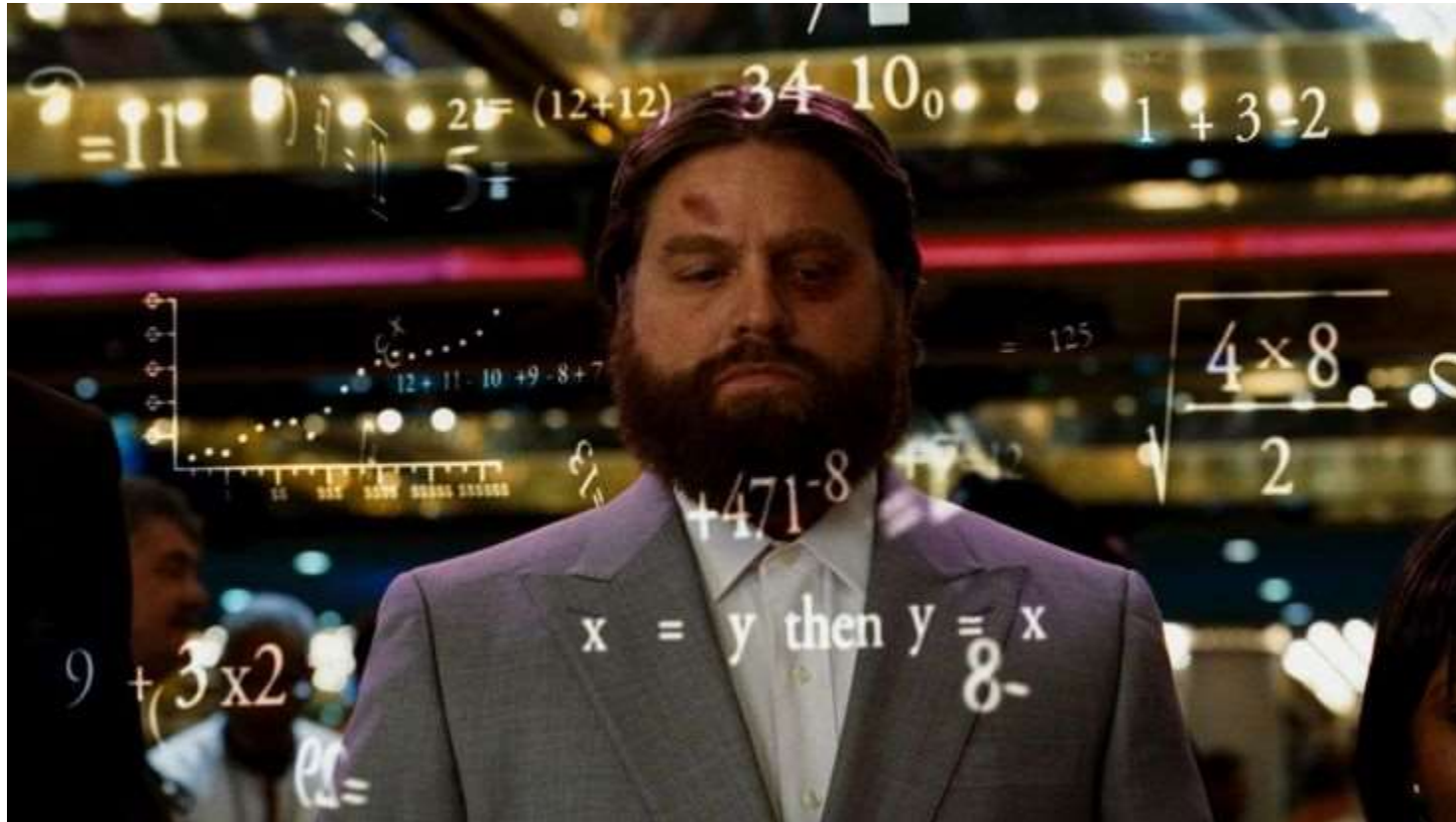


Эффективное обучение LLM

GPT-2 training memory footprint

1.5B параметров – на что уйдет память при обучении и сколько?



GPT-2 training memory footprint

На что уходит память?

- Параметры модели
- Градиенты
- Параметры оптимизатора
- Активации
- CUDA-контекст
- Данные в памяти

GPT-2 training memory footprint

Сколько?

- Параметры модели – $1.5B * 4 \text{ (fp32 = 4 байта)} = 6GB$
- Градиенты - $1.5B * 4 \text{ (fp32 = 4 байта)} = 6GB$
- Параметры оптимизатора (Adam) – $2 * 1.5B * 4 \text{ байта} = 12GB$
- Активации:
 - $12 * \text{hidden_dim} * \text{sequence_length} * \text{batch_size} * \text{num_layers} * \text{precision}$
 - $12 * 1600 * 1000 * 16 * 48 * 4 = 59GB$
- CUDA-контекст **~1GB**
- Данные в памяти – совсем мало

