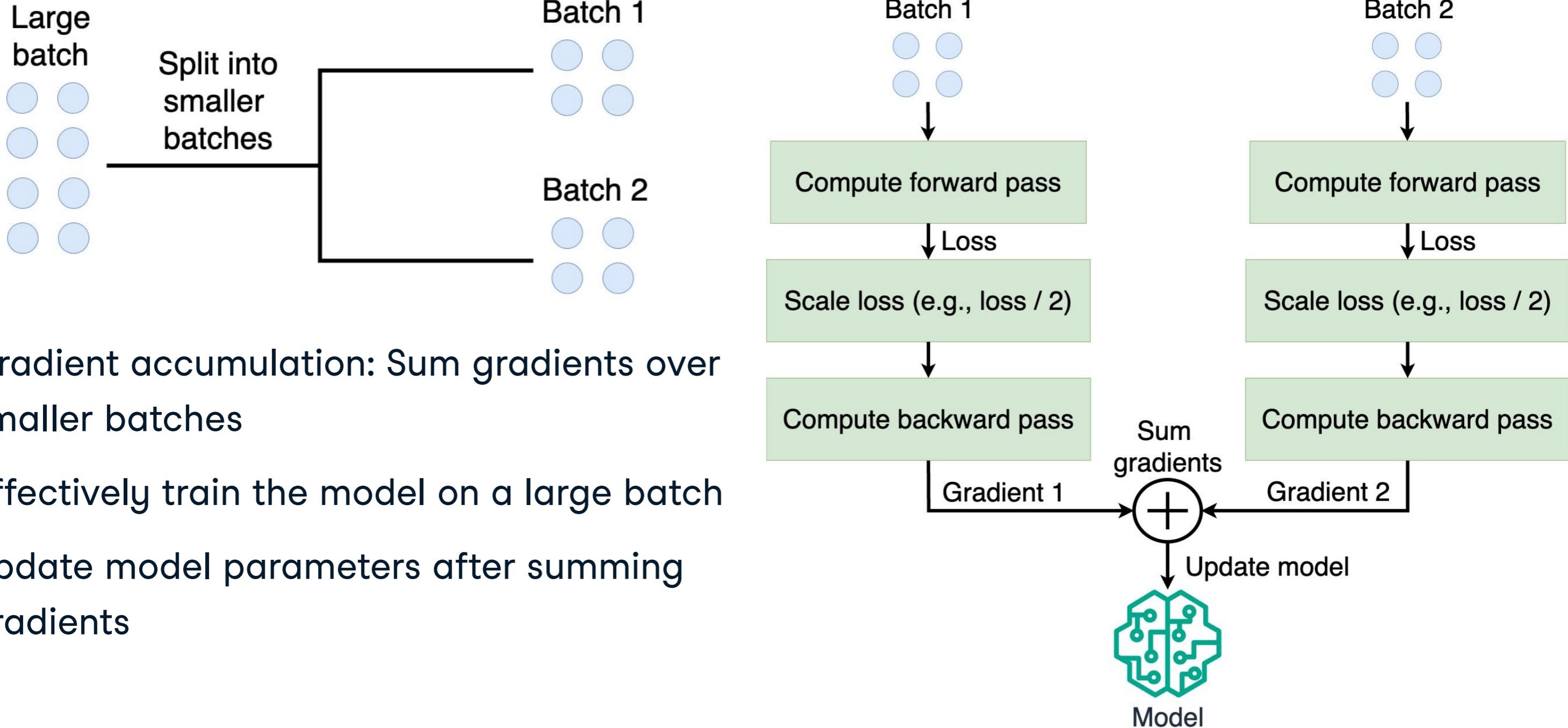
Computational
Efficiency

The problem with large batch sizes

- Large batch sizes: Robust gradient estimates for quicker learning
- GPU memory constrains batch sizes

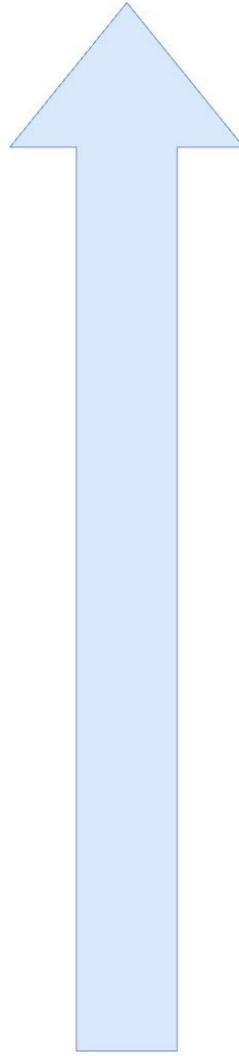


How does gradient accumulation work?



PyTorch, Accelerator, and Trainer

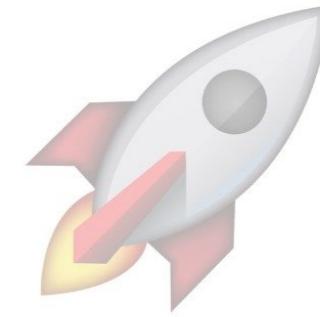
Ability to Customize



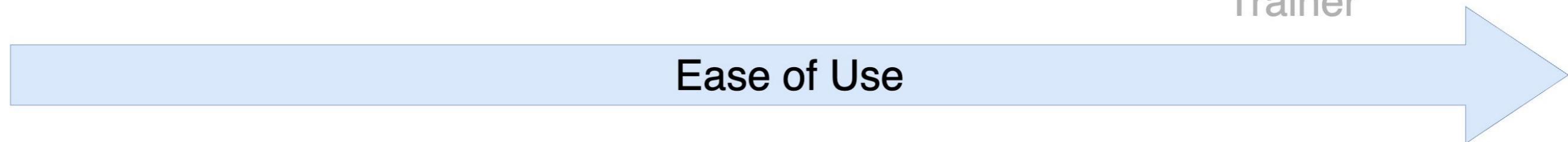
PyTorch



Accelerator



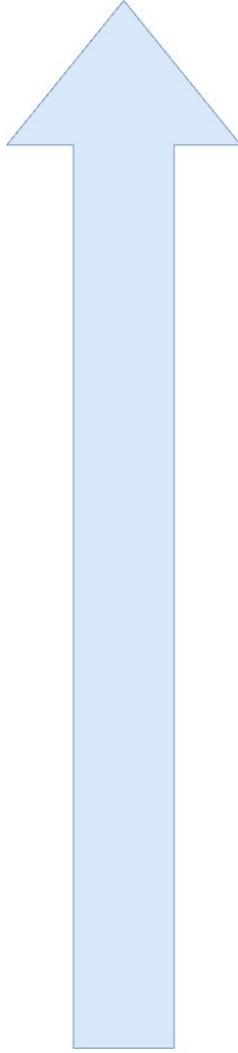
Trainer



Ease of Use

PyTorch, Accelerator, and Trainer

Ability to Customize



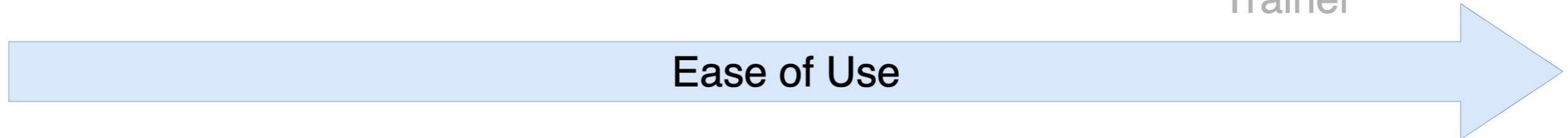
PyTorch



Accelerator



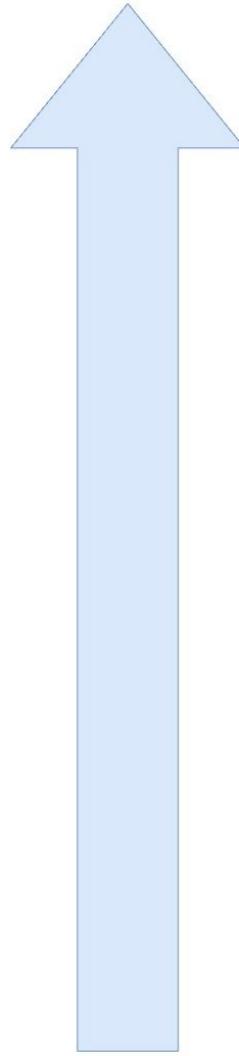
Trainer



Ease of Use

PyTorch, Accelerator, and Trainer

Ability to Customize



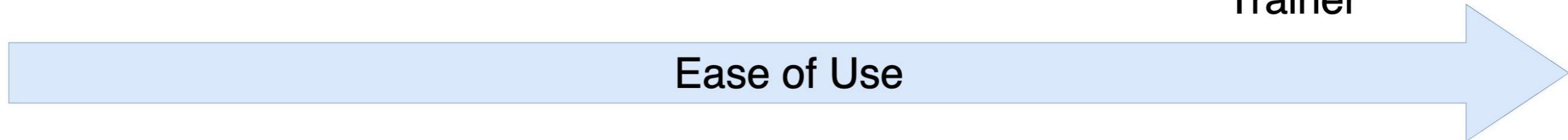
PyTorch



Accelerator



Trainer

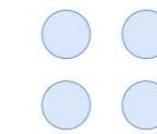


Ease of Use

Gradient accumulation with PyTorch

```
for index, batch in enumerate(dataloader):
    inputs, targets = (batch["input_ids"],
                        batch["labels"])
    inputs, targets = (inputs.to(device),
                       targets.to(device))
```

Batch 1



Compute forward pass

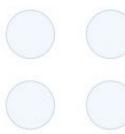
Loss

Scale loss (e.g., loss / 2)

Compute backward pass

Gradient 1

Batch 2



Compute forward pass

Loss

Scale loss (e.g., loss / 2)

Compute backward pass

Gradient 2



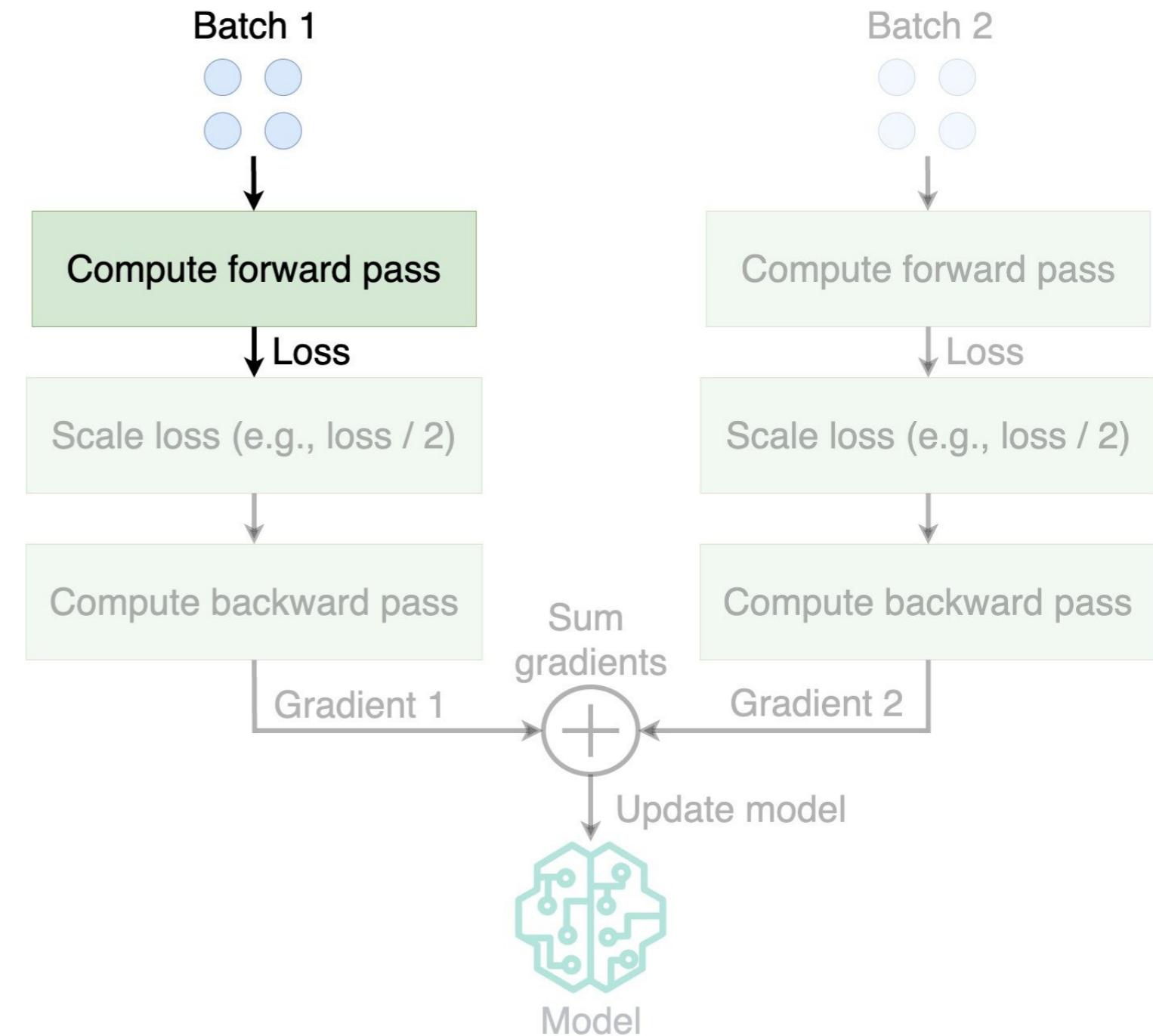
Model



Model

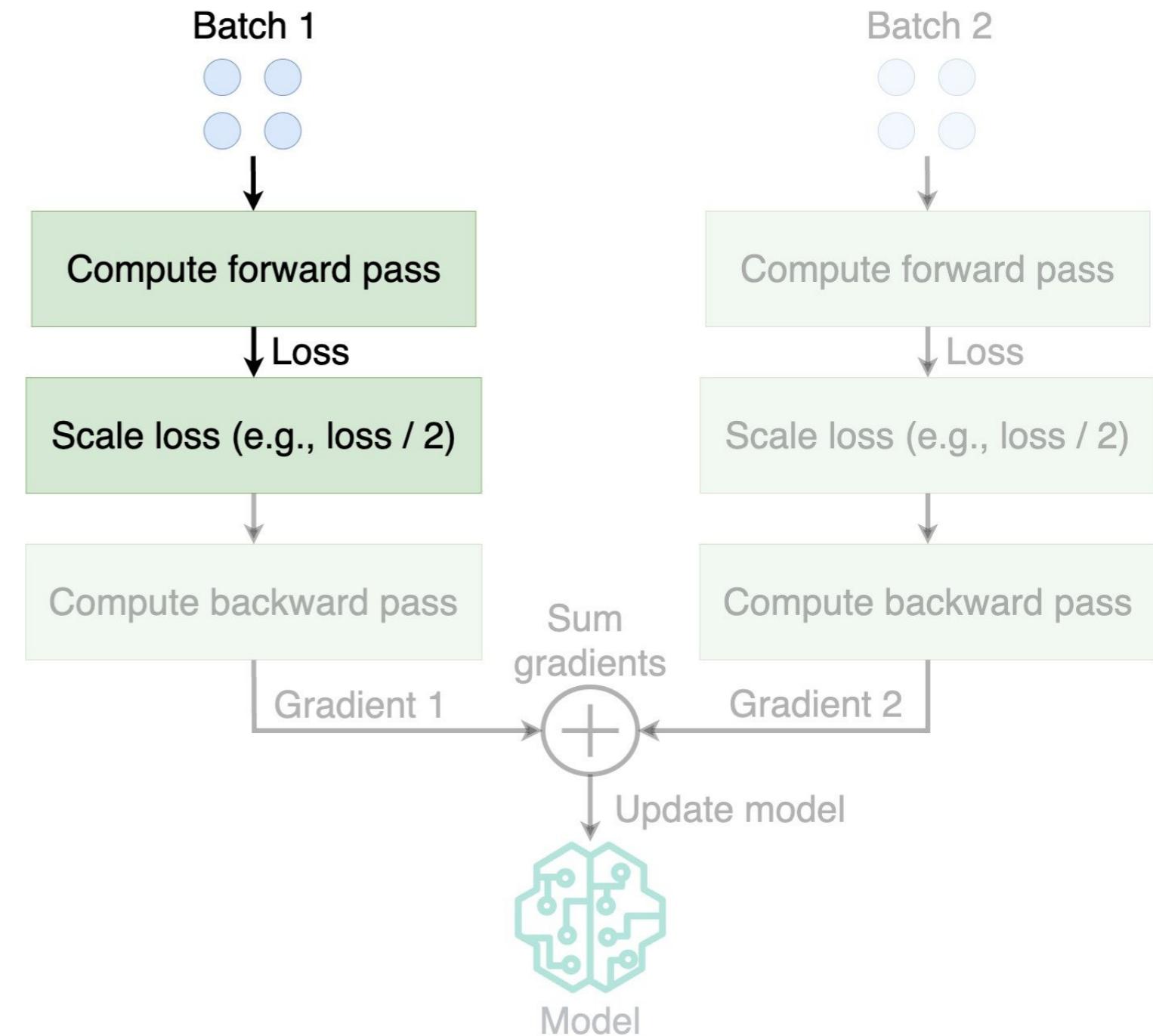
Gradient accumulation with PyTorch

```
for index, batch in enumerate(dataloader):
    inputs, targets = (batch["input_ids"],
                        batch["labels"])
    inputs, targets = (inputs.to(device),
                       targets.to(device))
    outputs = model(inputs, labels=targets)
    loss = outputs.loss
```



