

## 预训练BERT

利用BERT模型和从WikiText-2数据集生成的预训练样本，在WikiText-2数据集上对BERT进行预训练。

In [1]:

```
1 import torch
2 from torch import nn
3 from d2l import torch as d2l
```

首先，加载WikiText-2数据集作为小批量的预训练样本，用于遮蔽语言模型和下一句预测。批量大小是512，BERT输入序列的最大长度是64。注意，在原始BERT模型中，最大长度是512。

In [2]:

```
1 batch_size, max_len = 512, 64
2 train_iter, vocab = d2l.load_data_wiki(batch_size, max_len)
```

## 预训练BERT

原始BERT有两个不同模型尺寸的版本。基本模型（BERT<sub>BASE</sub>）使用12层（Transformer编码器块），768个隐藏单元（隐藏大小）和12个自注意力头。大模型（BERT<sub>LARGE</sub>）使用24层，1024个隐藏单元和16个自注意力头。值得注意的是，前者有1.1亿个参数，后者有3.4亿个参数。为了便于演示，下面定义了一个小的BERT，使用了2层、128个隐藏单元和2个自注意力头。

In [3]:

```
1 net = d2l.BERTModel(len(vocab), num_hiddens=128, norm_shape=[128],
2                       ffn_num_input=128, ffn_num_hiddens=256, num_heads=2,
3                       num_layers=2, dropout=0.2, key_size=128, query_size=128,
4                       value_size=128, hid_in_features=128, mlm_in_features=128,
5                       nsp_in_features=128)
6 devices = d2l.try_all_gpus()
7 loss = nn.CrossEntropyLoss()
```

在定义训练代码实现之前，定义了一个辅助函数 `_get_batch_loss_bert`。给定训练样本，该函数计算遮蔽语言模型和下一句子预测任务的损失。请注意，BERT预训练的最终损失是遮蔽语言模型损失和下一句预测损失的和。

In [4]:

```
1 def _get_batch_loss_bert(net, loss, vocab_size, tokens_X,
2                           segments_X, valid_lens_x,
3                           pred_positions_X, mlm_weights_X,
4                           mlm_Y, nsp_y):
5     # 前向传播
6     _, mlm_Y_hat, nsp_Y_hat = net(tokens_X, segments_X,
7                                   valid_lens_x.reshape(-1),
8                                   pred_positions_X)
9     # 计算遮蔽语言模型损失
10    mlm_l = loss(mlm_Y_hat.reshape(-1, vocab_size), mlm_Y.reshape(-1)) * \
11    mlm_weights_X.reshape(-1, 1)
12    mlm_l = mlm_l.sum() / (mlm_weights_X.sum() + 1e-8)
13    # 计算下一句子预测任务的损失
14    nsp_l = loss(nsp_Y_hat, nsp_y)
15    l = mlm_l + nsp_l
16    return mlm_l, nsp_l, l
```

通过调用上述两个辅助函数，下面的 `train_bert` 函数定义了 WikiText-2 ( `train_iter` ) 数据集上预训练 BERT ( `net` ) 的过程。训练 BERT 可能需要很长时间。以下函数的输入 `num_steps` 指定了训练的迭代步数，而不是指定训练的轮数。