Итого: **~84GB**

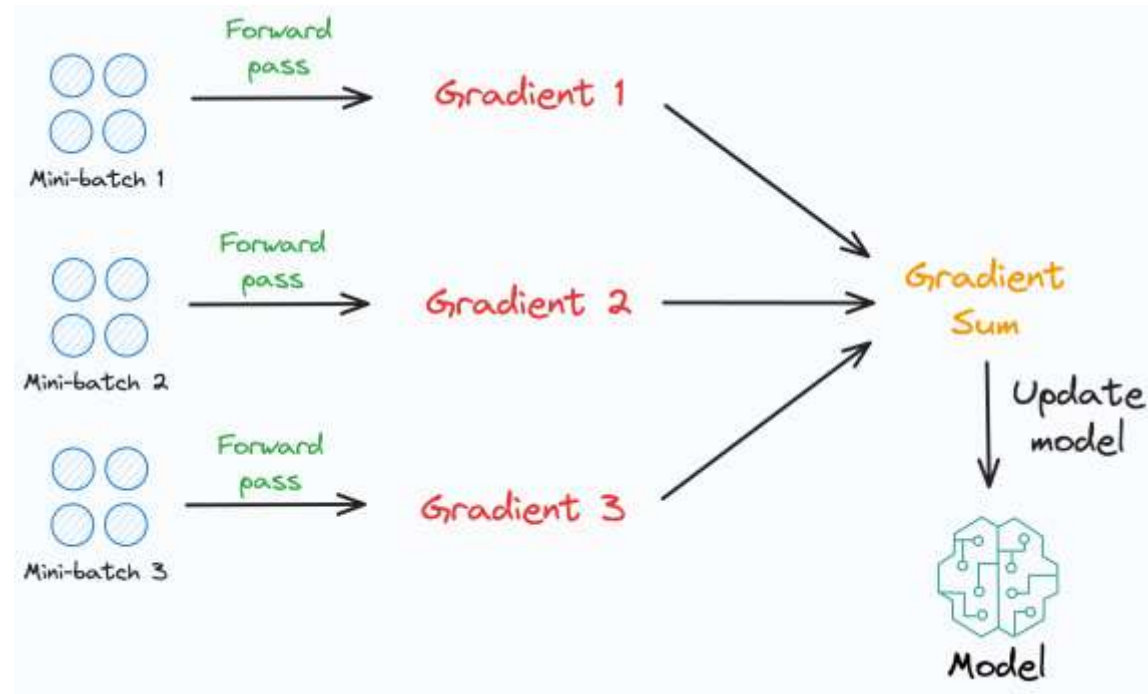
GPT-2 training memory footprint



Как оптимизировать?

- **Уменьшить batch-size** (линейно влияет на активации)
- **Уменьшить seq len** (линейно влияет на активации)
- **AMP** – уменьшает память на активации вдвое ускоряет обучение
- **Gradient checkpointing** – снижает память на активации, замедляет обучение
- **Квантизация***
- **Offload на CPU**
- **PEFT**
- **Распределенное обучение**

Gradient accumulation



```
1 for epoch in range(num_epochs):
2     train_acc = torchmetrics.Accuracy(
3         task="multiclass", num_classes=2).to(fabric.device)
4
5     for batch_idx, batch in enumerate(train_loader):
6         model.train()
7
8         ### FORWARD AND BACK PROP
9         outputs = model(
10             batch["input_ids"],
11             attention_mask=batch["attention_mask"],
12             labels=batch["label"])
13
14         outputs["loss"] = outputs["loss"] / accumulation_steps
15         fabric.backward(outputs["loss"])
16
17         ### UPDATE MODEL PARAMETERS
18         if not batch_idx % accumulation_steps:
19             optimizer.step()
20             optimizer.zero_grad()
21
22         ### LOGGING
23         model.eval()
24         with torch.no_grad():
25             predicted_labels = torch.argmax(outputs["logits"], 1)
26             train_acc.update(predicted_labels, batch["label"])
```

Automatic Mixed Precision

AMP автоматически определяет, какие операции следует выполнять в FP16, а какие — в FP32. Операции разделяются на три категории:

- Безопасные для FP16: линейные слои, свертки — выполняются в половинной точности для максимальной скорости
- Требующие FP32: операции редукции (сумма, среднее), softmax, loss функции — остаются в полной точности для стабильности
- Смешанные операции: автоматически приводятся к нужному типу данных

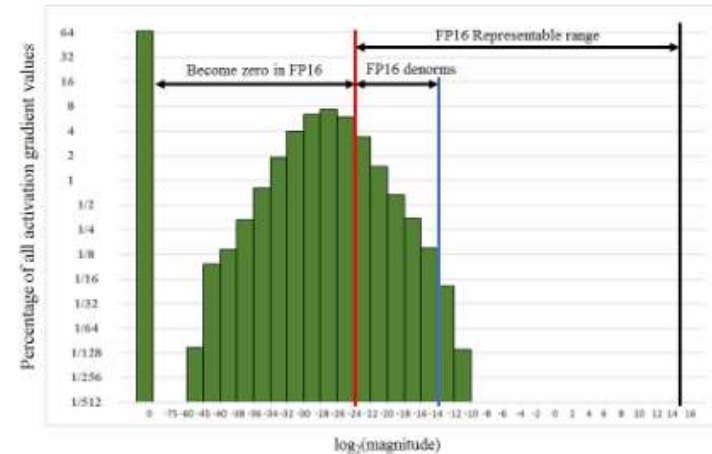
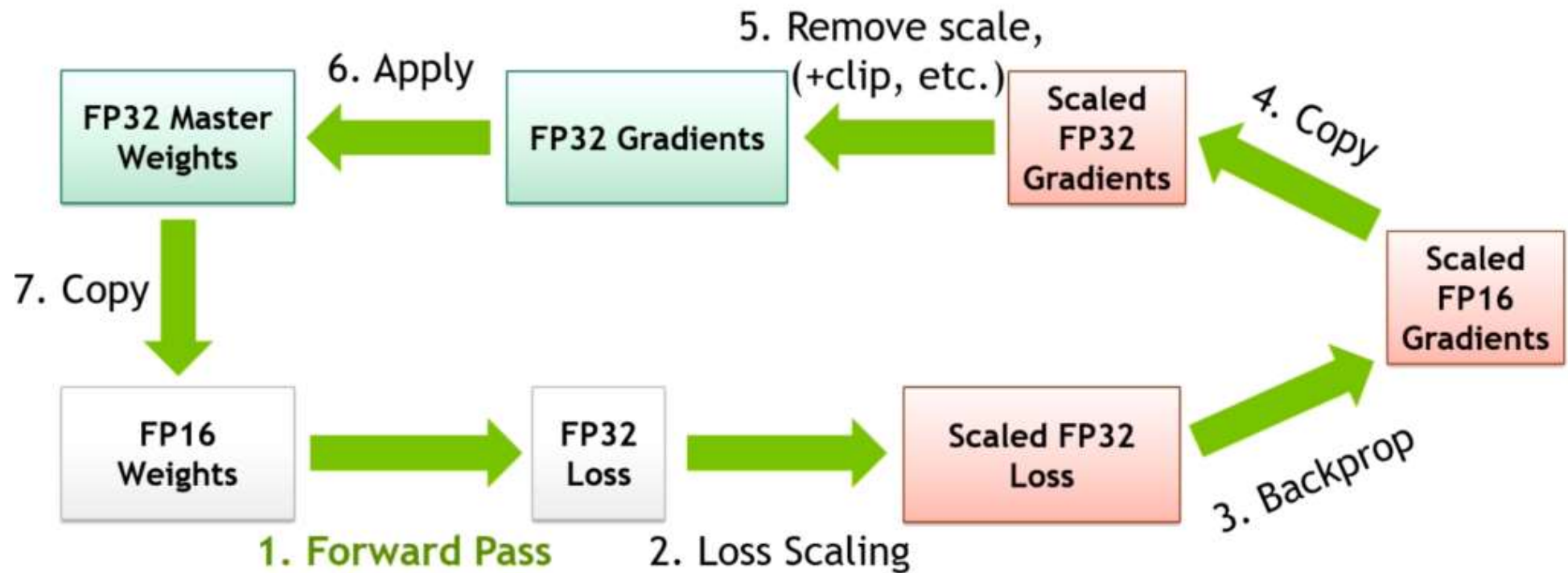


Figure 3: Histogram of activation gradient values during the training of Multibox SSD network. Note that the bins on the x-axis cover varying ranges and there's a separate bin for zeros. For example, 2% of the values are in the $[2^{-34}, 2^{-32})$ range, 2% of values are in the $[2^{-24}, 2^{-23})$ range, and 67% of values are zero.