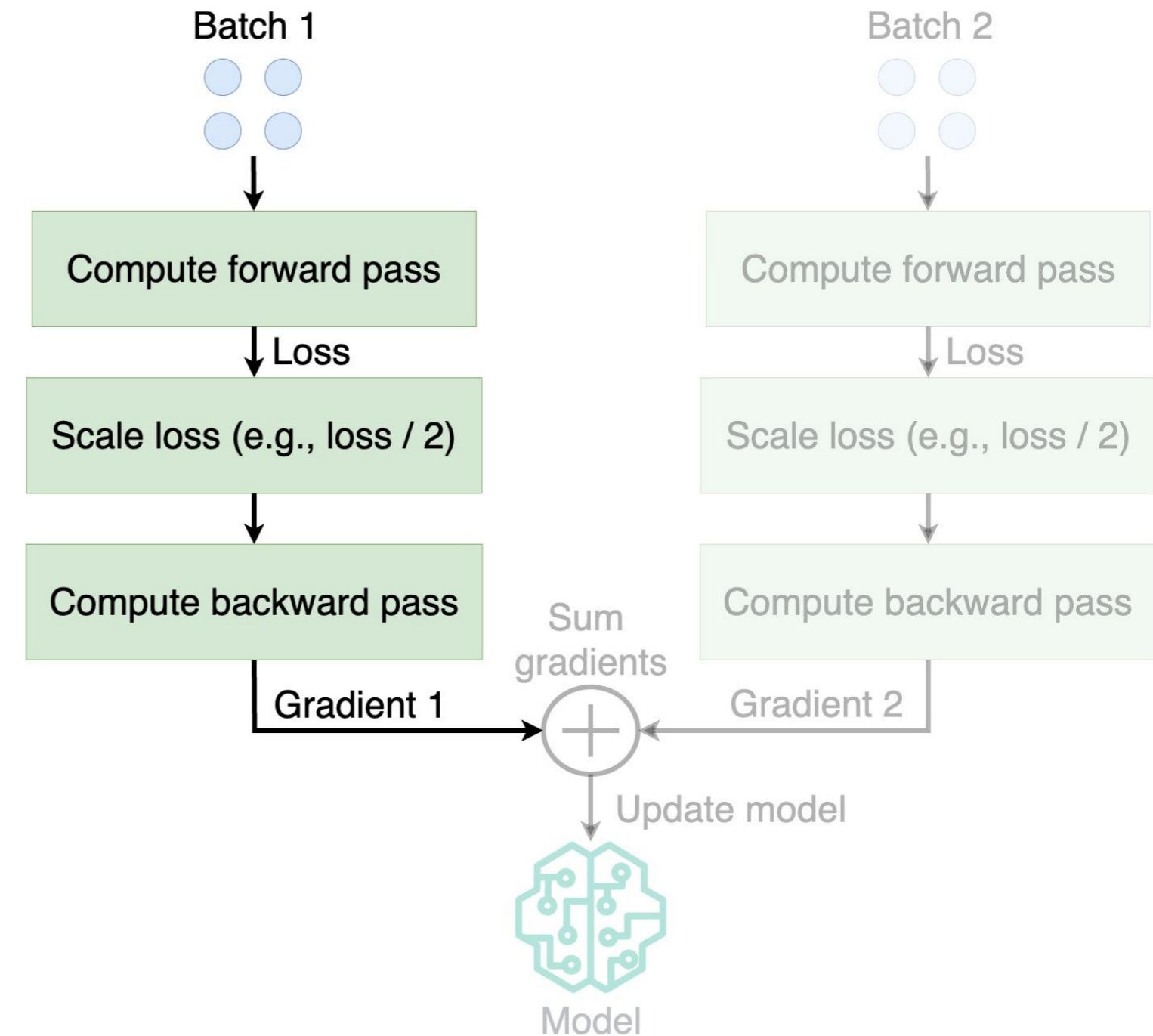
Gradient accumulation with PyTorch

```
for index, batch in enumerate(dataloader):
    inputs, targets = (batch["input_ids"],
                        batch["labels"])
    inputs, targets = (inputs.to(device),
                       targets.to(device))
    outputs = model(inputs, labels=targets)
    loss = outputs.loss
    loss = loss / gradient_accumulation_steps
```



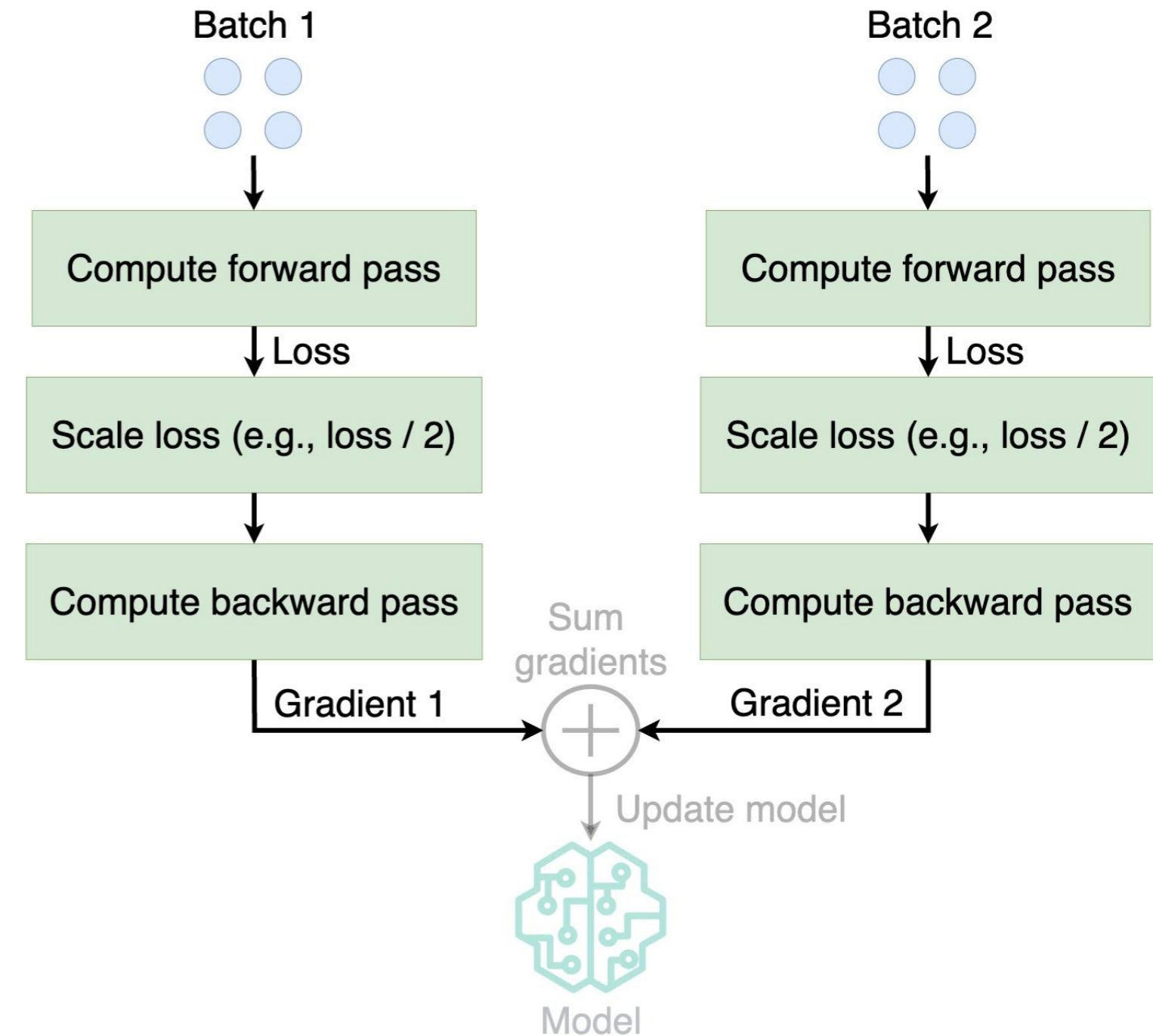
Gradient accumulation with PyTorch

```
for index, batch in enumerate(dataloader):
    inputs, targets = (batch["input_ids"],
                        batch["labels"])
    inputs, targets = (inputs.to(device),
                       targets.to(device))
    outputs = model(inputs, labels=targets)
    loss = outputs.loss
    loss = loss / gradient_accumulation_steps
    loss.backward()
```



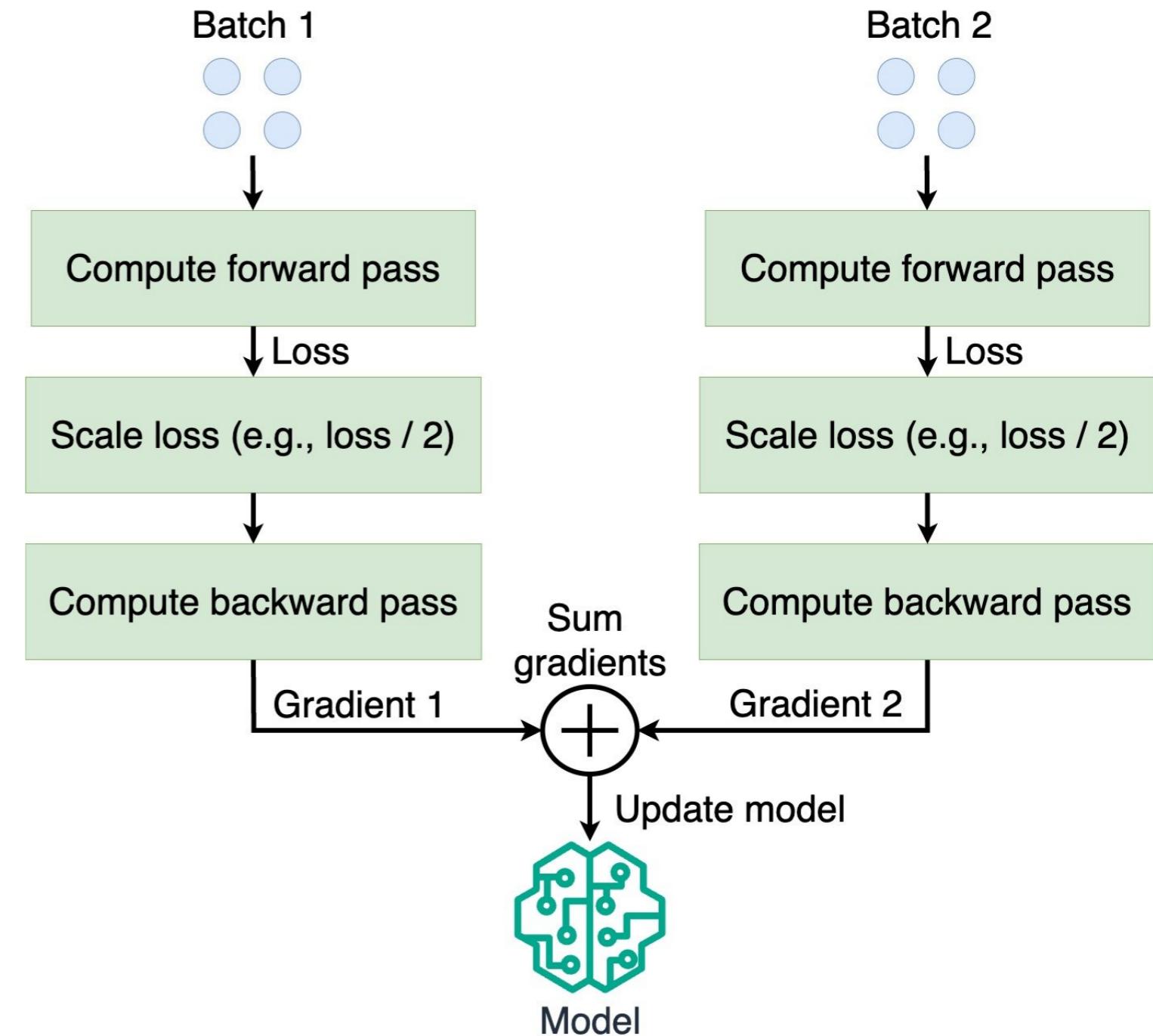
Gradient accumulation with PyTorch

```
for index, batch in enumerate(dataloader):
    inputs, targets = (batch["input_ids"],
                        batch["labels"])
    inputs, targets = (inputs.to(device),
                       targets.to(device))
    outputs = model(inputs, labels=targets)
    loss = outputs.loss
    loss = loss / gradient_accumulation_steps
    loss.backward()
    if ((index + 1) % gradient_accumulation_steps == 0):
```



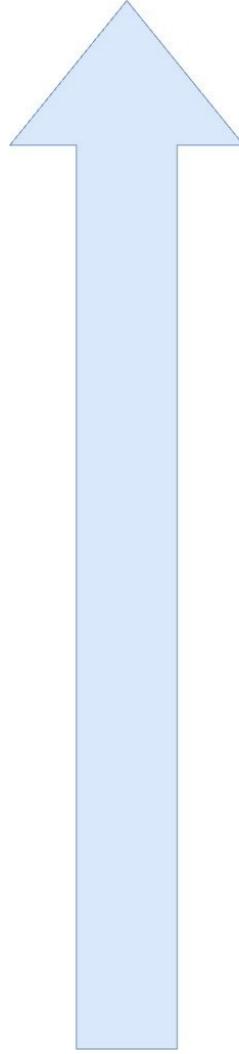
Gradient accumulation with PyTorch

```
for index, batch in enumerate(dataloader):
    inputs, targets = (batch["input_ids"],
                        batch["labels"])
    inputs, targets = (inputs.to(device),
                       targets.to(device))
    outputs = model(inputs, labels=targets)
    loss = outputs.loss
    loss = loss / gradient_accumulation_steps
    loss.backward()
    if ((index + 1) % gradient_accumulation_steps == 0):
        optimizer.step()
        lr_scheduler.step()
        optimizer.zero_grad()
```



From PyTorch to Accelerator

Ability to Customize



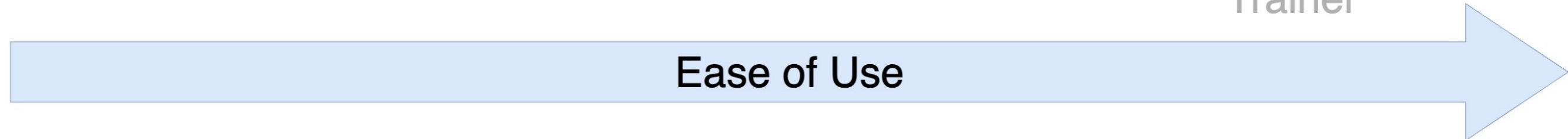
PyTorch



Accelerator



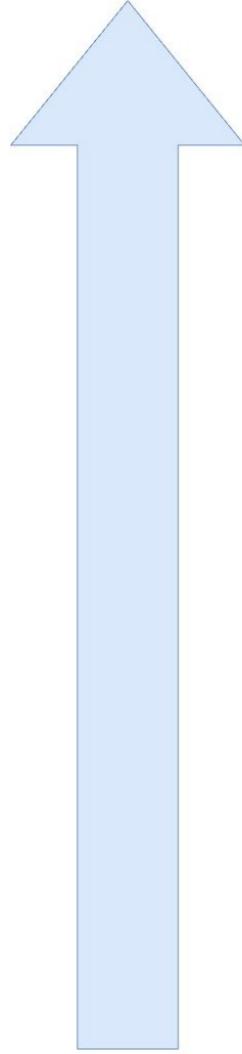
Trainer



Ease of Use

From PyTorch to Accelerator

Ability to Customize



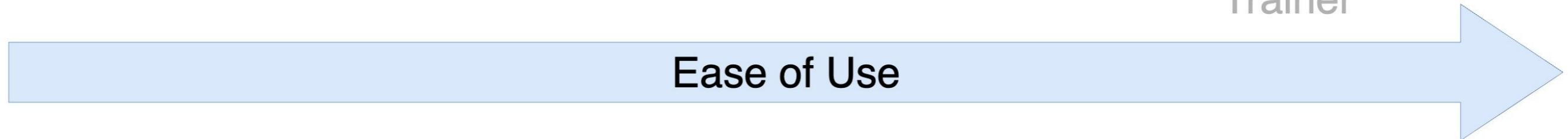
PyTorch



Accelerator



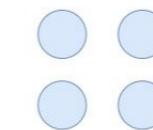
Trainer



Gradient accumulation with Accelerator

```
accelerator = \
    Accelerator(gradient_accumulation_steps=2)
for index, batch in enumerate(dataloader):
    inputs, targets = (batch["input_ids"],
                        batch["labels"])
```

Batch 1



Compute forward pass

Loss

Scale loss (e.g., loss / 2)

Compute backward pass

Gradient 1

Sum
gradients

Update model



Model

Batch 2



Compute forward pass

Loss

Scale loss (e.g., loss / 2)