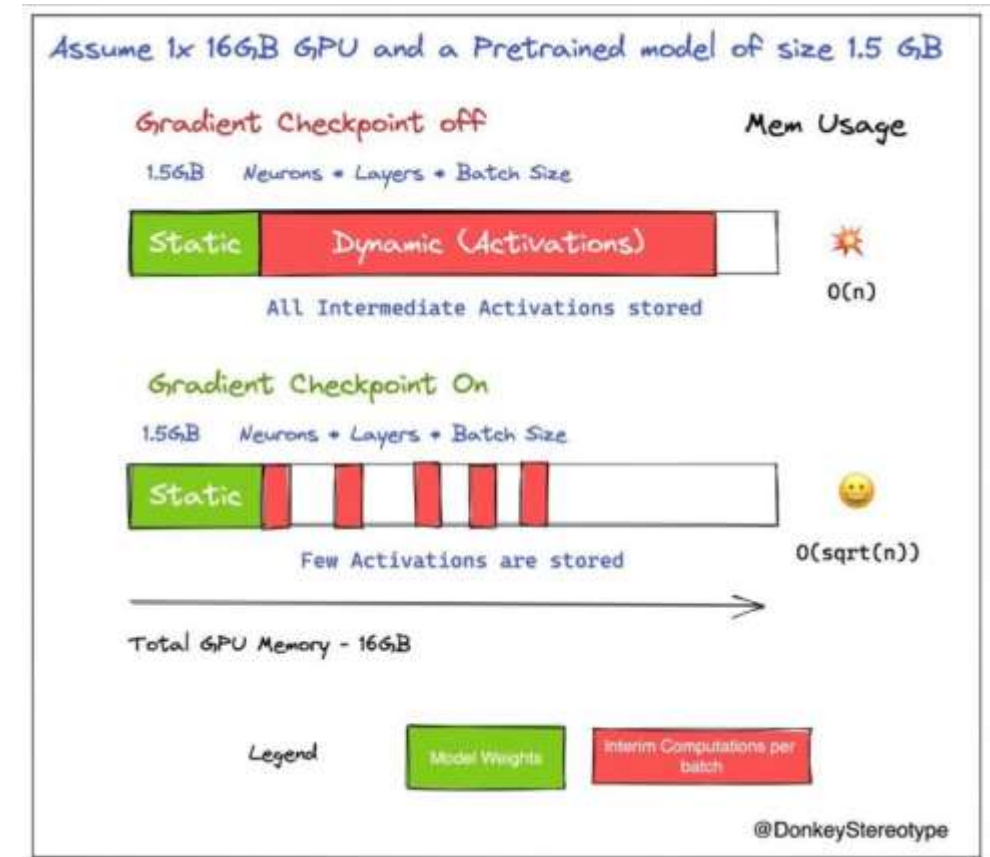
Automatic Mixed Precision

MIXED PRECISION TRAINING



Gradient checkpointing

Сеть разбивается на сегменты, и вместо сохранения всех активаций, сохраняются только активации определенных слоев (чекпойнтов)



GPT-2 training memory footprint – v2

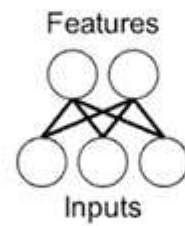
Сколько после оптимизаций?

- Параметры модели – $1.5B * 2 \text{ (fp16 = 2 байта)} = 3GB$
 - + **мастер-копия весов в fp32** = 6GB
- Градиенты - $1.5B * 2 \text{ (fp16 = 2 байта)} = 3GB$
- Параметры оптимизатора (Adam) – $2 * 1.5B * 4 \text{ байта} = 12GB$
- Активации:
 - $6 * \text{hidden_dim} * \text{sequence_length} * \text{batch_size} * \text{num_layers} * \text{precision}$
 $6 * 1600 * 1000 * 2 * 48 * 2 = 1.8GB$
- CUDA-контекст **~1GB**
- Данные в памяти – совсем мало

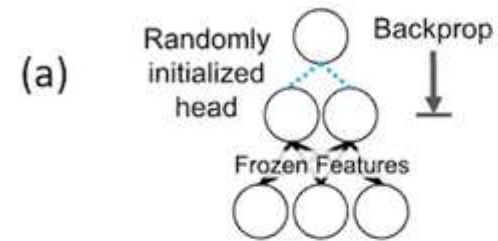
Итого: **26.8GB**

Linear probing

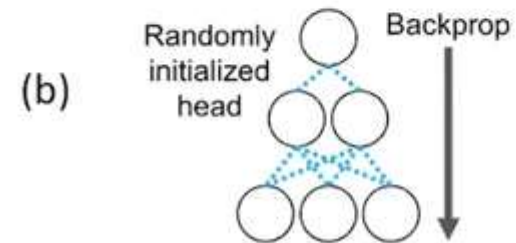
Pretraining



Linear probing

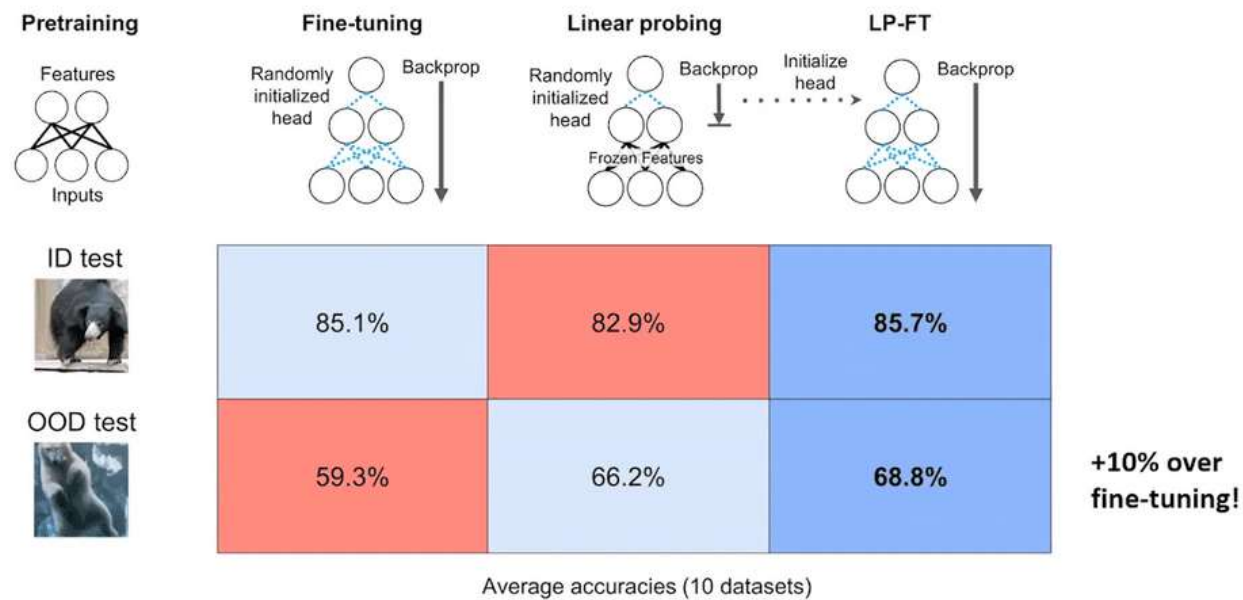


Fine-tuning



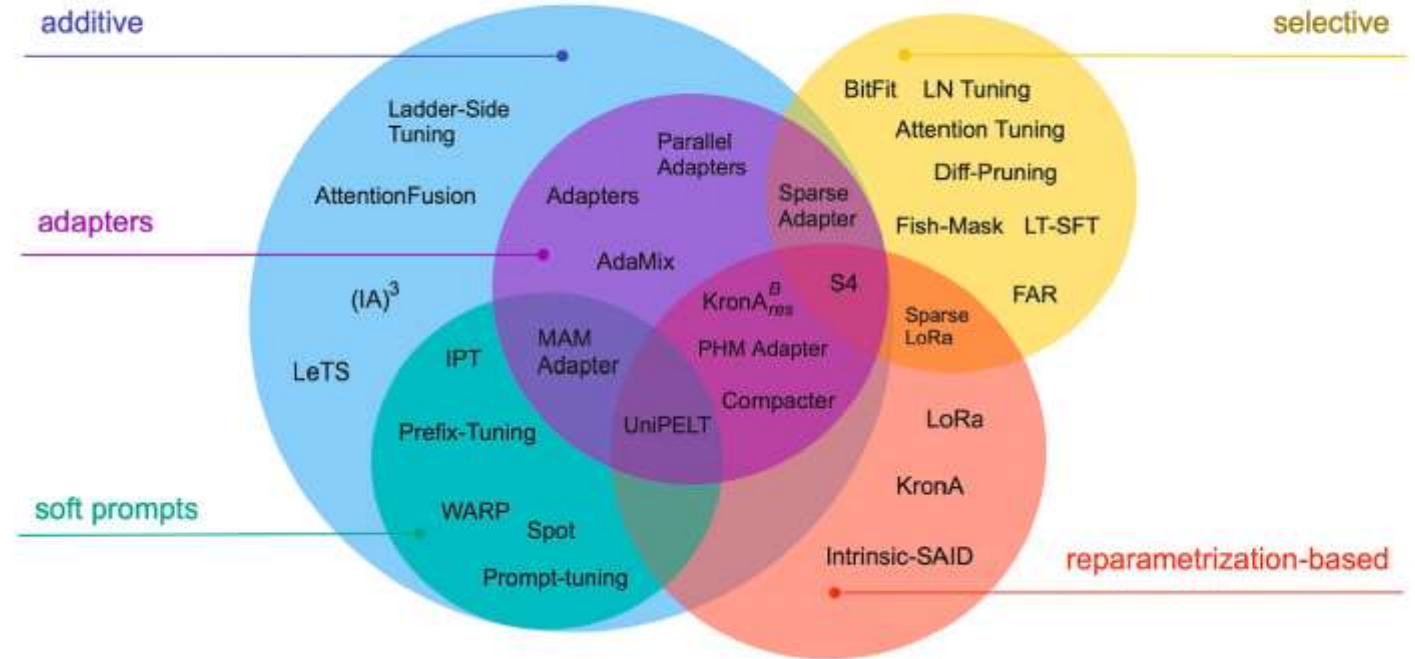
LP-FT

Improving fine-tuning



Разнообразие PEFT методов

- **Адаптерные методы** - добавляют слои между слоями исходной модели
- **Методы низкоранговой адаптации** – добавляют обучаемые матрицы низкоранговой декомпозиции в слои исходной модели
- **Методы на основе промптов** – добавляют обучаемые эмбединги в разные места*
- **Селективные методы** – обучаются только выбранные параметры