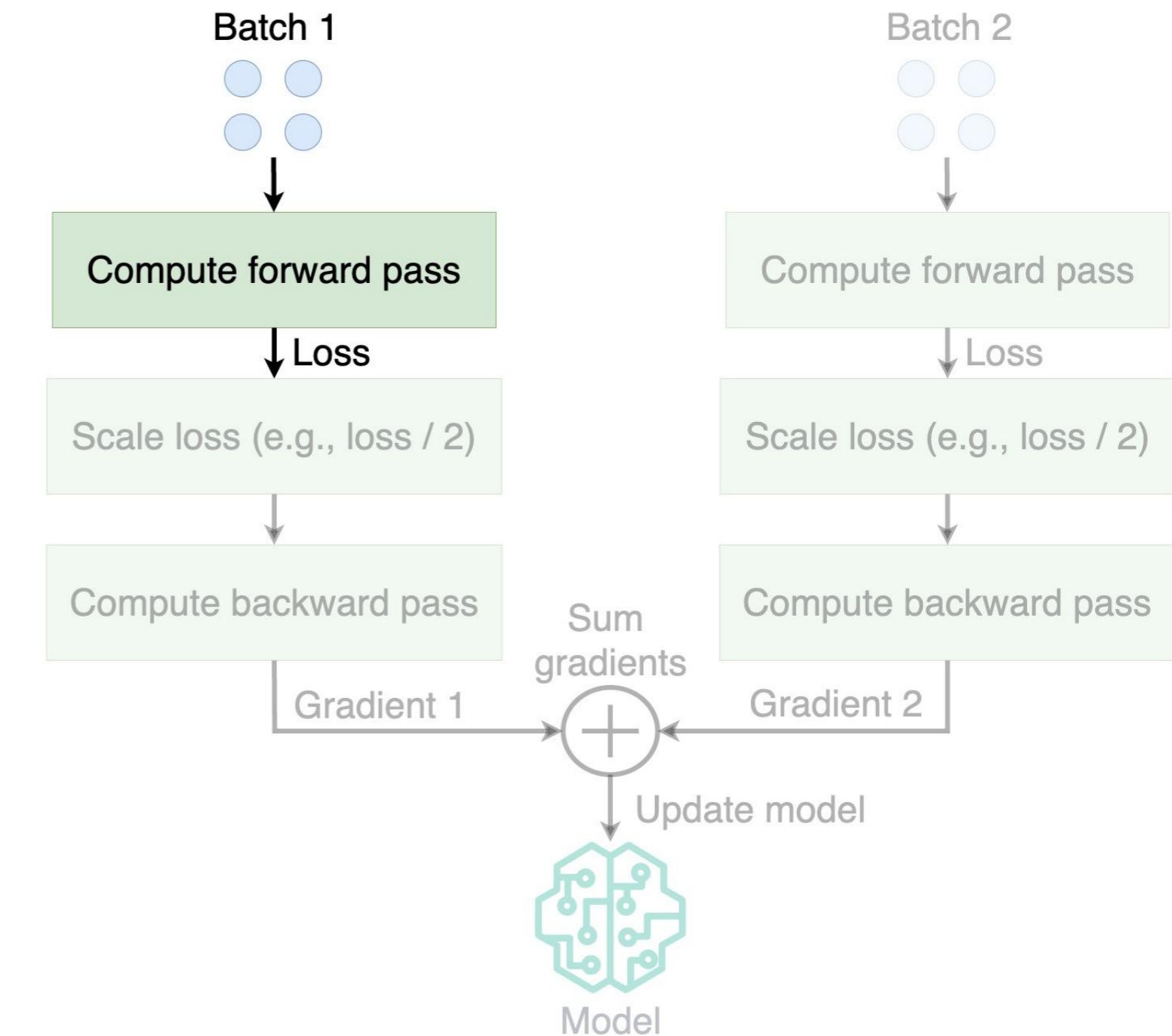
Compute backward pass

Gradient 2

Gradient accumulation with Accelerator

```
accelerator = \
    Accelerator(gradient_accumulation_steps=2)

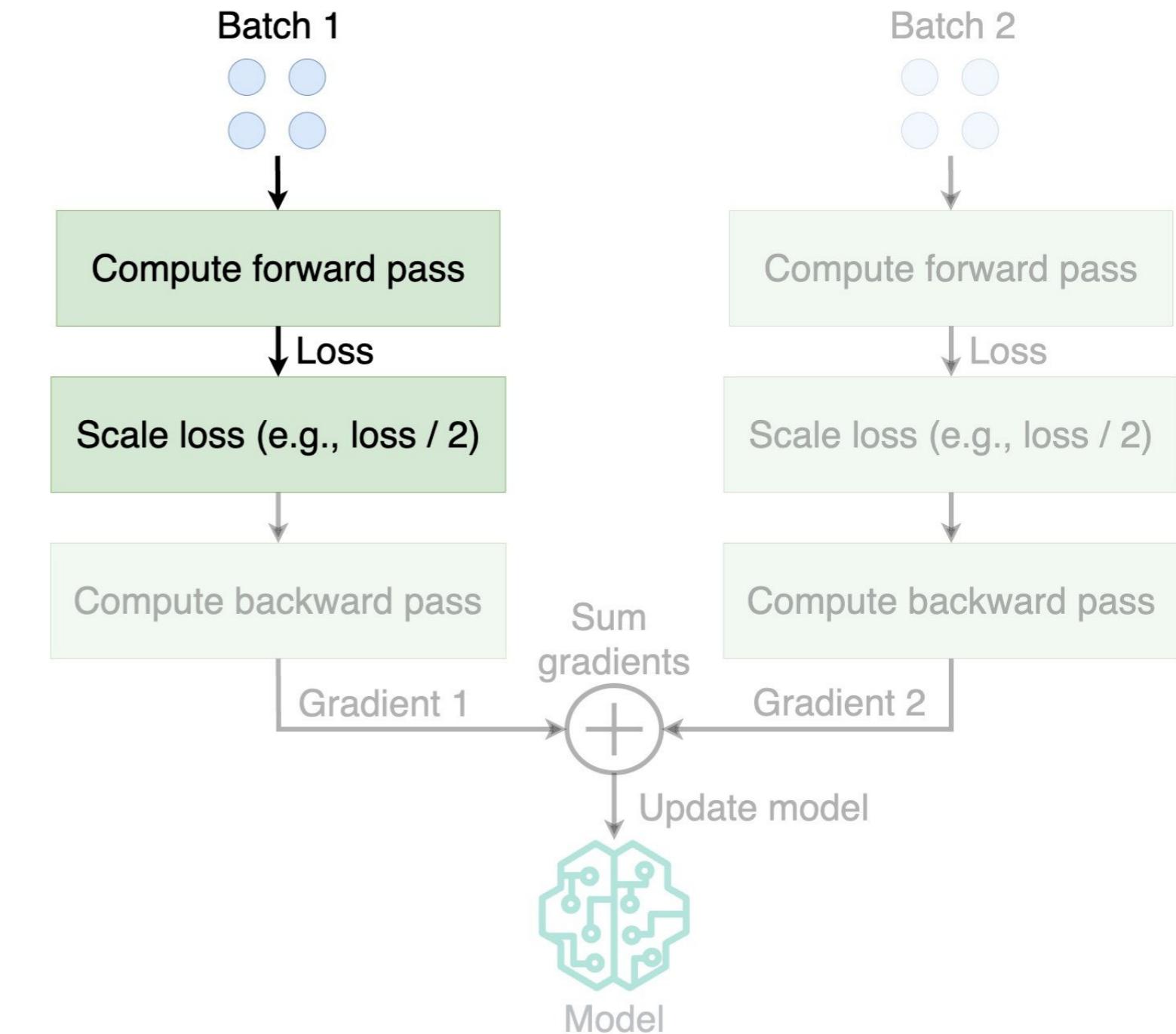
for index, batch in enumerate(dataloader):
    inputs, targets = (batch["input_ids"],
                        batch["labels"])
    outputs = model(inputs,
                     labels=targets)
    loss = outputs.loss
```



Gradient accumulation with Accelerator

```
accelerator = \
    Accelerator(gradient_accumulation_steps=2)

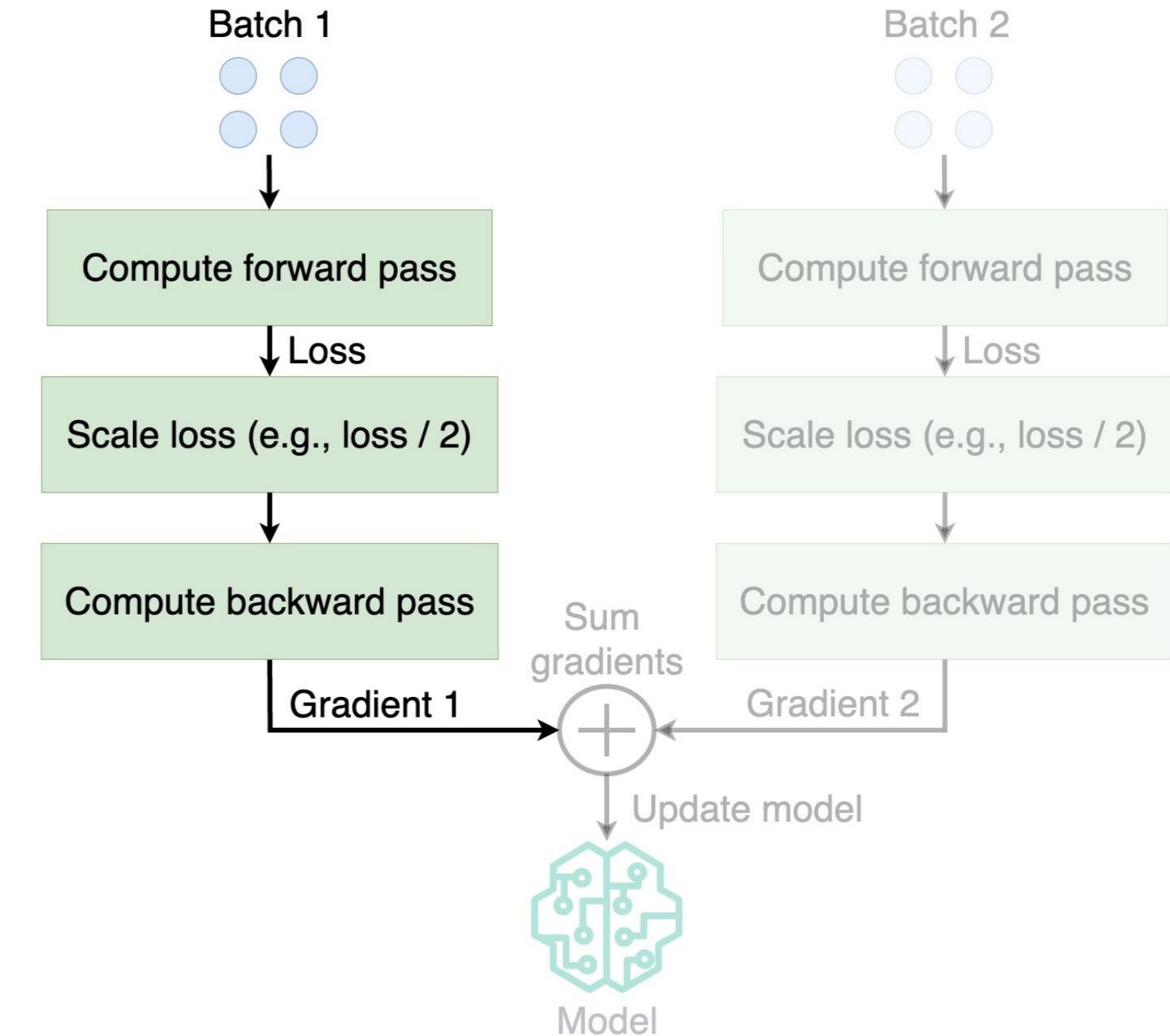
for index, batch in enumerate(dataloader):
    with accelerator.accumulate(model):
        inputs, targets = (batch["input_ids"],
                            batch["labels"])
        outputs = model(inputs,
                         labels=targets)
        loss = outputs.loss
```



Gradient accumulation with Accelerator

```
accelerator = \
    Accelerator(gradient_accumulation_steps=2)

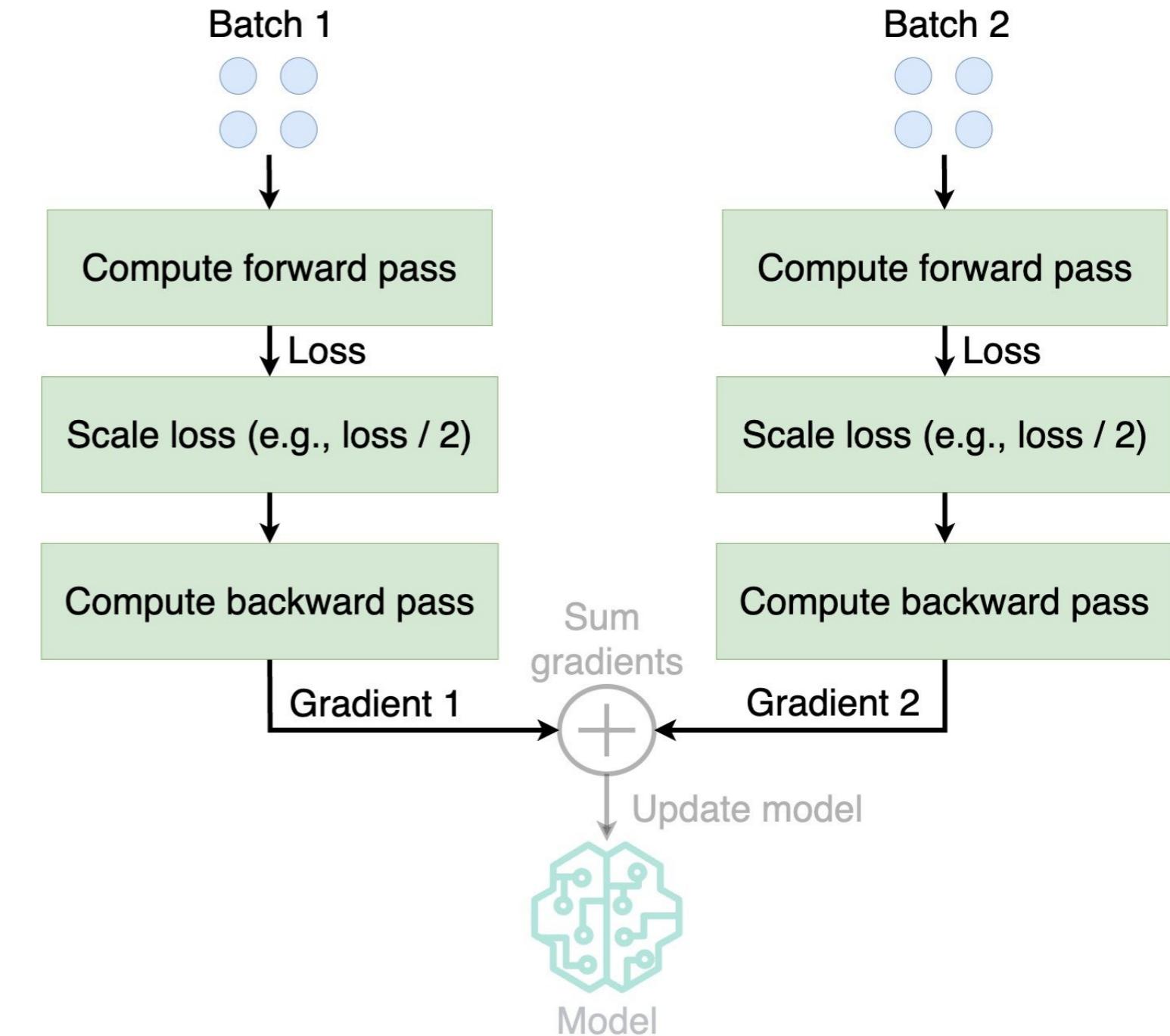
for index, batch in enumerate(dataloader):
    with accelerator.accumulate(model):
        inputs, targets = (batch["input_ids"],
                            batch["labels"])
        outputs = model(inputs,
                         labels=targets)
        loss = outputs.loss
        accelerator.backward(loss)
```



Gradient accumulation with Accelerator

```
accelerator = \
    Accelerator(gradient_accumulation_steps=2)

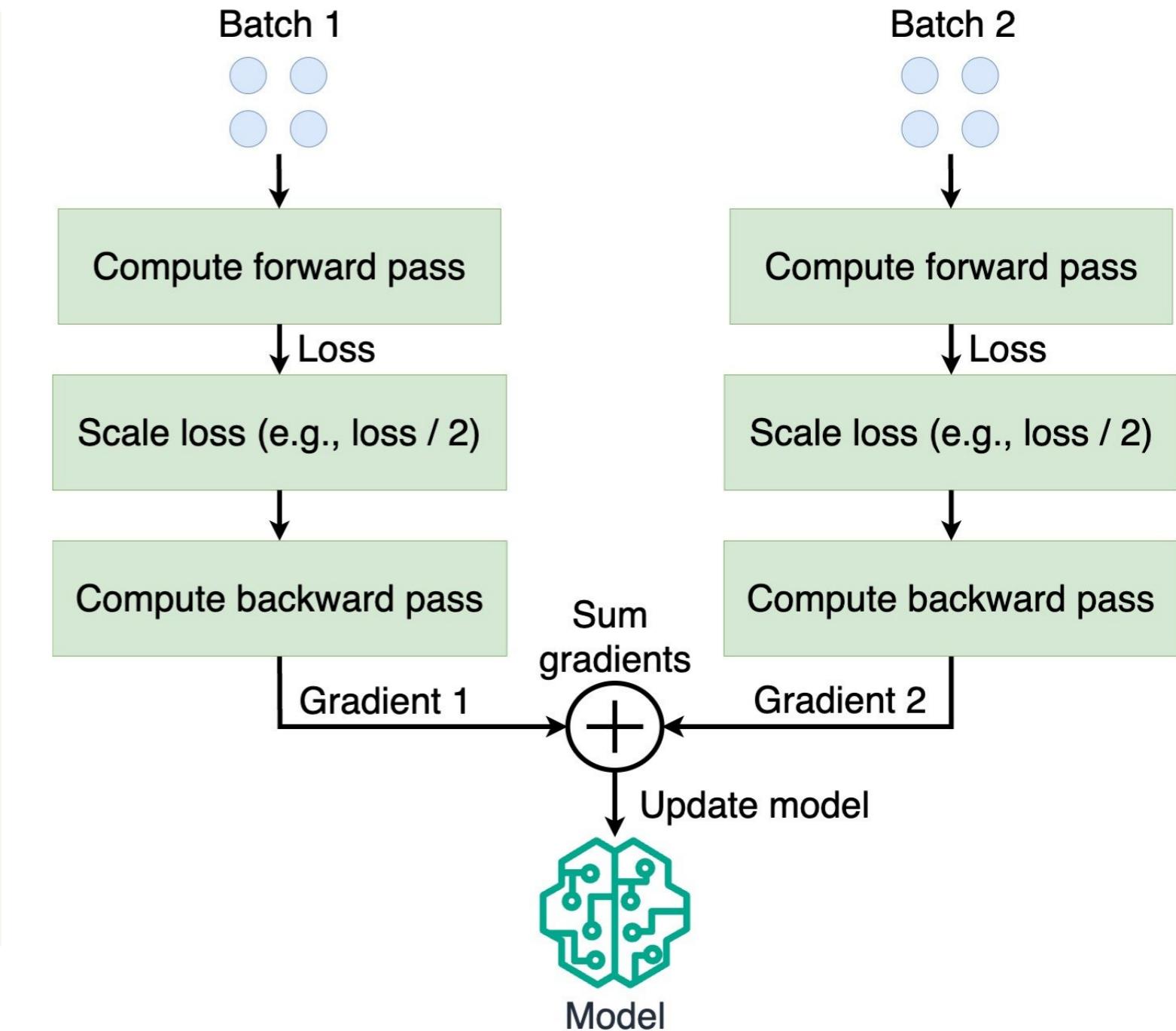
for index, batch in enumerate(dataloader):
    with accelerator.accumulate(model):
        inputs, targets = (batch["input_ids"],
                            batch["labels"])
        outputs = model(inputs,
                         labels=targets)
        loss = outputs.loss
        accelerator.backward(loss)
```