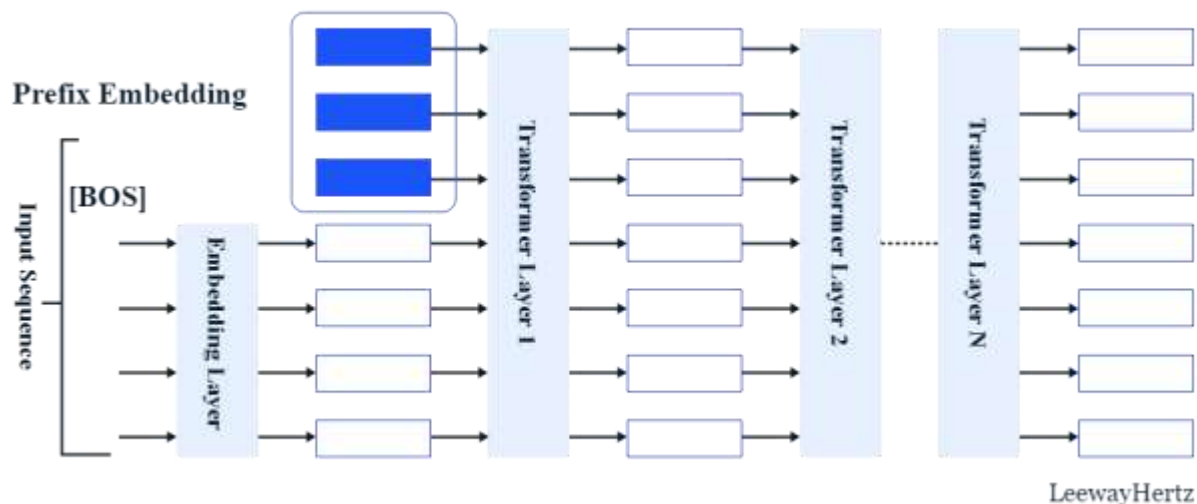


PEFT – Prompt Tuning



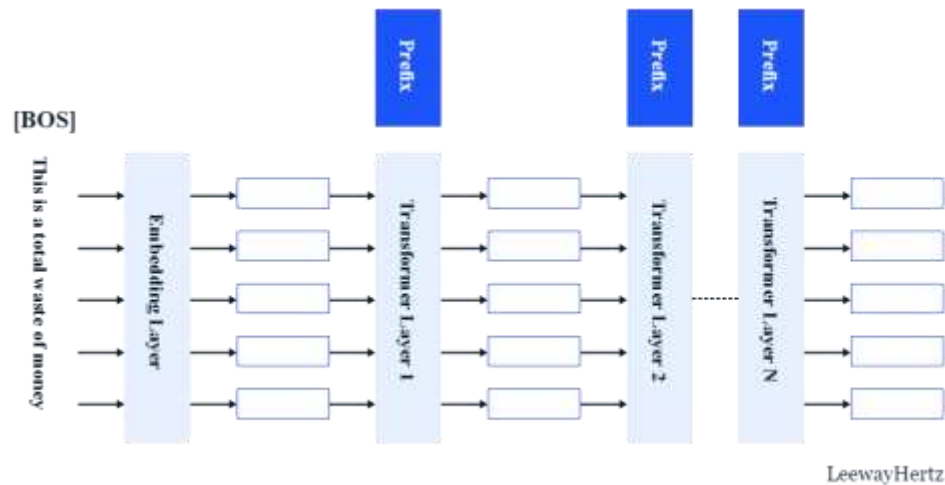
```
1 1) "Translate the English sentence '{english_sentence}' into German: {german_translation}"  
2  
3 2) "English: '{english_sentence}' | German: {german_translation}"  
4  
5 3) "From English to German: '{english_sentence}' -> {german_translation}"
```

PEFT – Prompt Tuning



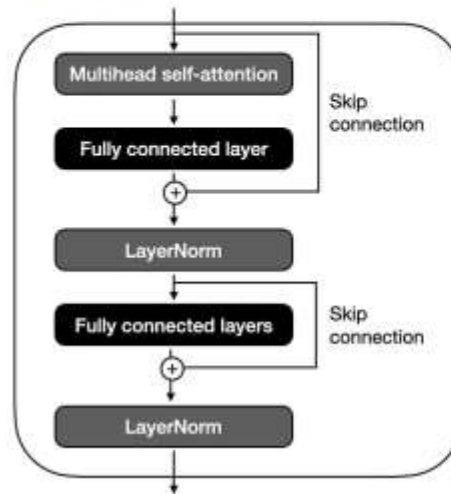
```
1 soft_prompt = torch.nn.Parameter( # Make tensor trainable
2     torch.rand(num_tokens, embed_dim)) # Initialize soft prompt tensor
3
4 def input_with_soft_prompt(x, soft_prompt):
5     x = concatenate([soft_prompt, x], # Prepend soft prompt to input
6                     dim=seq_len)
7     return x
8
9 # train soft prompt tensor via gradient descent
10 train(model(input_with_soft_prompt(x)))
11
12 # use model with soft prompts
13 model(input_with_soft_prompt(x))
```


PEFT – Prefix Tuning

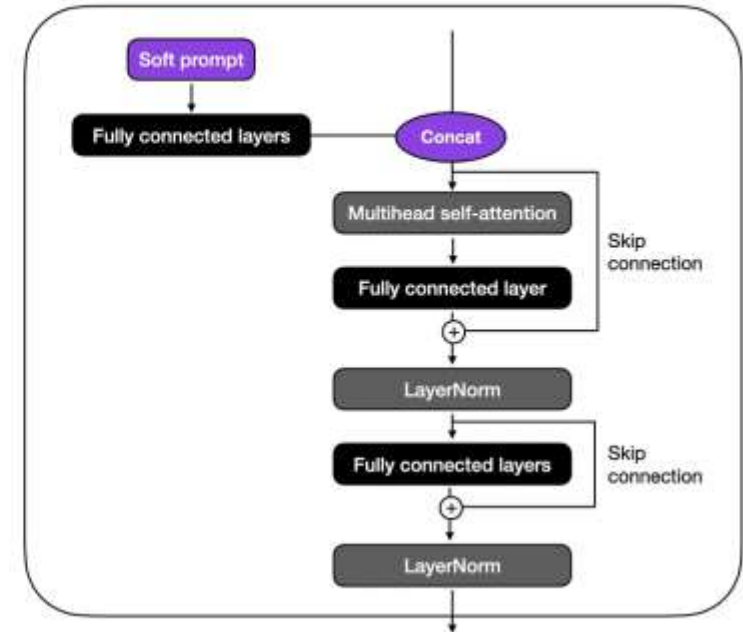


```
1 def transformer_block_with_prefix(x, soft_prompt):
2     soft_prompt = FullyConnectedLayers(soft_prompt) # Prefix
3     x = concatenate([soft_prompt, x],                # Prefix
4                     dim=seq_len)                    # Prefix
5     residual = x
6     x = self_attention(x)
7     x = LayerNorm(x + residual)
8     residual = x
9     x = FullyConnectedLayers(x)
10    x = LayerNorm(x + residual)
11    return x
```

REGULAR TRANSFORMER BLOCK

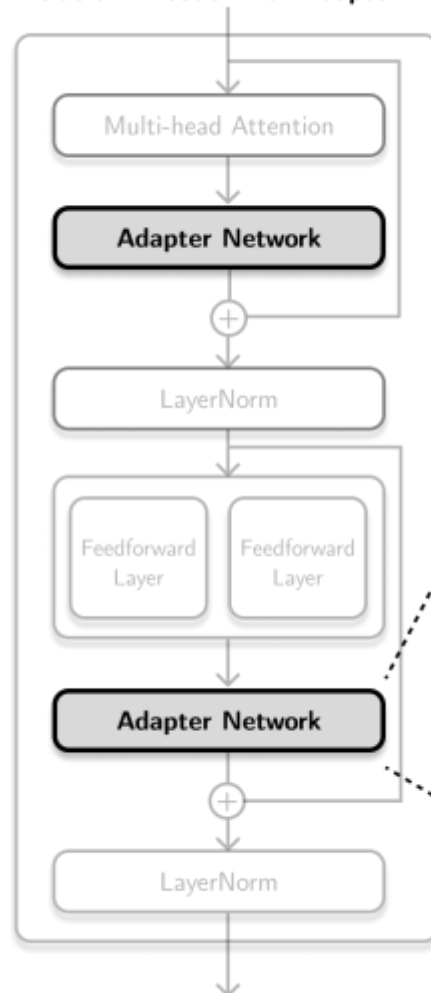


TRANSFORMER BLOCK WITH PREFIX



PEFT - Адаптеры

Inside an Encoder with Adapter



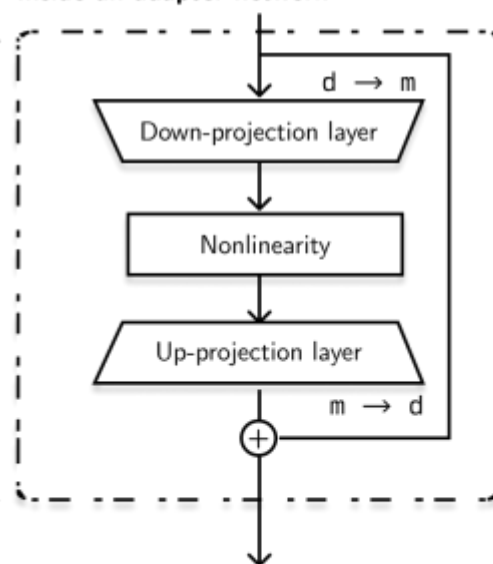
1

Add adapters inside the transformer layer. In Housby et al.'s (2019) work, the adapter networks are placed after each transformer sub-layer.

2

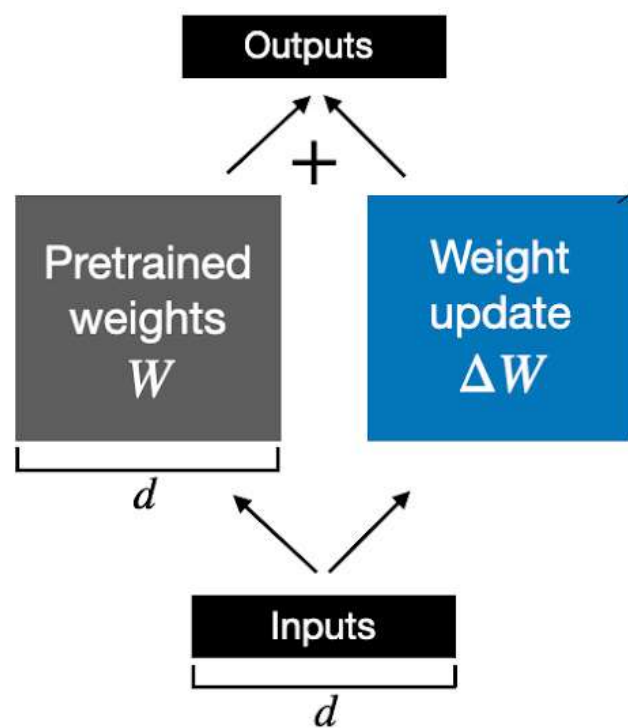
Adapters are feedforward networks. They down-project an input size of d into m , apply nonlinearity, then project it back into d .