Gradient accumulation with Accelerator

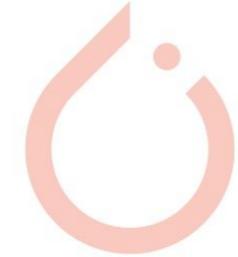
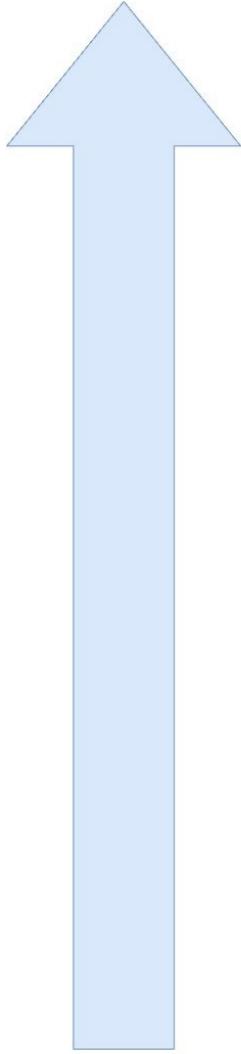
```
accelerator = \
    Accelerator(gradient_accumulation_steps=2)

for index, batch in enumerate(dataloader):
    with accelerator.accumulate(model):
        inputs, targets = (batch["input_ids"],
                            batch["labels"])
        outputs = model(inputs,
                         labels=targets)
        loss = outputs.loss
        accelerator.backward(loss)
        optimizer.step()
        lr_scheduler.step()
        optimizer.zero_grad()
```



From Accelerator to Trainer

Ability to Customize



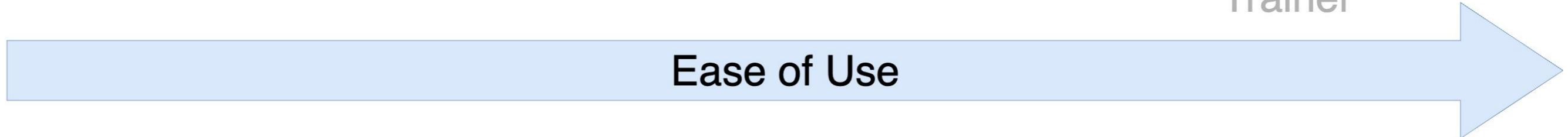
PyTorch



Accelerator



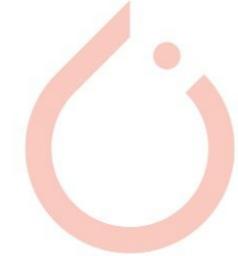
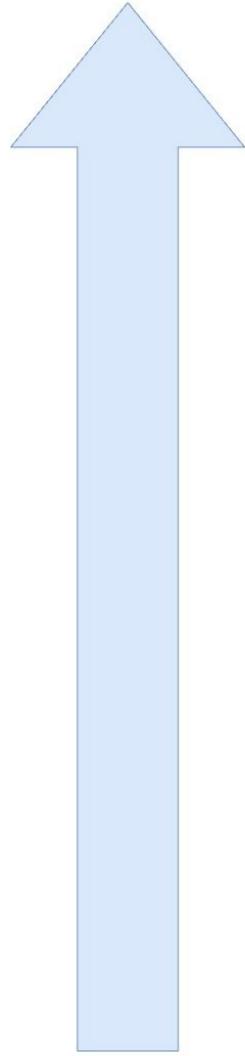
Trainer



Ease of Use

From Accelerator to Trainer

Ability to Customize



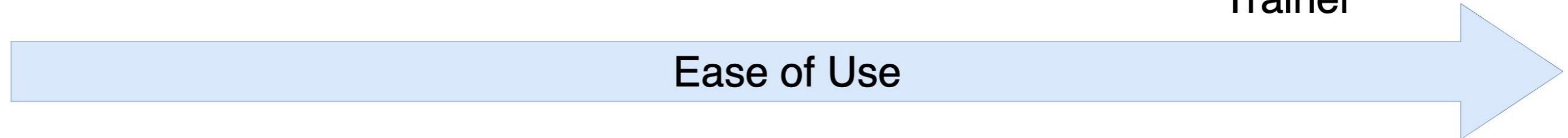
PyTorch



Accelerator



Trainer



Gradient accumulation with Trainer

```
training_args = TrainingArguments(output_dir="../results",
                                  evaluation_strategy="epoch",
                                  gradient_accumulation_steps=2)

trainer = Trainer(model=model,
                  args=training_args,
                  train_dataset=dataset["train"],
                  eval_dataset=dataset["validation"],
                  compute_metrics=compute_metrics)

trainer.train()
```

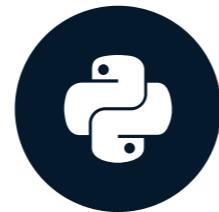
```
{'epoch': 1.0, 'eval_loss': 0.73, 'eval_accuracy': 0.03, 'eval_f1': 0.05}
{'epoch': 2.0, 'eval_loss': 0.68, 'eval_accuracy': 0.19, 'eval_f1': 0.25}
```

Let's practice!

EFFICIENT AI MODEL TRAINING WITH PYTORCH

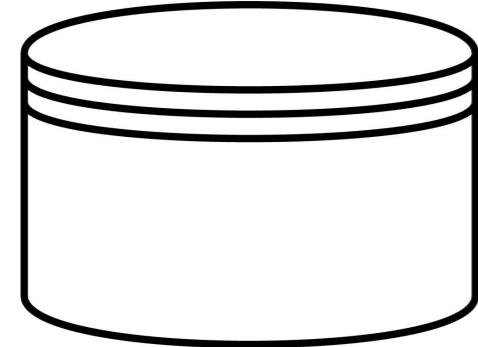
Gradient checkpointing and local SGD

EFFICIENT AI MODEL TRAINING WITH PYTORCH

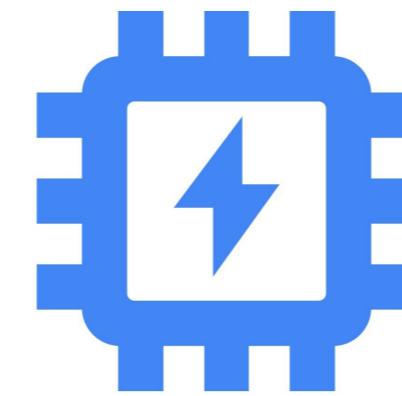


Dennis Lee
Data Engineer

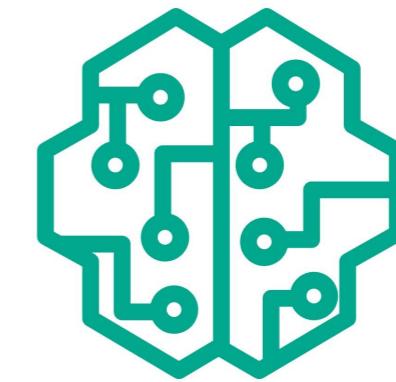
Improving training efficiency



Memory
Efficiency

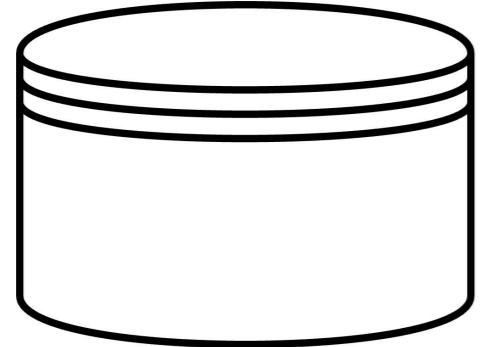


Communication
Efficiency

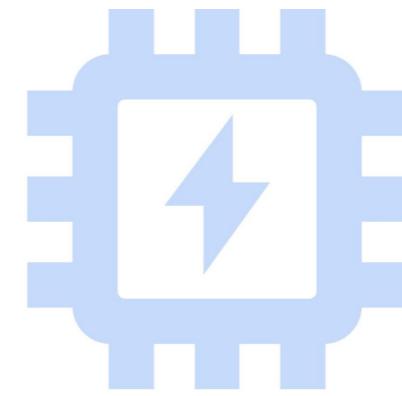


Computational
Efficiency

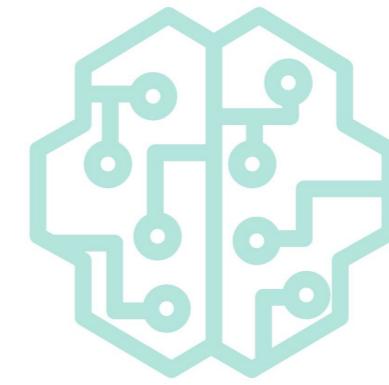
Gradient checkpointing improves memory efficiency



Memory
Efficiency

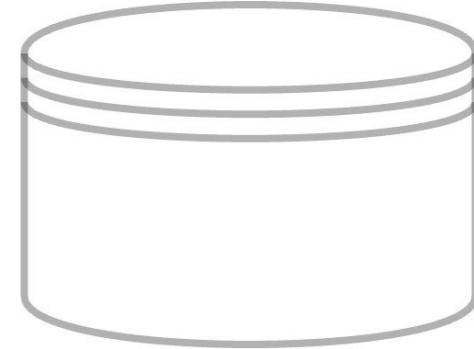


Communication
Efficiency

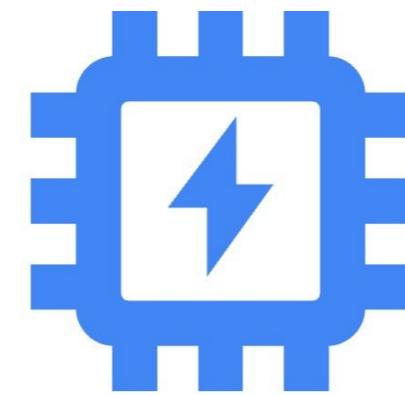


Computational
Efficiency

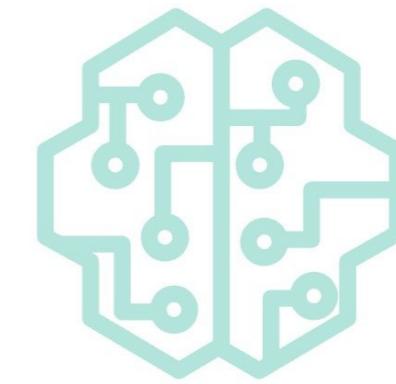
Local SGD addresses communication efficiency



Memory
Efficiency



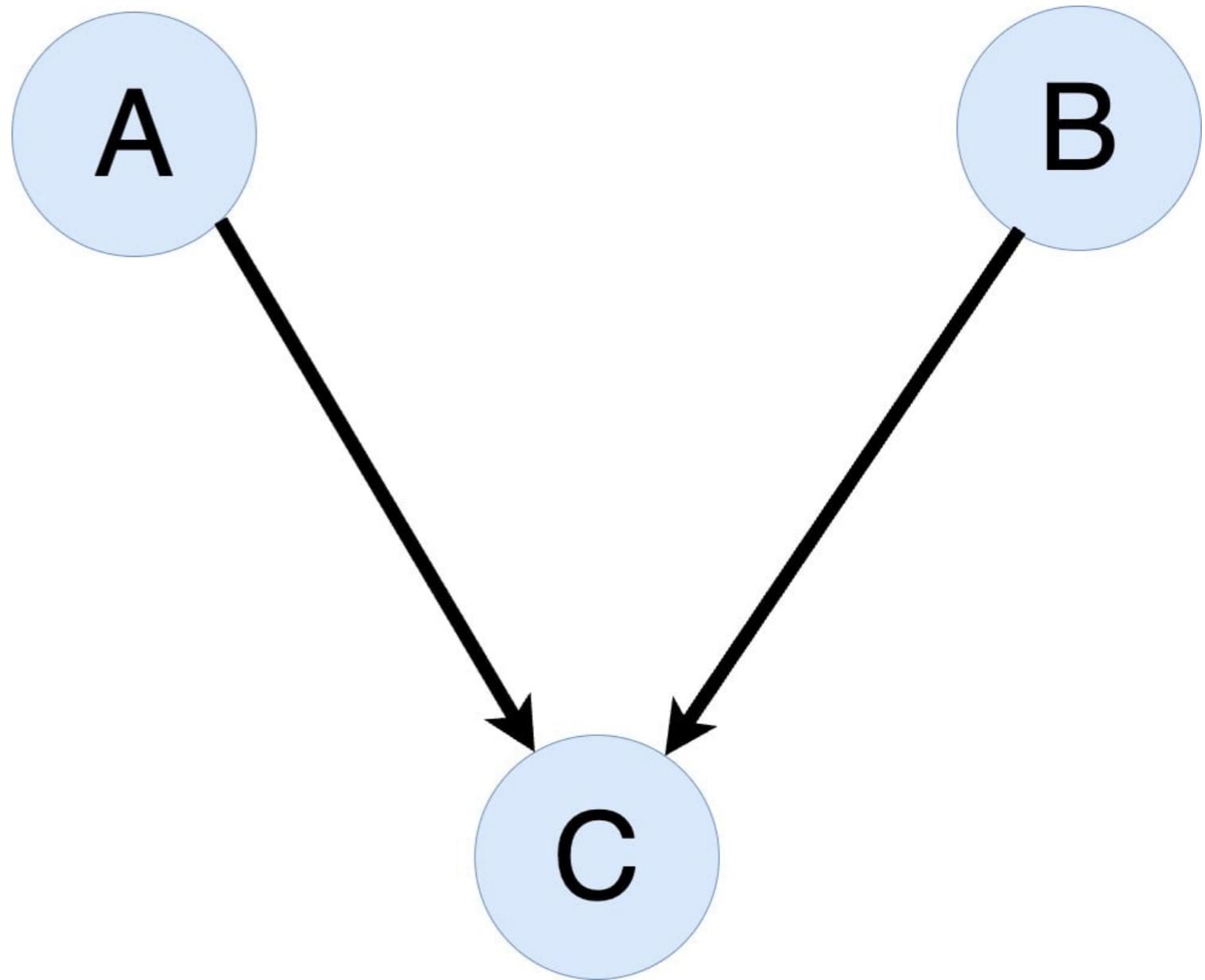
Communication
Efficiency



Computational
Efficiency

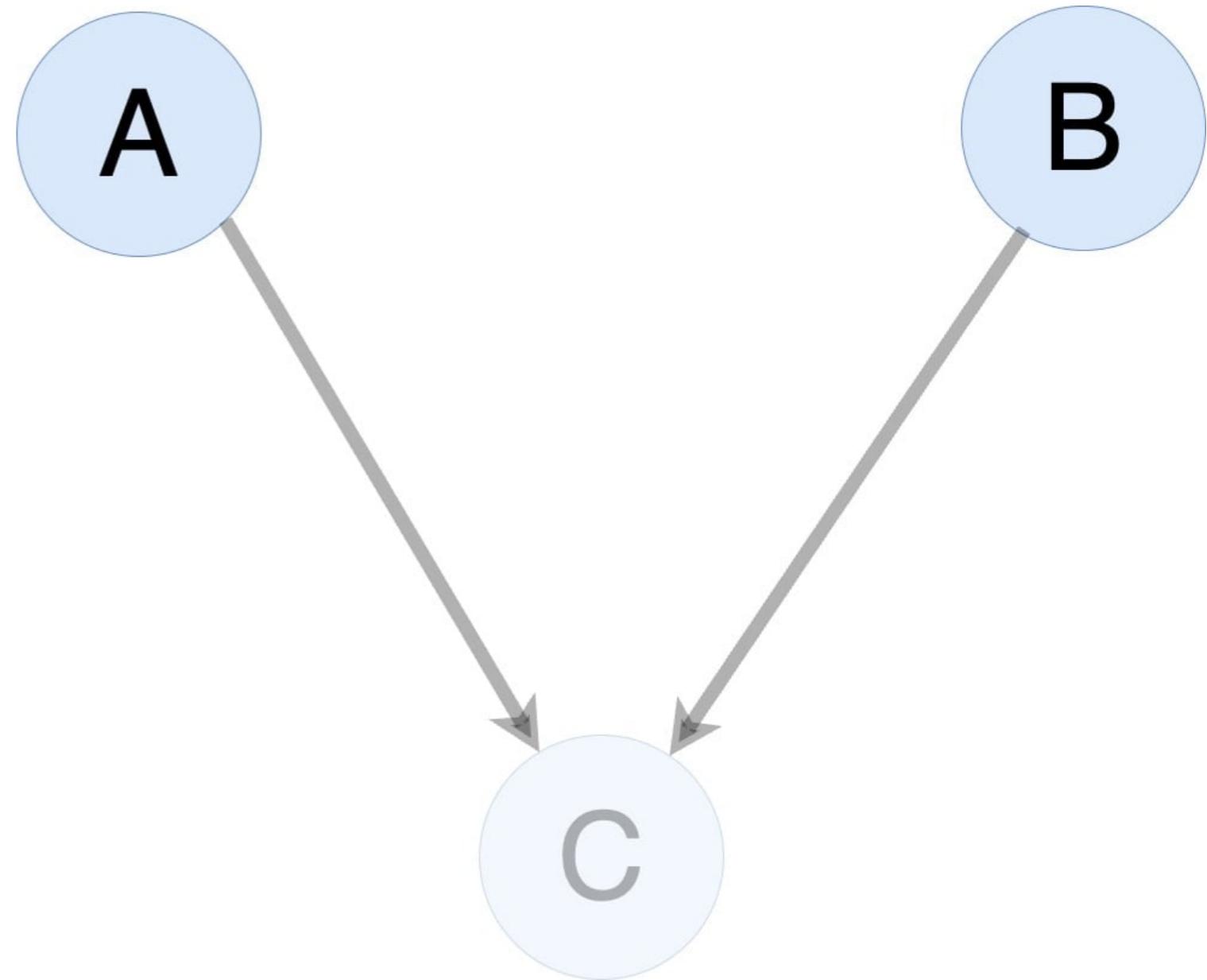
What is gradient checkpointing?

- Gradient checkpointing: reduce memory by selecting which activations to save
- Example: compute $A + B = C$



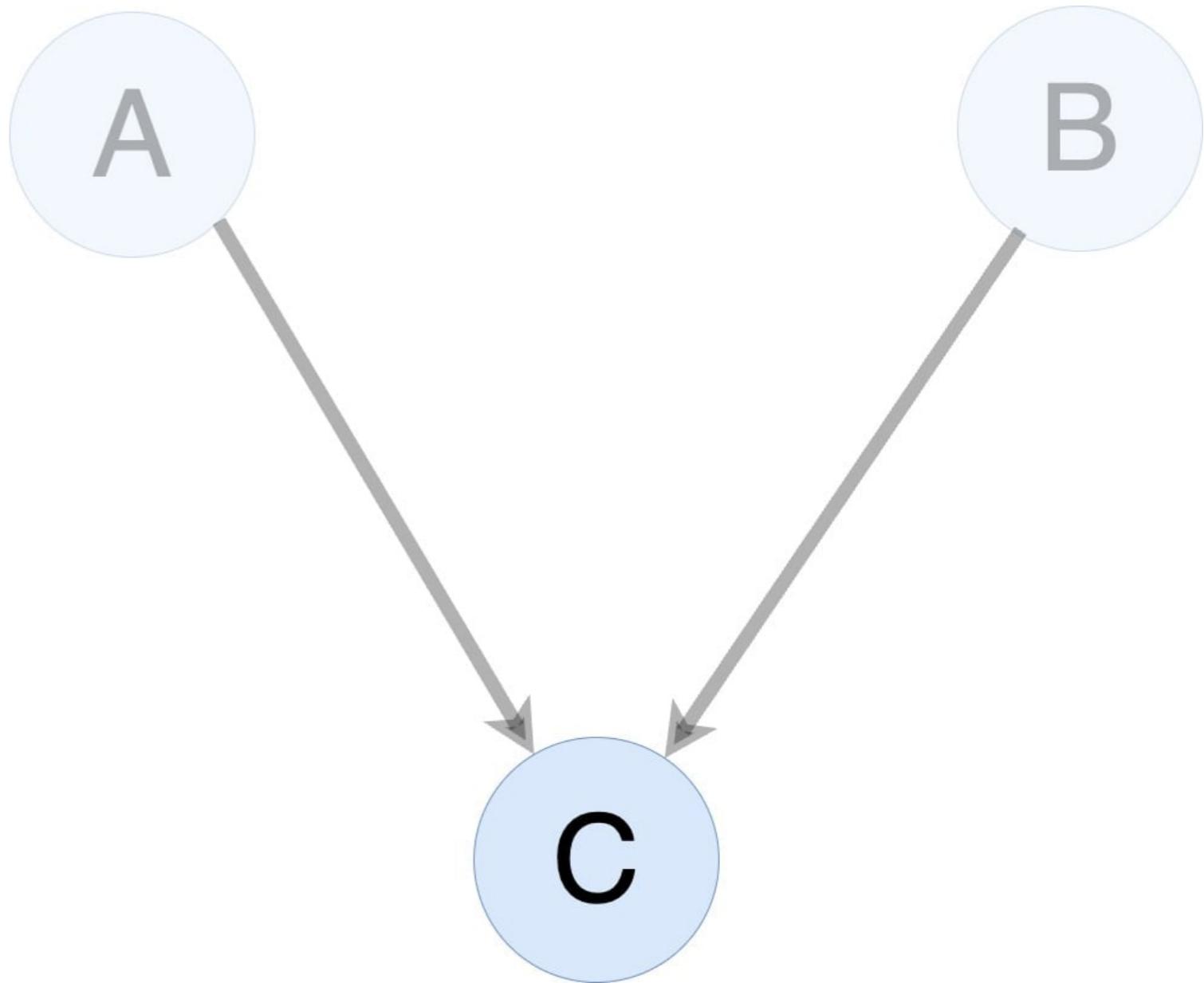
What is gradient checkpointing?

- Gradient checkpointing: reduce memory by selecting which activations to save
- Example: compute $A + B = C$
 - First compute A, B, then compute C



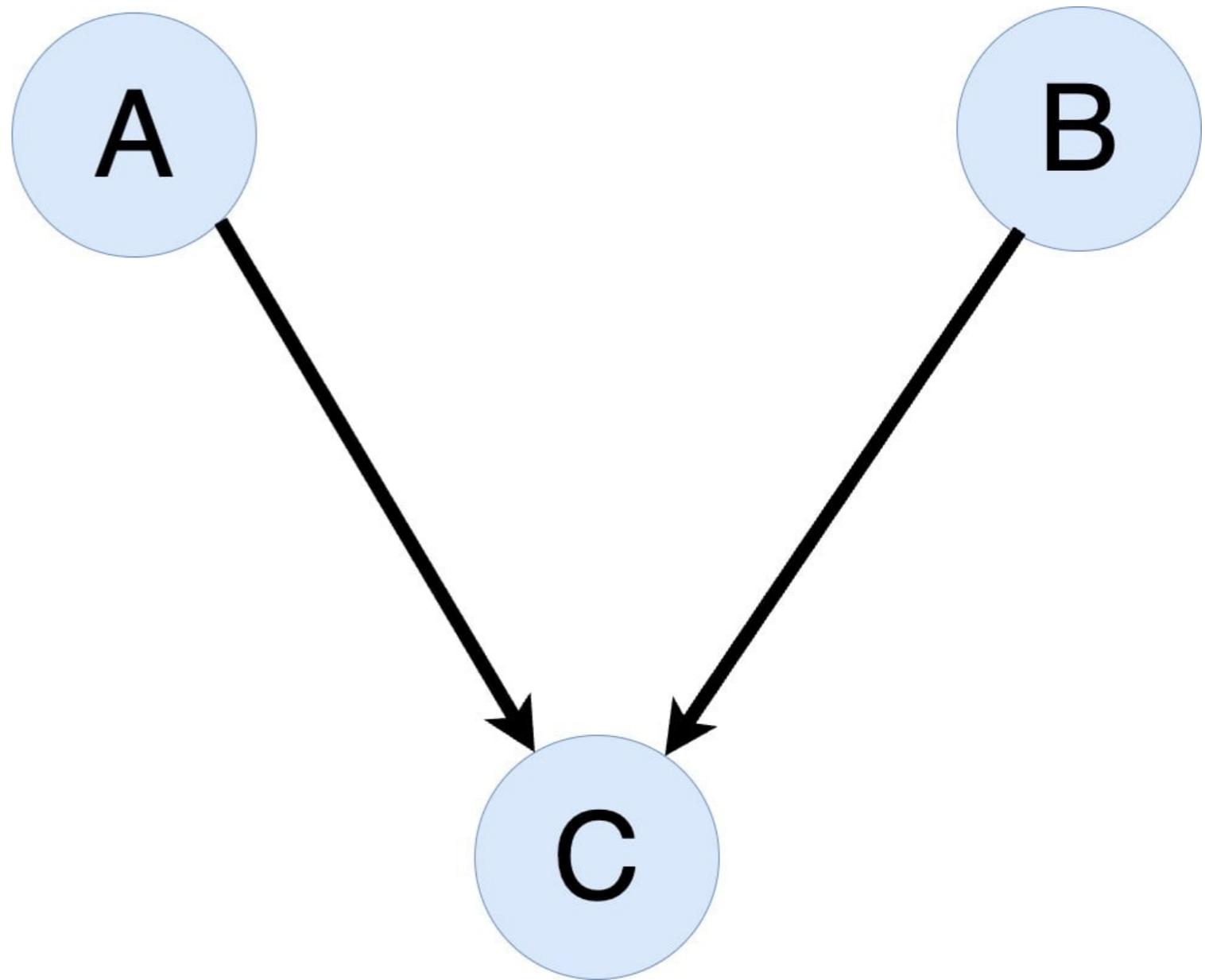
What is gradient checkpointing?

- Gradient checkpointing: reduce memory by selecting which activations to save
- Example: compute $A + B = C$
 - First compute A, B, then compute C
 - A, B not needed for rest of forward pass
- Should we save or remove A and B?



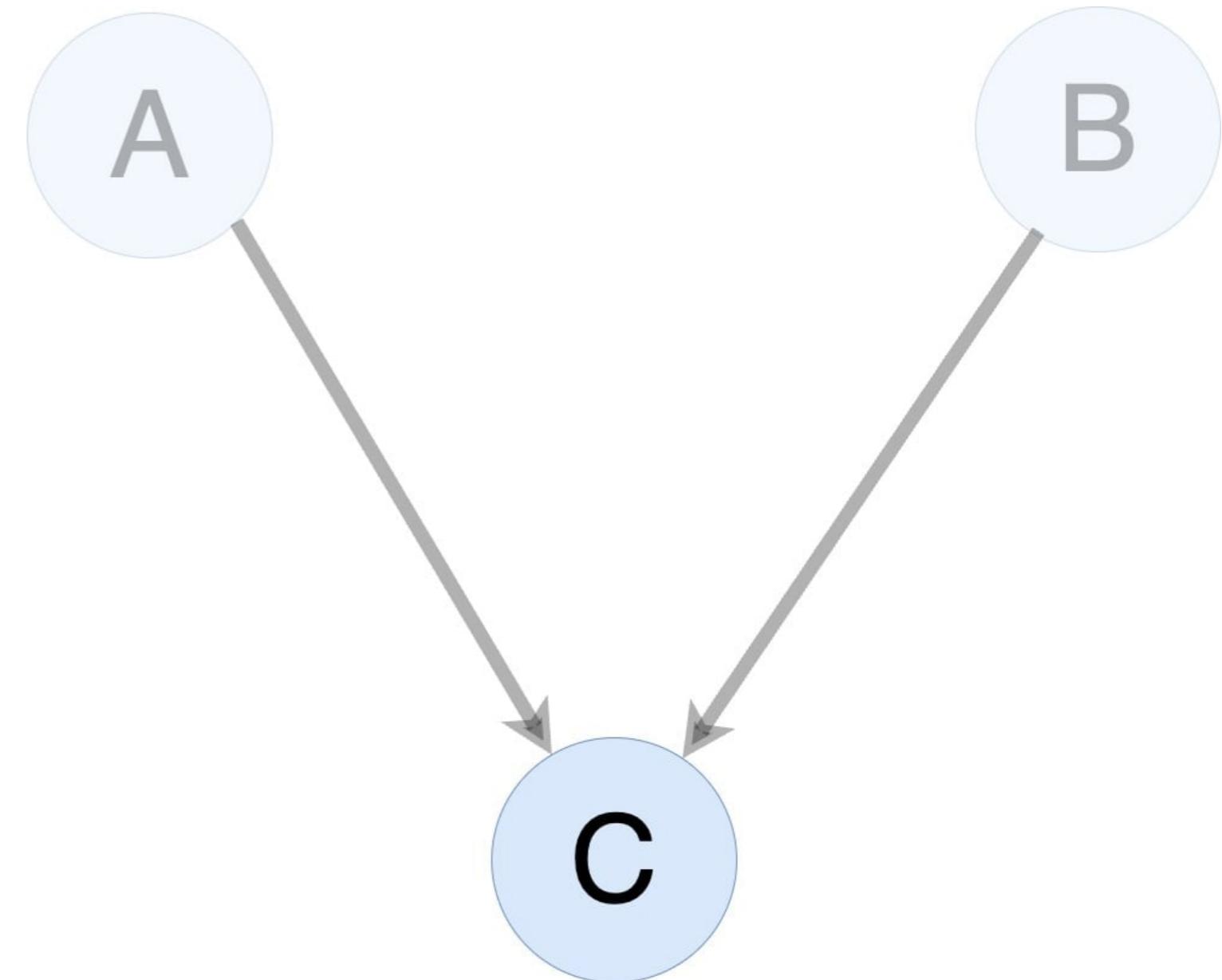
What is gradient checkpointing?

- Gradient checkpointing: reduce memory by selecting which activations to save
- Example: compute $A + B = C$
 - First compute A, B, then compute C
 - A, B not needed for rest of forward pass
- Should we save or remove A and B?
 - No gradient checkpointing: save A, B



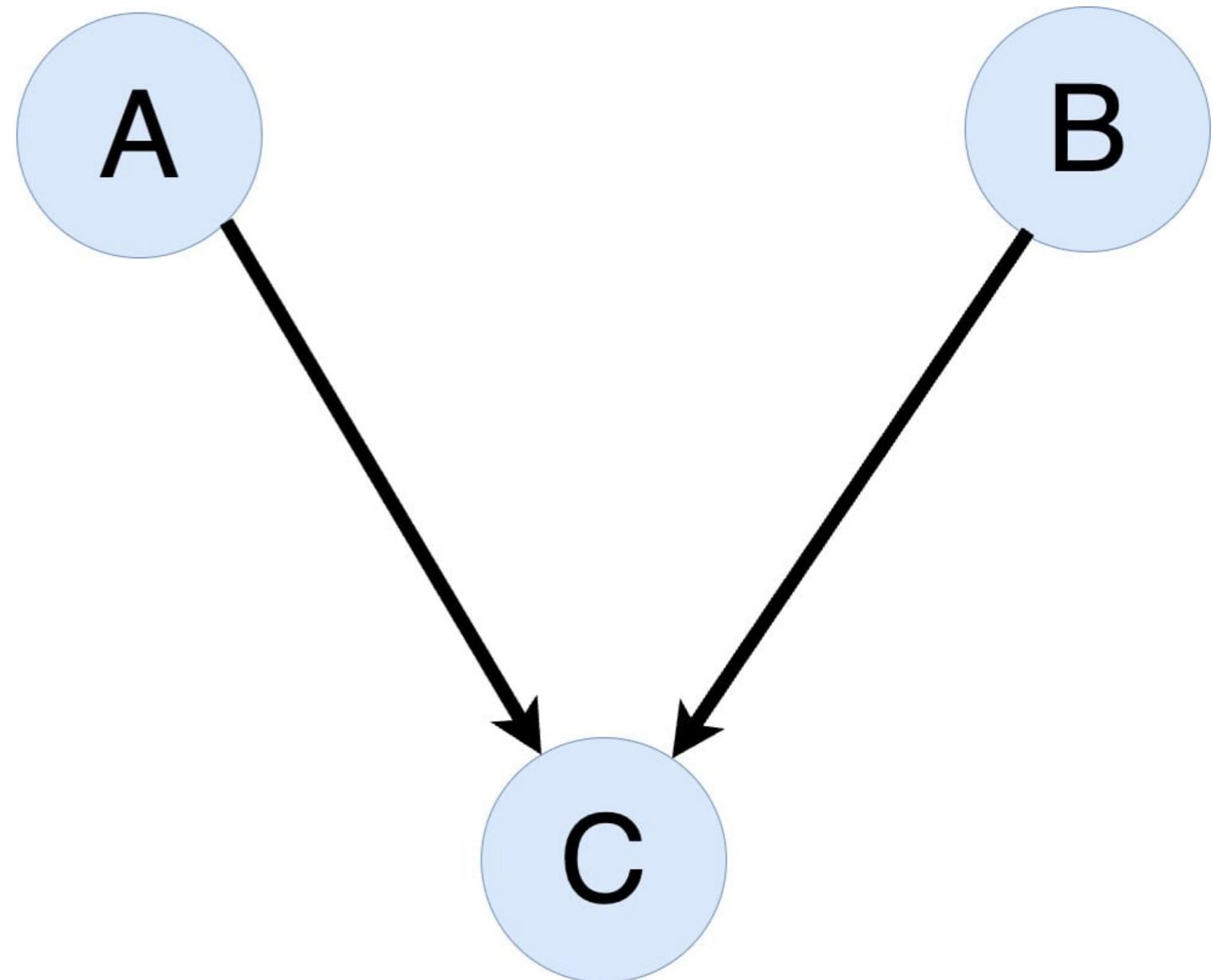
What is gradient checkpointing?