Inside an adapter network

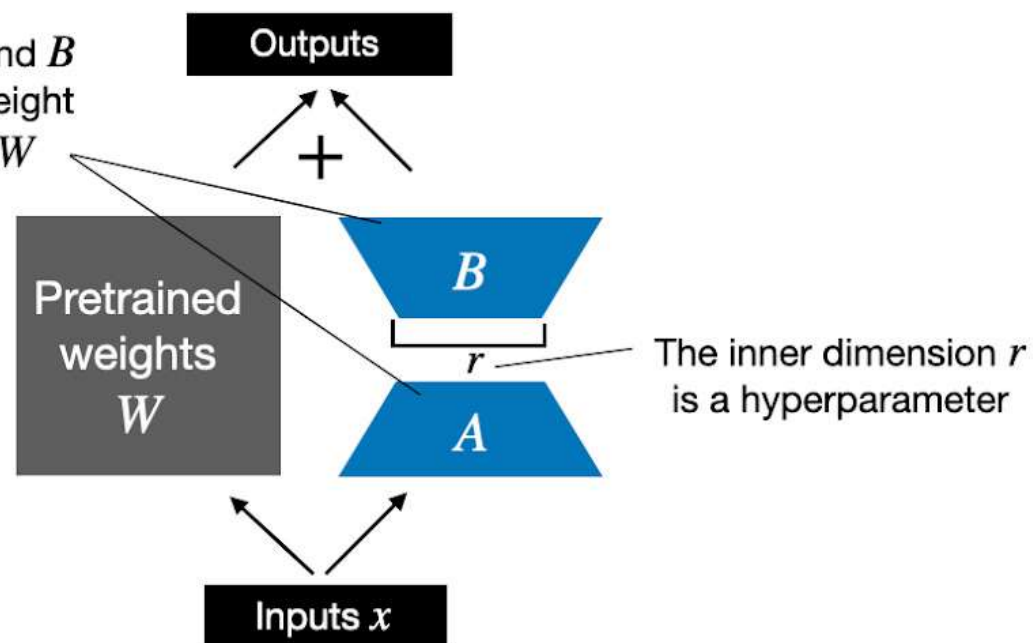


PEFT - LoRA

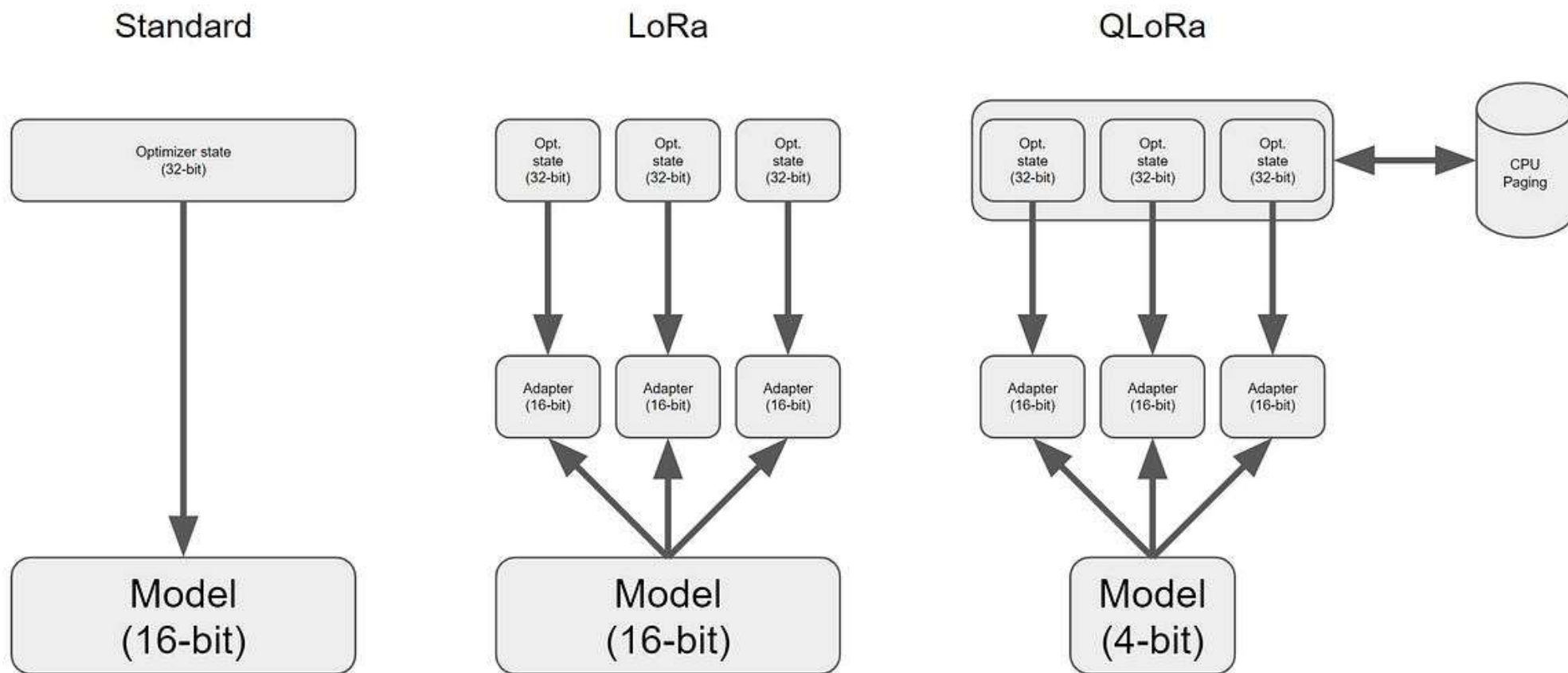
Weight update in **regular finetuning**



Weight update in **LoRA**

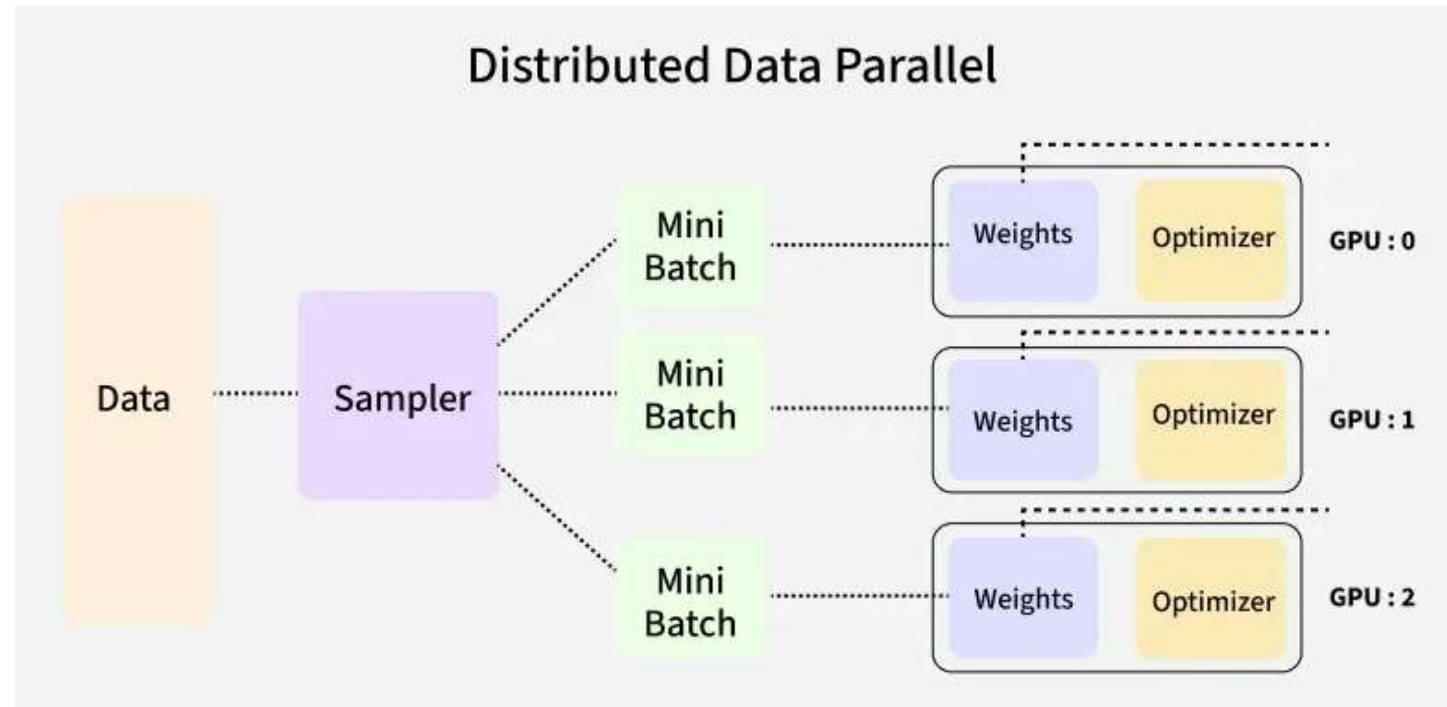


PEFT - QLoRA

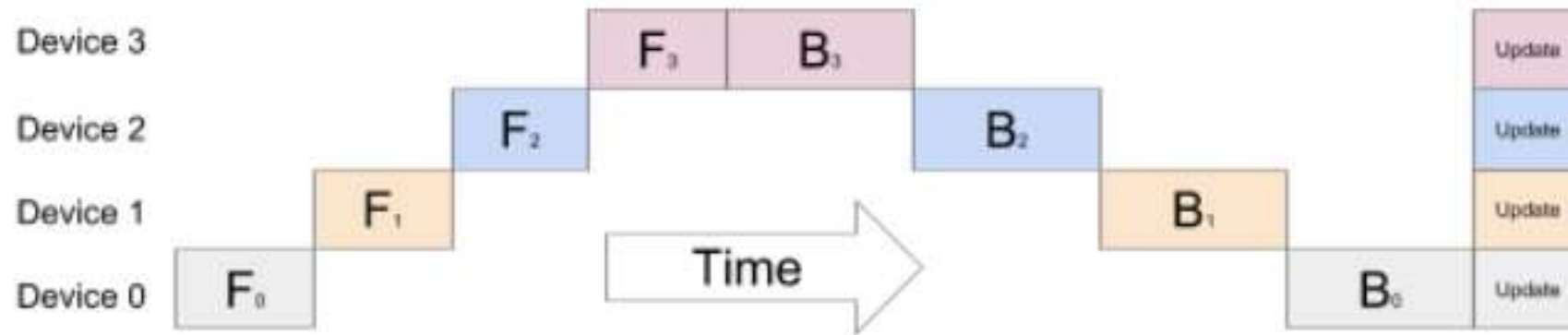


Data Parallel

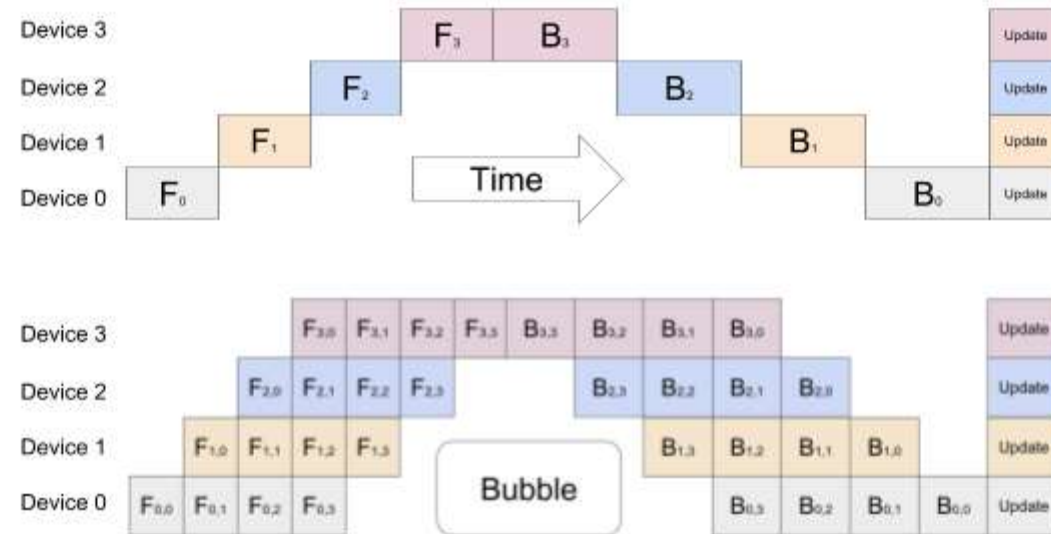
- Input разбивается и рассылается по GPU
- Параллельный FP и BP
- Градиенты собираются и усредняются на GPU0
- Параметры обновляются на GPU0 и копируются всем остальным