- Gradient checkpointing: reduce memory by selecting which activations to save
- Example: compute $A + B = C$
 - First compute A, B, then compute C
 - A, B not needed for rest of forward pass
- Should we save or remove A and B?
 - No gradient checkpointing: save A, B
 - Gradient checkpointing: remove A, B



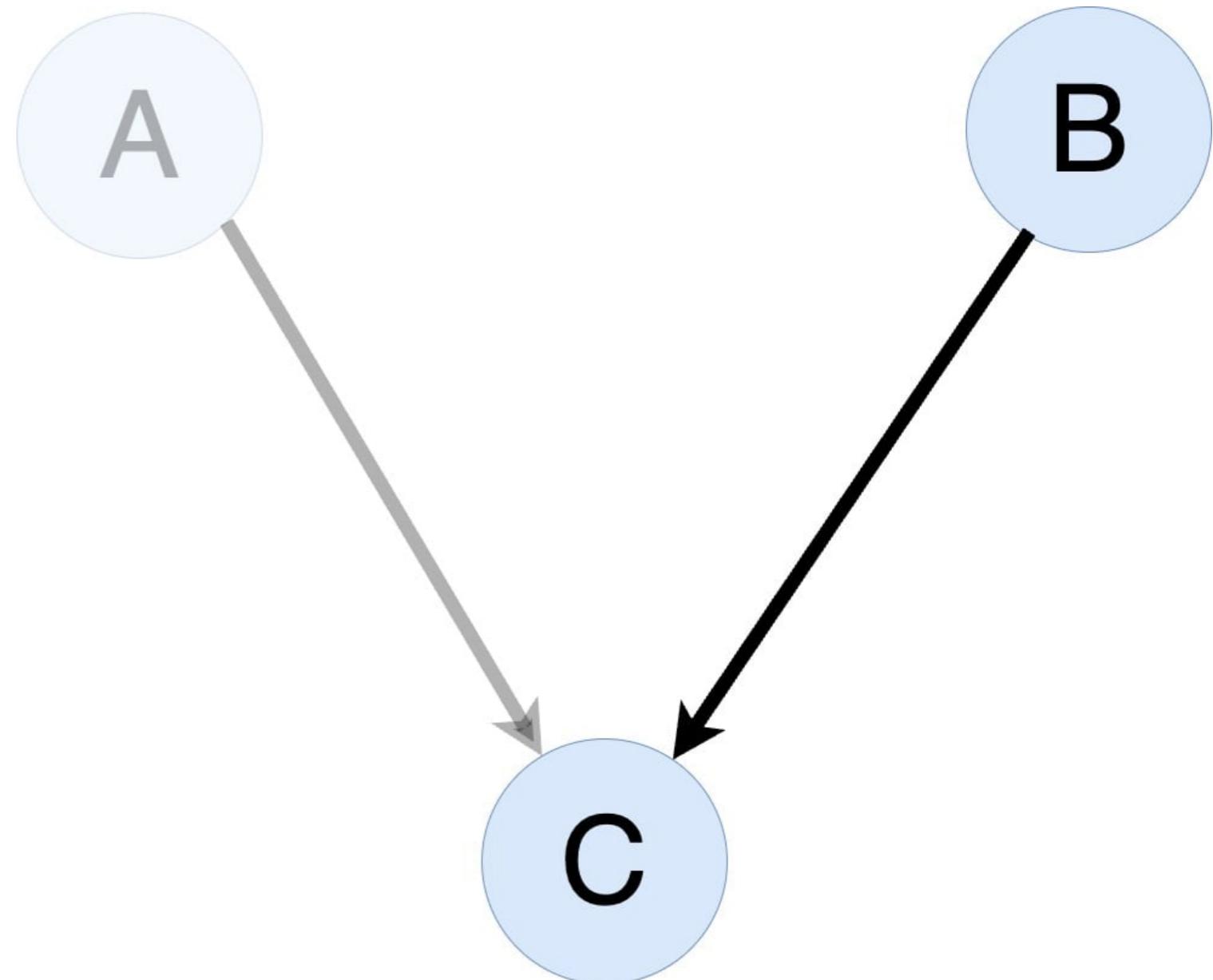
What is gradient checkpointing?

- Gradient checkpointing: reduce memory by selecting which activations to save
- Example: compute $A + B = C$
 - First compute A, B, then compute C
 - A, B not needed for rest of forward pass
- Should we save or remove A and B?
 - No gradient checkpointing: save A, B
 - Gradient checkpointing: remove A, B
 - Recompute A, B during backward pass



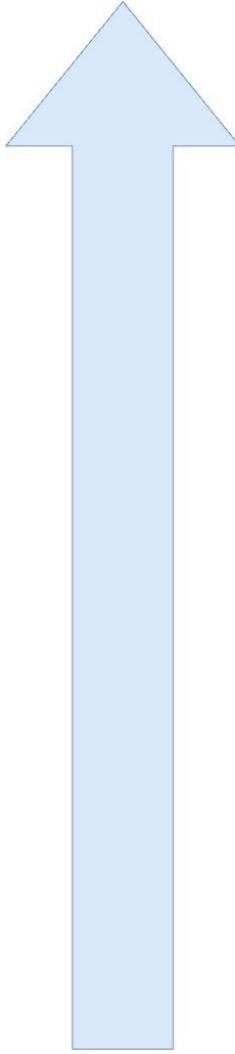
What is gradient checkpointing?

- Gradient checkpointing: reduce memory by selecting which activations to save
- Example: compute $A + B = C$
 - First compute A, B, then compute C
 - A, B not needed for rest of forward pass
- Should we save or remove A and B?
 - No gradient checkpointing: save A, B
 - Gradient checkpointing: remove A, B
 - Recompute A, B during backward pass
 - If B is expensive to recompute, save it



Trainer and Accelerator

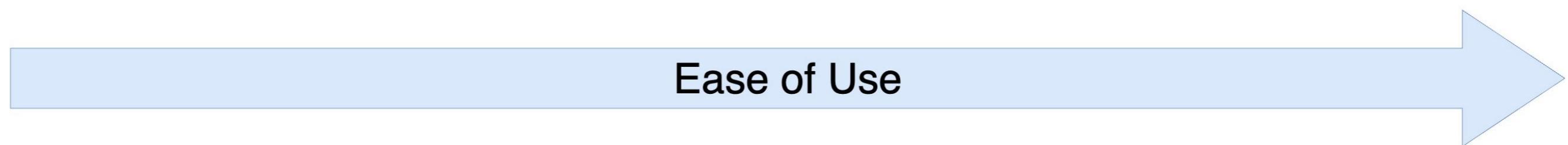
Ability to Customize



Accelerator



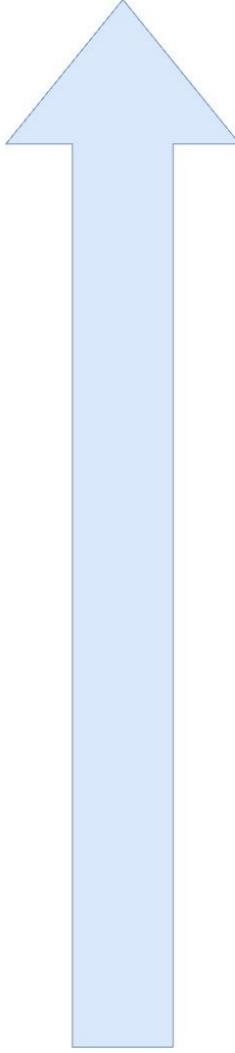
Trainer



Ease of Use

Trainer and Accelerator

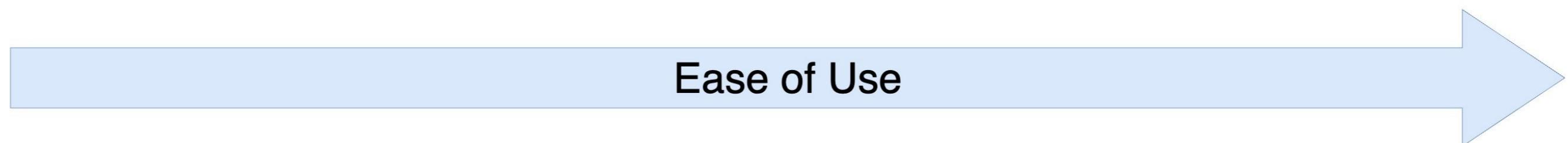
Ability to Customize



Accelerator



Trainer



Ease of Use

Gradient checkpointing with Trainer

```
training_args = TrainingArguments(output_dir="../results",  
                                  evaluation_strategy="epoch",  
                                  gradient_accumulation_steps=4)
```

Gradient checkpointing with Trainer

```
training_args = TrainingArguments(output_dir="../results",
                                  evaluation_strategy="epoch",
                                  gradient_accumulation_steps=4,
                                  gradient_checkpointing=True)

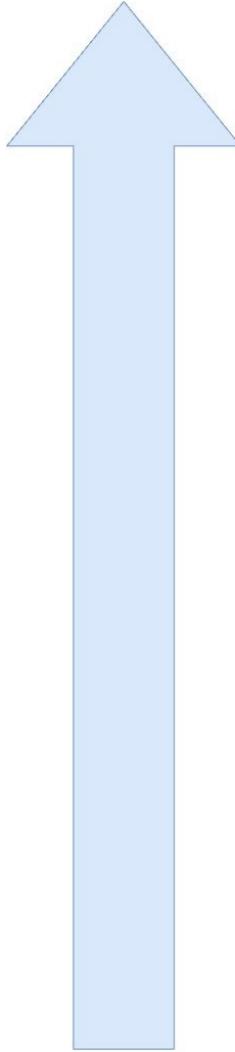
trainer = Trainer(model=model,
                   args=training_args,
                   train_dataset=dataset["train"],
                   eval_dataset=dataset["validation"],
                   compute_metrics=compute_metrics)

trainer.train()
```

```
{'epoch': 1.0, 'eval_loss': 0.73, 'eval_accuracy': 0.03, 'eval_f1': 0.05}
```

From Trainer to Accelerator

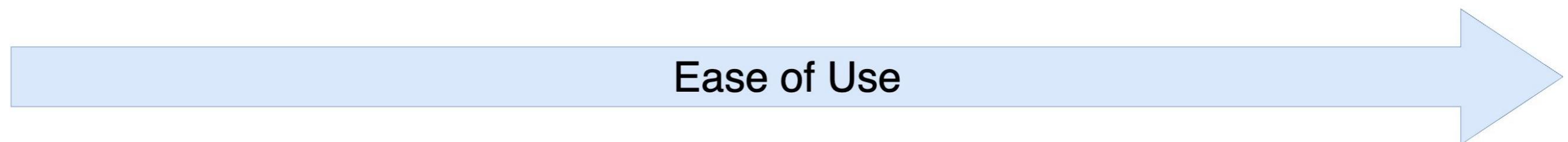
Ability to Customize



Accelerator



Trainer



Ease of Use

Gradient checkpointing with Accelerator

```
accelerator = Accelerator(gradient_accumulation_steps=2)

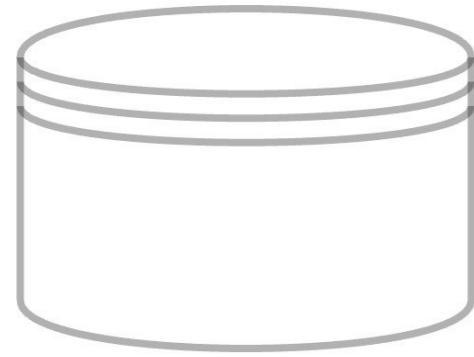
for index, batch in enumerate(dataloader):
    with accelerator.accumulate(model):
        inputs, targets = batch["input_ids"], batch["labels"]
        outputs = model(inputs, labels=targets)
        loss = outputs.loss
        accelerator.backward(loss)
        optimizer.step()
        lr_scheduler.step()
        optimizer.zero_grad()
```

Gradient checkpointing with Accelerator

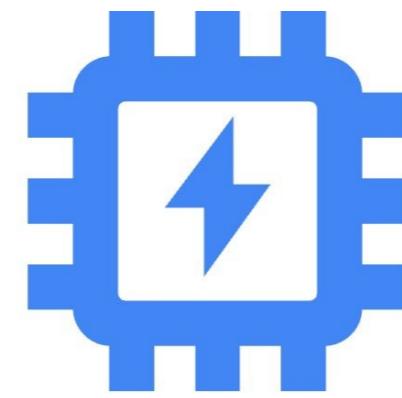
```
accelerator = Accelerator(gradient_accumulation_steps=2)
model.gradient_checkpointing_enable()

for index, batch in enumerate(dataloader):
    with accelerator.accumulate(model):
        inputs, targets = batch["input_ids"], batch["labels"]
        outputs = model(inputs, labels=targets)
        loss = outputs.loss
        accelerator.backward(loss)
        optimizer.step()
        lr_scheduler.step()
        optimizer.zero_grad()
```

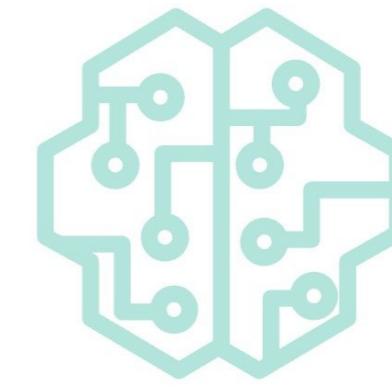
Local SGD improves communication efficiency



Memory
Efficiency

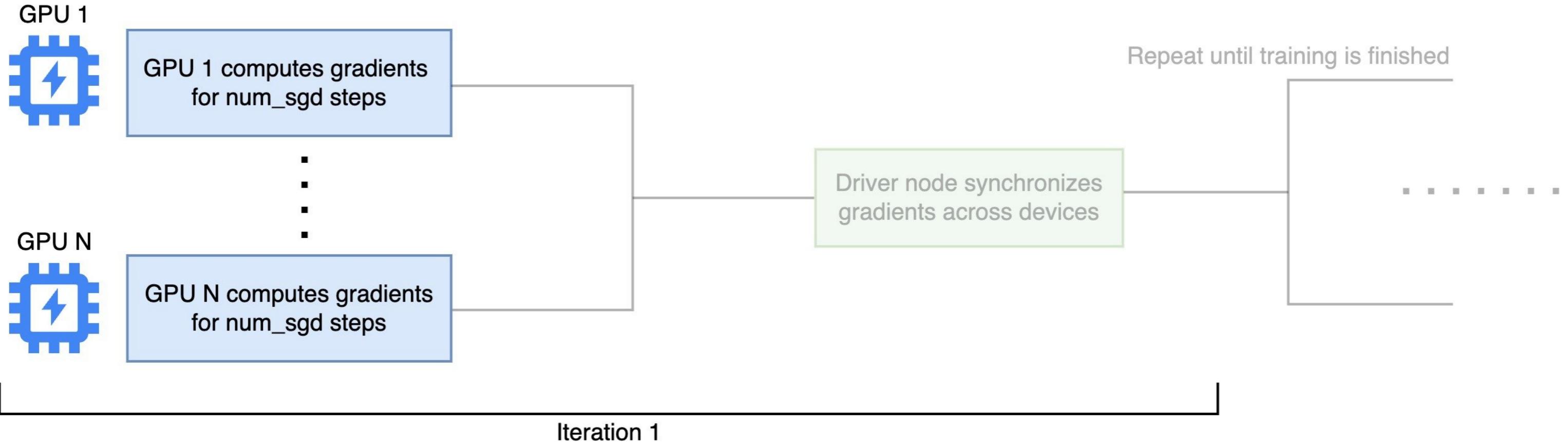


Communication
Efficiency



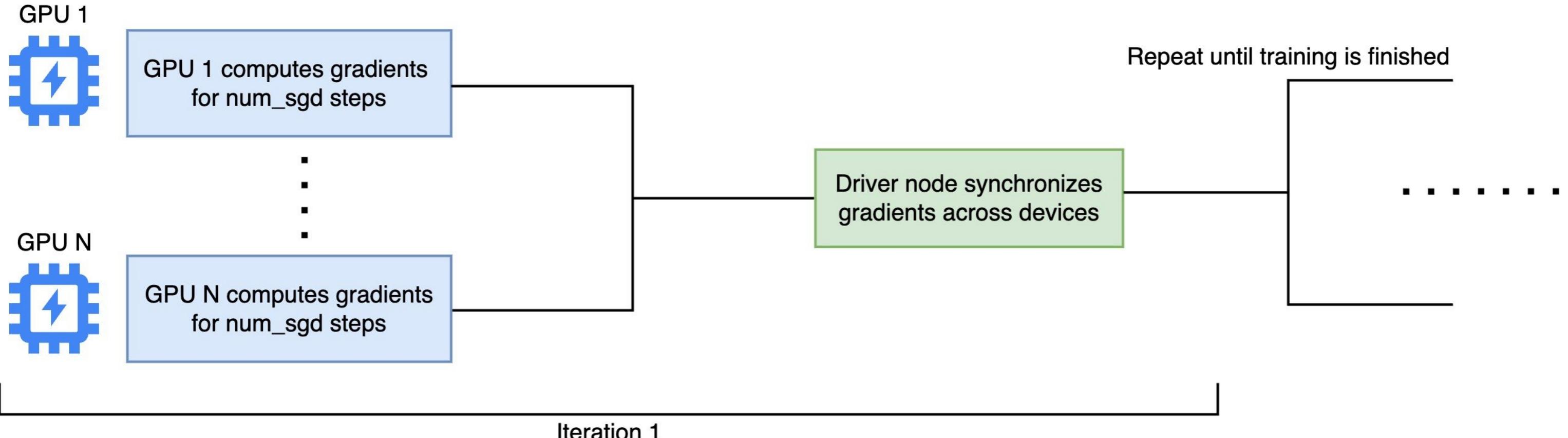
Computational
Efficiency

What is local SGD?



- Each device computes gradients in parallel

What is local SGD?



- Each device computes gradients in parallel
- Gradient synchronization: Driver node updates model parameters on each device
- Local SGD: Reduce frequency of gradient synchronization

Local SGD with Accelerator

```
for index, batch in enumerate(dataloader):
    with accelerator.accumulate(model):
        inputs, targets = batch["input_ids"], batch["labels"]
        outputs = model(inputs, labels=targets)
        loss = outputs.loss
        accelerator.backward(loss)
        optimizer.step()
        lr_scheduler.step()
        optimizer.zero_grad()
```

Local SGD with Accelerator

```
from accelerate.local_sgd import LocalSGD

with LocalSGD(accelerator=accelerator, model=model, local_sgd_steps=8,
              enabled=True) as local_sgd:
    for index, batch in enumerate(dataloader):
        with accelerator.accumulate(model):
            inputs, targets = batch["input_ids"], batch["labels"]
            outputs = model(inputs, labels=targets)
            loss = outputs.loss
            accelerator.backward(loss)
            optimizer.step()
            lr_scheduler.step()
            optimizer.zero_grad()
```

Local SGD with Accelerator

```
from accelerate.local_sgd import LocalSGD

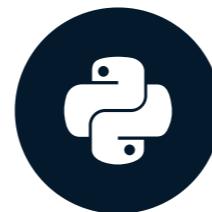
with LocalSGD(accelerator=accelerator, model=model, local_sgd_steps=8,
              enabled=True) as local_sgd:
    for index, batch in enumerate(dataloader):
        with accelerator.accumulate(model):
            inputs, targets = batch["input_ids"], batch["labels"]
            outputs = model(inputs, labels=targets)
            loss = outputs.loss
            accelerator.backward(loss)
            optimizer.step()
            lr_scheduler.step()
            optimizer.zero_grad()
            local_sgd.step()
```

Let's practice!