- Повторять до готовности



Naive model parallelism

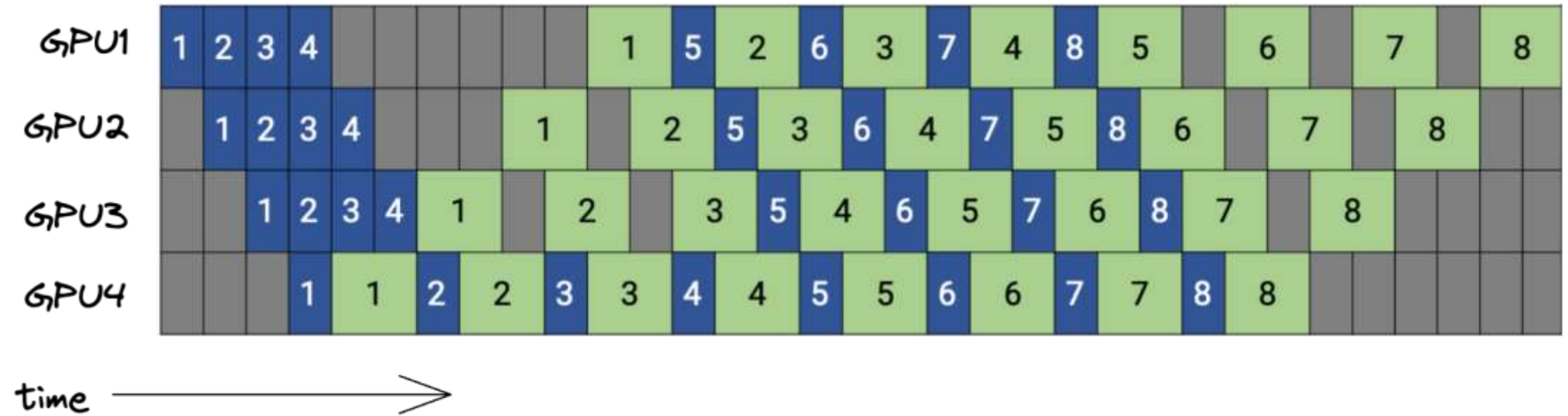


Pipeline model parallelism - GPipe



Top: The naive model parallelism strategy leads to severe underutilization due to the sequential nature of the network. Only one accelerator is active at a time. Bottom: GPipe divides the input mini-batch into smaller micro-batches, enabling different accelerators to work on separate micro-batches at the same time.

Pipeline model parallelism - PipeDream



Sharded training - ZeRO