EFFICIENT AI MODEL TRAINING WITH PYTORCH

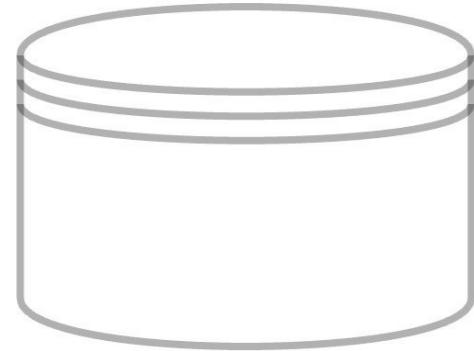
Mixed precision training

EFFICIENT AI MODEL TRAINING WITH PYTORCH

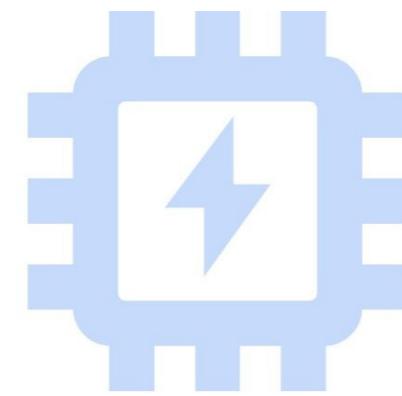


Dennis Lee
Data Engineer

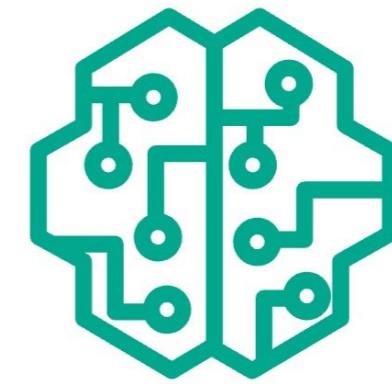
Mixed precision training accelerates computation



Memory
Efficiency



Communication
Efficiency

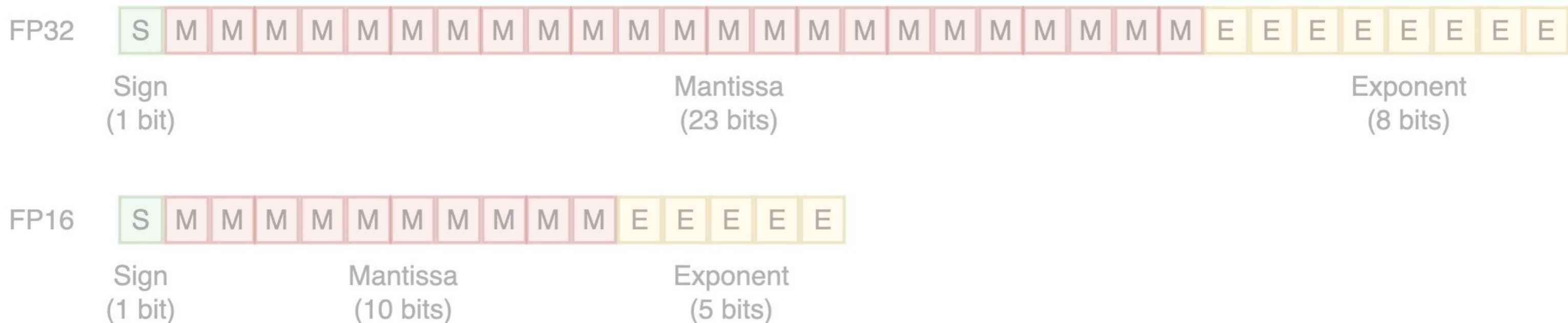


Computational
Efficiency

Faster calculations with less precision

$$+ \boxed{1.25} \times \boxed{10^{-1}}$$

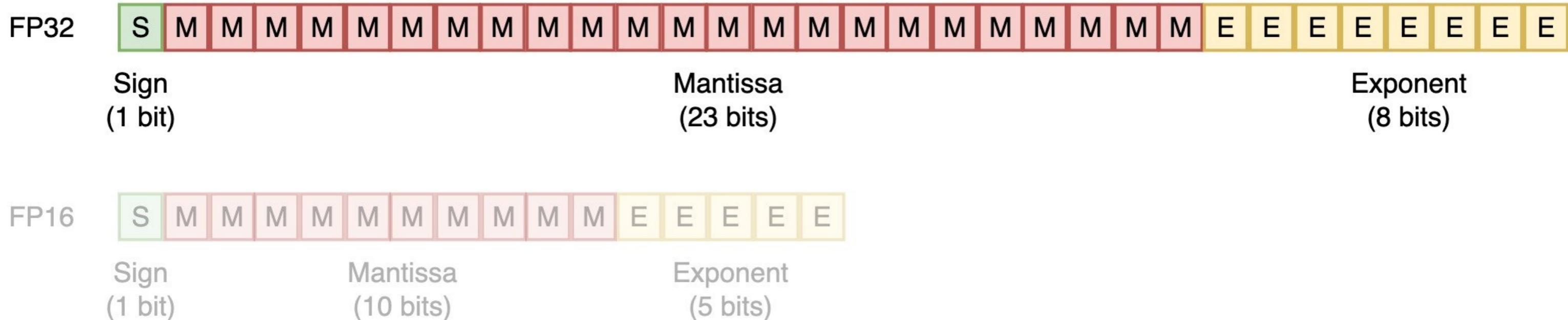
Sign Mantissa Exponent



Faster calculations with less precision

$$+ \boxed{1.25} \times \boxed{10^{-1}}$$

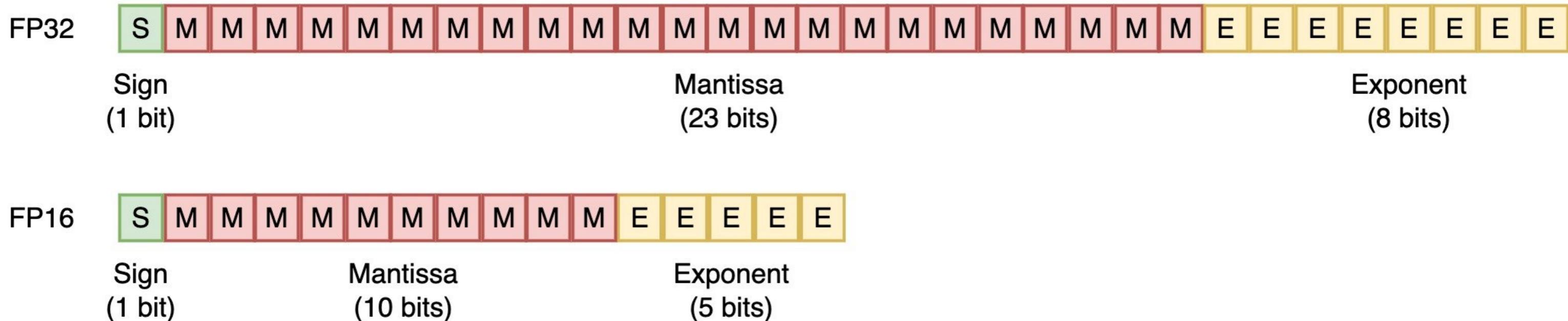
Sign Mantissa Exponent



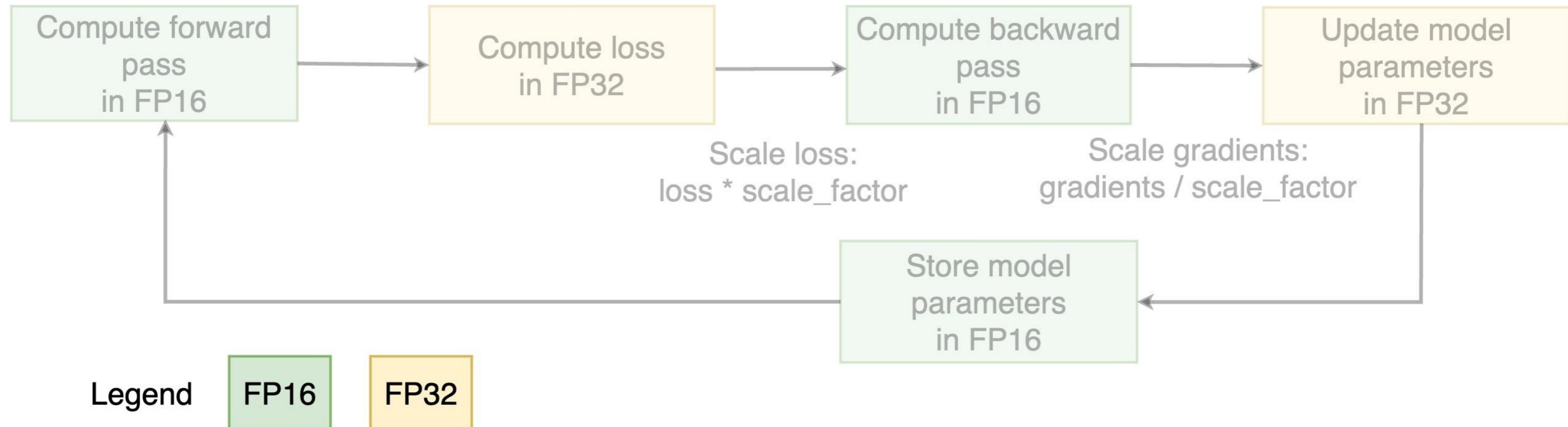
Faster calculations with less precision

+ 1.25×10^{-1}

Sign Mantissa Exponent

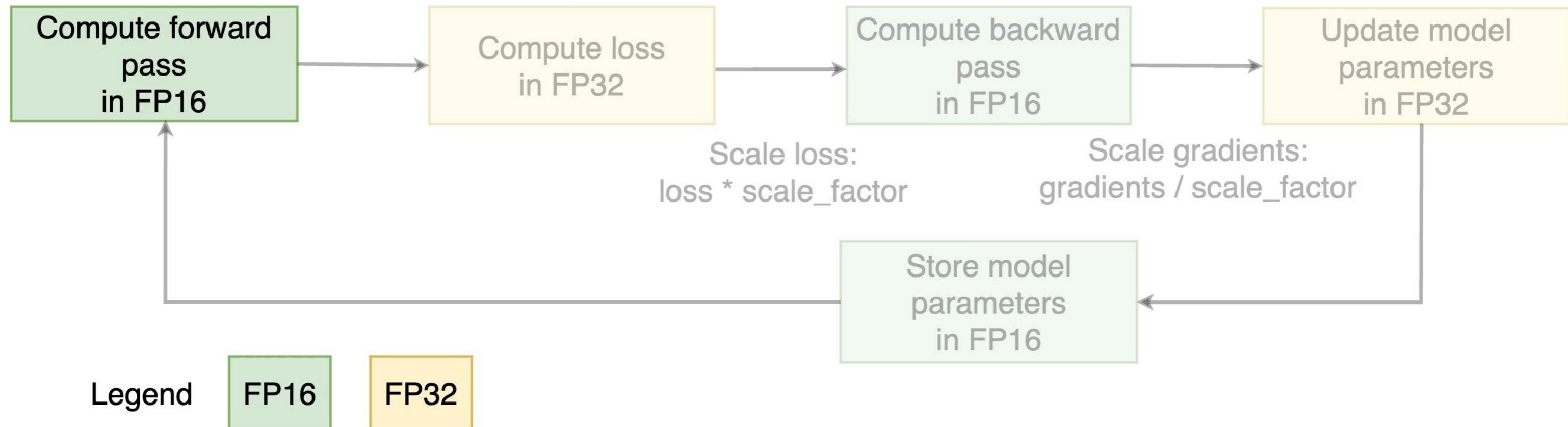


What is mixed precision training?



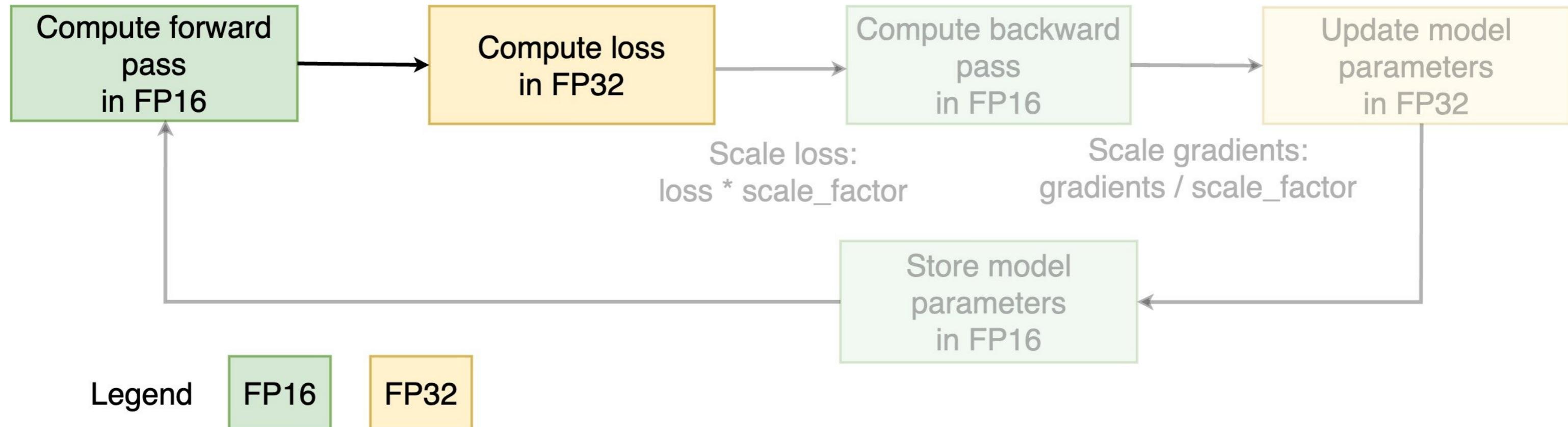
- Mixed precision training: combine FP16, FP32 computations to speed up training

What is mixed precision training?



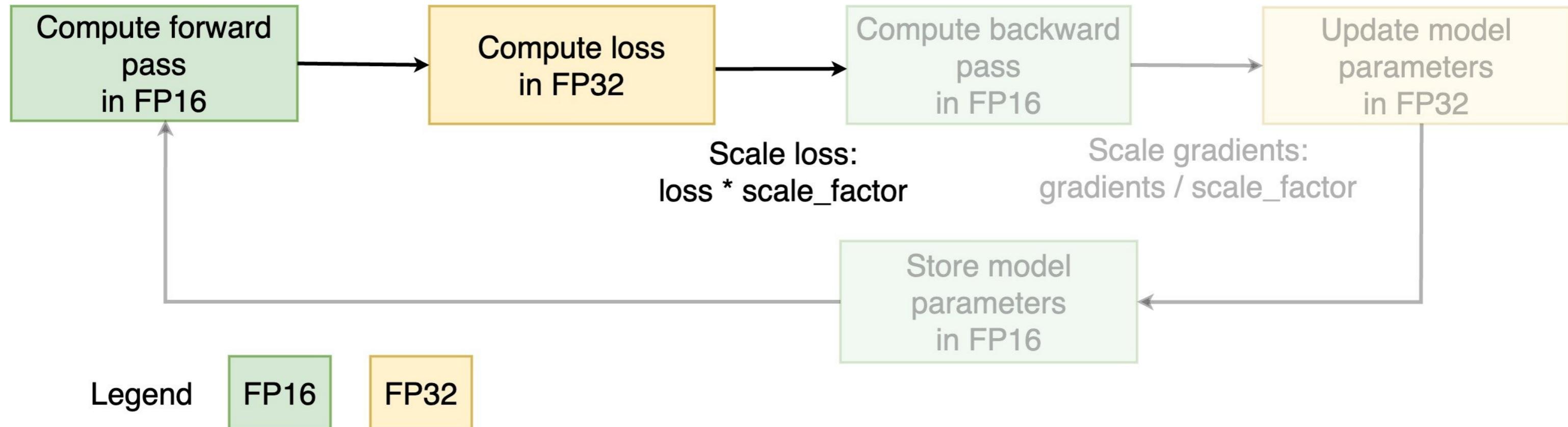
- Mixed precision training: combine FP16, FP32 computations to speed up training

What is mixed precision training?



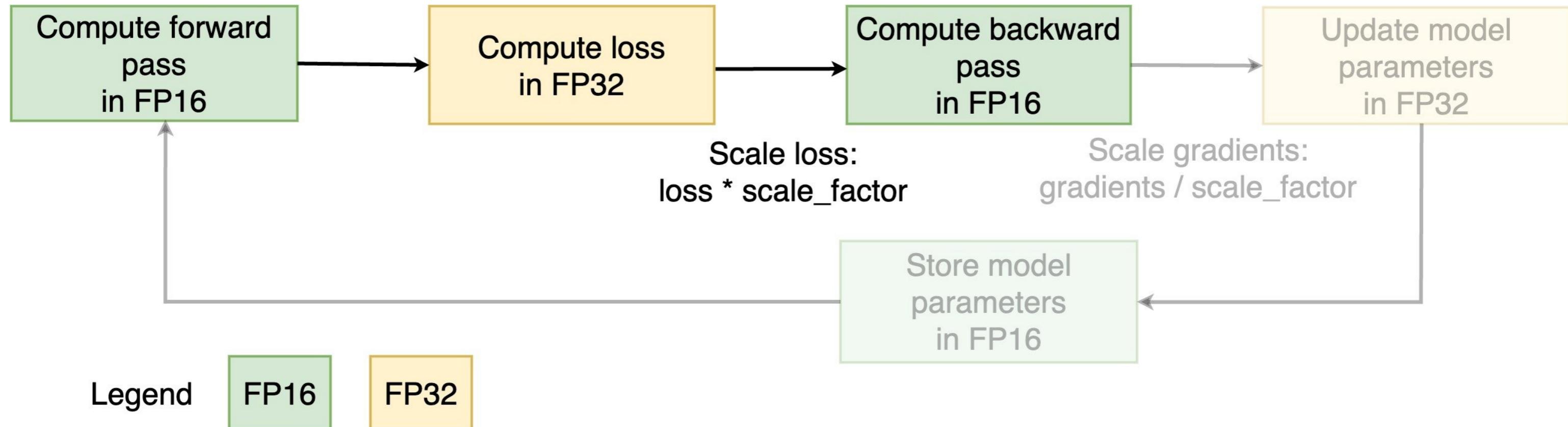
- Mixed precision training: combine FP16, FP32 computations to speed up training
- Underflow: number vanishes to 0 because it falls below precision

What is mixed precision training?



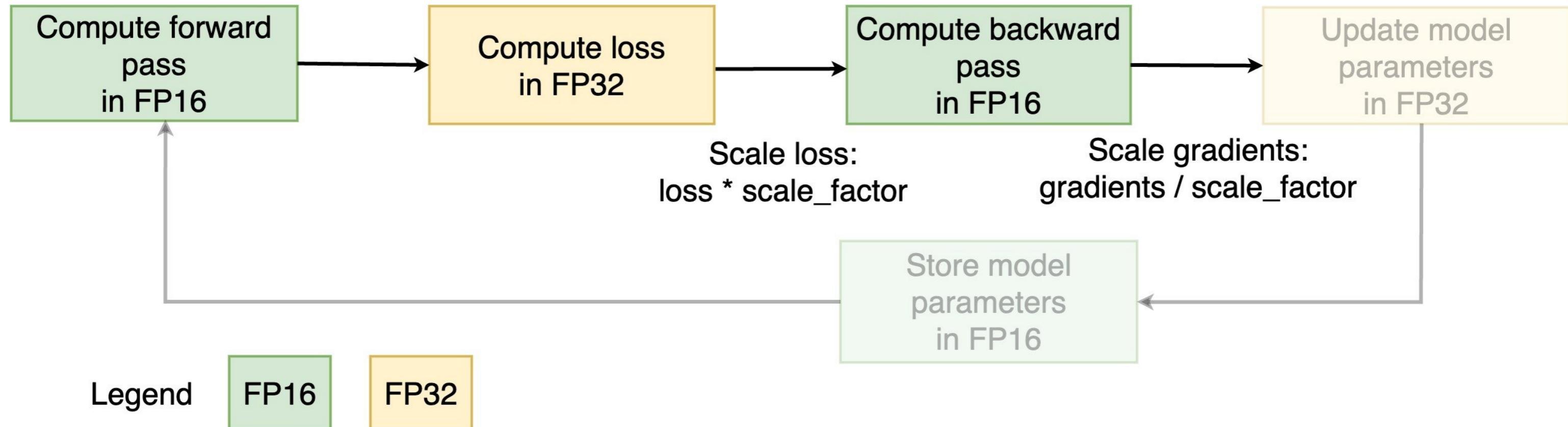
- Mixed precision training: combine FP16, FP32 computations to speed up training
- Underflow: number vanishes to 0 because it falls below precision
- Scale loss to prevent underflow

What is mixed precision training?



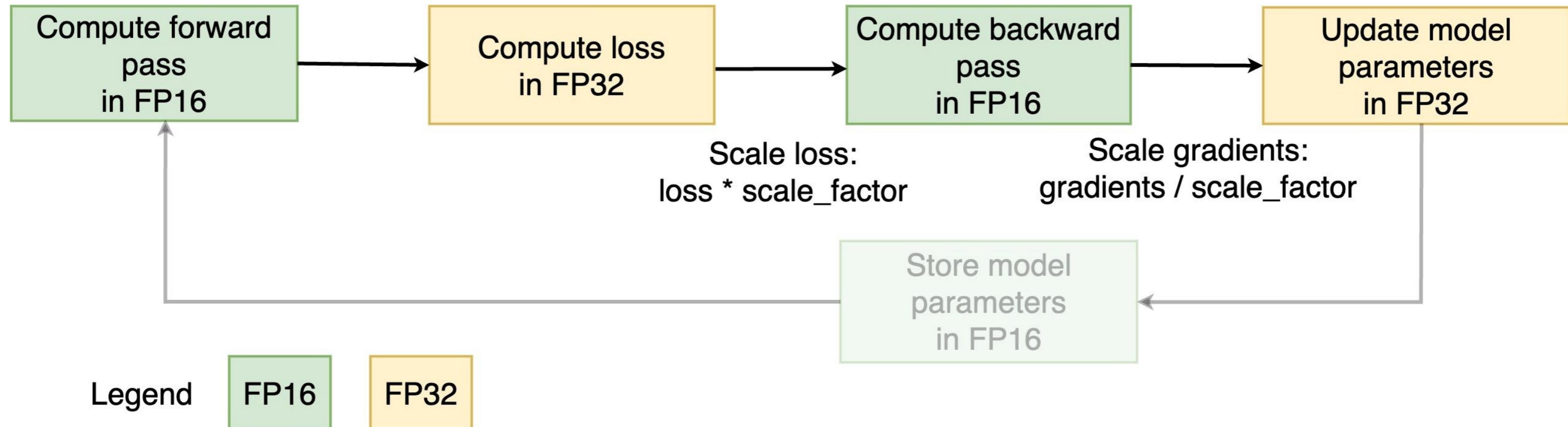
- Mixed precision training: combine FP16, FP32 computations to speed up training
- Underflow: number vanishes to 0 because it falls below precision
- Scale loss to prevent underflow

What is mixed precision training?



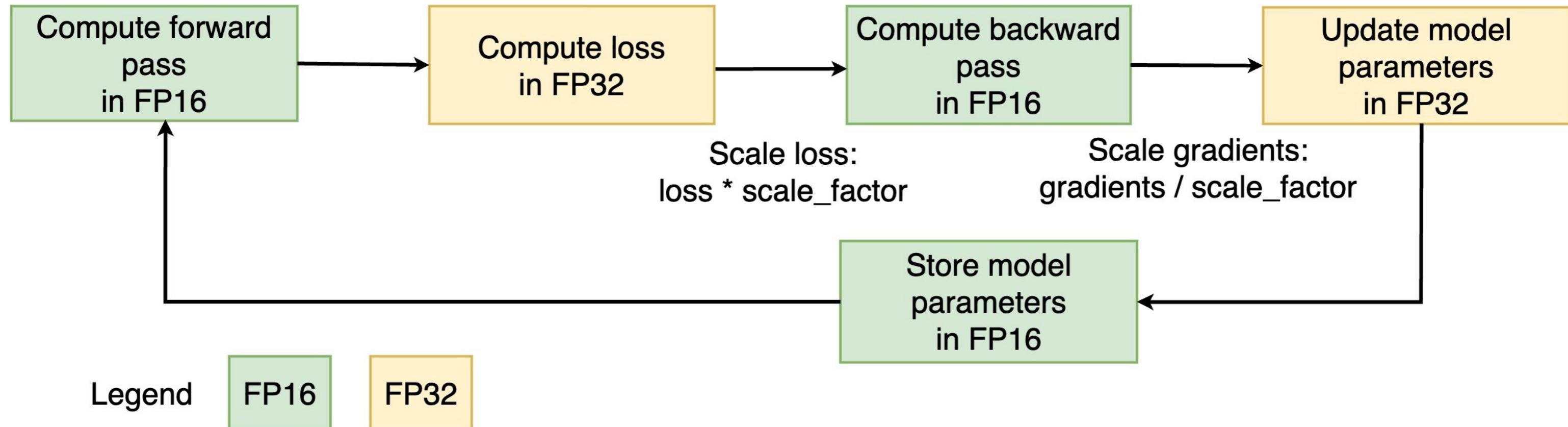
- Mixed precision training: combine FP16, FP32 computations to speed up training
- Underflow: number vanishes to 0 because it falls below precision
- Scale loss to prevent underflow

What is mixed precision training?



- Mixed precision training: combine FP16, FP32 computations to speed up training
- Underflow: number vanishes to 0 because it falls below precision
- Scale loss to prevent underflow

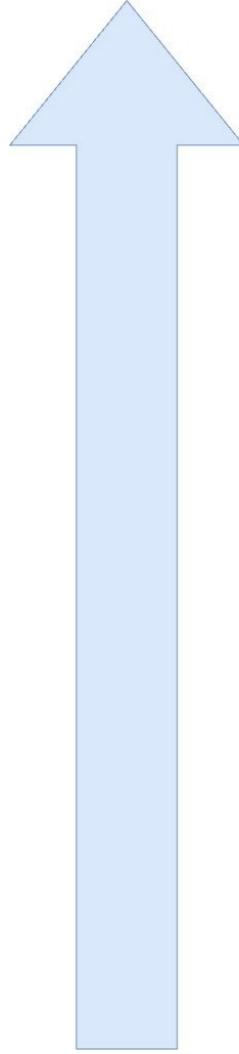
What is mixed precision training?



- Mixed precision training: combine FP16, FP32 computations to speed up training
- Underflow: number vanishes to 0 because it falls below precision
- Scale loss to prevent underflow

PyTorch implementation

Ability to Customize



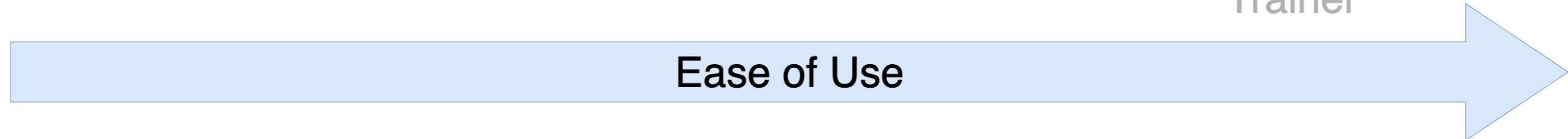
PyTorch



Accelerator



Trainer



Ease of Use

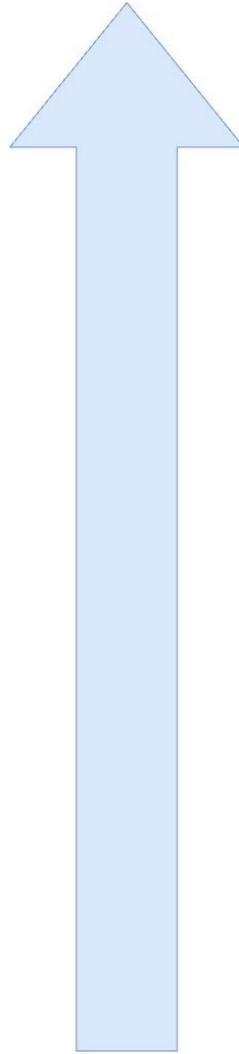
Mixed precision training with PyTorch

```
scaler = torch.amp.GradScaler()

for batch in train_dataloader:
    inputs, targets = batch["input_ids"], batch["labels"]
    with torch.autocast(device_type="cpu", dtype=torch.float16):
        outputs = model(inputs, labels=targets)
        loss = outputs.loss
    scaler.scale(loss).backward()
    scaler.step(optimizer)
    scaler.update()
    optimizer.zero_grad()
```

From PyTorch to Accelerator

Ability to Customize



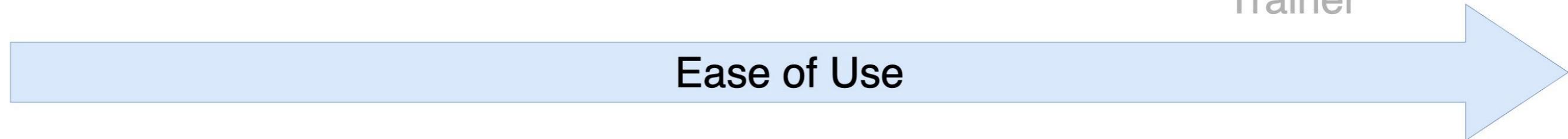
PyTorch



Accelerator



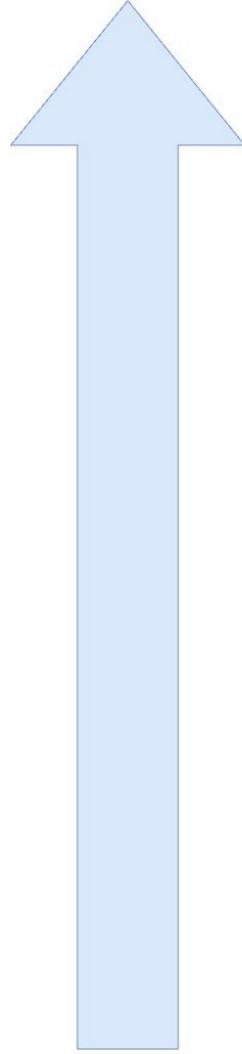
Trainer



Ease of Use

From PyTorch to Accelerator

Ability to Customize



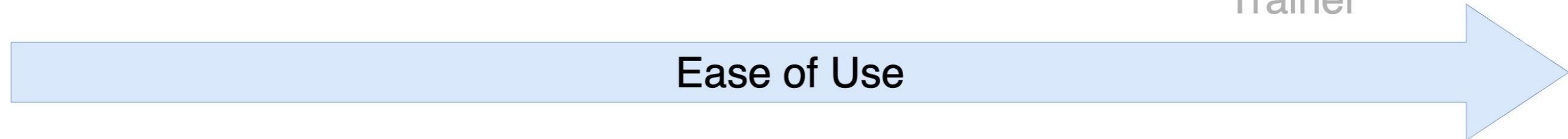
PyTorch



Accelerator



Trainer



Ease of Use

Mixed precision training with Accelerator

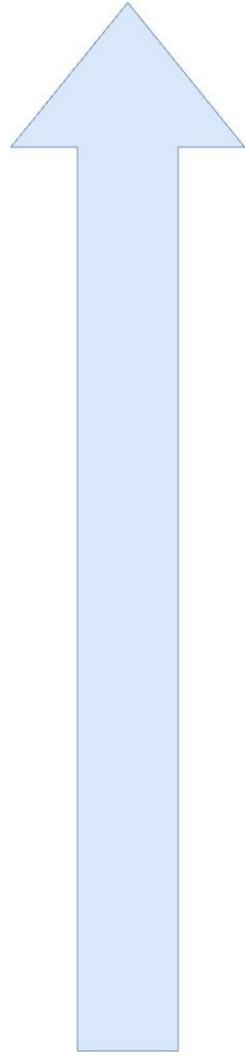
```
accelerator = Accelerator(mixed_precision="fp16")

model, optimizer, train_dataloader, lr_scheduler = \
    accelerator.prepare(model, optimizer, train_dataloader, lr_scheduler)

for batch in train_dataloader:
    inputs, targets = batch["input_ids"], batch["labels"]
    outputs = model(inputs, labels=targets)
    loss = outputs.loss
    accelerator.backward(loss)
    optimizer.step()
    optimizer.zero_grad()
```

From Accelerator to Trainer

Ability to Customize



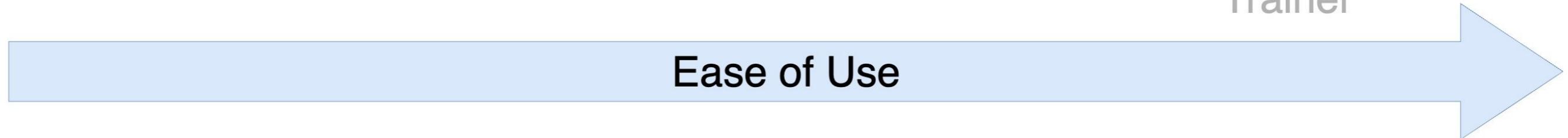
PyTorch



Accelerator



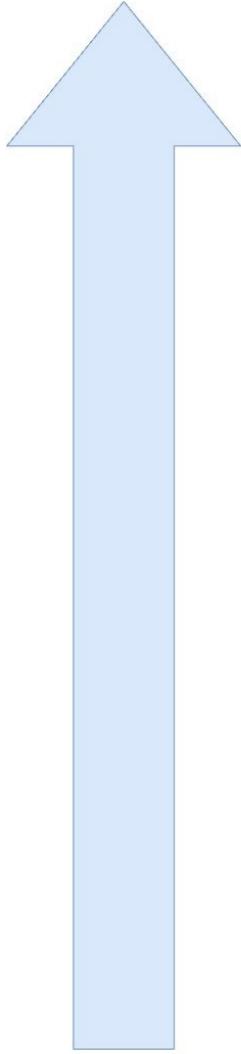
Trainer



Ease of Use

From Accelerator to Trainer

Ability to Customize



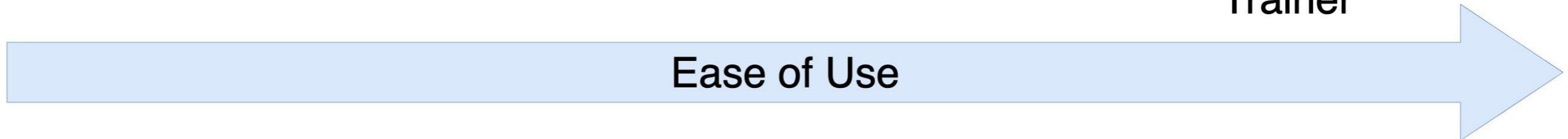
PyTorch



Accelerator



Trainer



Ease of Use

Mixed precision training with Trainer

```
training_args = TrainingArguments(  
    output_dir="../results",  
    evaluation_strategy="epoch",  
    fp16=True  
)  
trainer = Trainer(  
    model=model,  
    args=training_args,  
    train_dataset=dataset["train"],  
    eval_dataset=dataset["validation"],  
    compute_metrics=compute_metrics,  
)  
trainer.train()
```

Let's practice!

EFFICIENT AI MODEL TRAINING WITH PYTORCH