In [5]:

```

1 def train_bert(train_iter, net, loss, vocab_size, devices, num_steps):
2     net = nn.DataParallel(net, device_ids=devices).to(devices[0])
3     trainer = torch.optim.Adam(net.parameters(), lr=0.01)
4     step, timer = 0, d2l.Timer()
5     animator = d2l.Animator(xlabel='step', ylabel='loss',
6                             xlim=[1, num_steps], legend=['mlm', 'nsp'])
7     # 遮蔽语言模型损失的和, 下一句预测任务损失的和, 句子对的数量, 计数
8     metric = d2l.Accumulator(4)
9     num_steps_reached = False
10    while step < num_steps and not num_steps_reached:
11        for tokens_X, segments_X, valid_lens_x, pred_positions_X, \
12            mlm_weights_X, mlm_Y, nsp_y in train_iter:
13            tokens_X = tokens_X.to(devices[0])
14            segments_X = segments_X.to(devices[0])
15            valid_lens_x = valid_lens_x.to(devices[0])
16            pred_positions_X = pred_positions_X.to(devices[0])
17            mlm_weights_X = mlm_weights_X.to(devices[0])
18            mlm_Y, nsp_y = mlm_Y.to(devices[0]), nsp_y.to(devices[0])
19            trainer.zero_grad()
20            timer.start()
21            mlm_l, nsp_l, l = _get_batch_loss_bert(
22                net, loss, vocab_size, tokens_X, segments_X, valid_lens_x,
23                pred_positions_X, mlm_weights_X, mlm_Y, nsp_y)
24            l.backward()
25            trainer.step()
26            metric.add(mlm_l, nsp_l, tokens_X.shape[0], 1)
27            timer.stop()
28            animator.add(step + 1,
29                          (metric[0] / metric[3], metric[1] / metric[3]))
30            step += 1
31            if step == num_steps:
32                num_steps_reached = True
33                break
34
35    print(f'MLM loss {metric[0] / metric[3]:.3f}, '
36          f'NSP loss {metric[1] / metric[3]:.3f}')
37    print(f'{metric[2] / timer.sum():.1f} sentence pairs/sec on '
38          f'{str(devices)}')

```

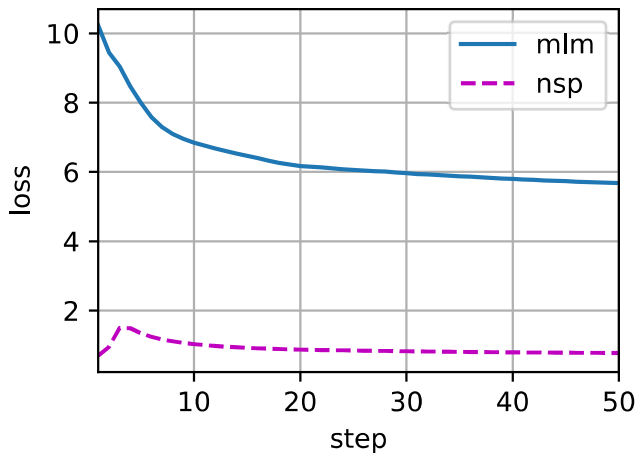
在预训练过程中, 可以绘制出遮蔽语言模型损失和下一句预测损失。

In [6]:

```
1 train_bert(train_iter, net, loss, len(vocab), devices, 50)
```

MLM loss 5.679, NSP loss 0.773

1003.4 sentence pairs/sec on [device(type='cuda', index=0), device(type='cuda', index=1), device(type='cuda', index=2), device(type='cuda', index=3)]



## 用BERT表示文本

在预训练BERT之后，可以用它来表示单个文本、文本对或其中的任何词元。下面的函数返回 `tokens_a` 和 `tokens_b` 中所有词元的BERT (`net`) 表示。

In [7]:

```
1 def get_bert_encoding(net, tokens_a, tokens_b=None):
2     tokens, segments = d2l.get_tokens_and_segments(tokens_a, tokens_b)
3     token_ids = torch.tensor(vocab[tokens], device=devices[0]).unsqueeze(0)
4     segments = torch.tensor(segments, device=devices[0]).unsqueeze(0)
5     valid_len = torch.tensor(len(tokens), device=devices[0]).unsqueeze(0)
6     encoded_X, _, _ = net(token_ids, segments, valid_len)
7     return encoded_X
```

考虑“a crane is flying”这句话。插入特殊标记“<cls>”（用于分类）和“<sep>”（用于分隔）后，BERT输入序列的长度为6。因为零是“<cls>”词元，`encoded_text[:, 0, :]` 是整个输入语句的BERT表示。为了评估一词多义词元“crane”，还打印出了该词元的BERT表示的前三个元素。

In [8]:

```
1 tokens_a = ['a', 'crane', 'is', 'flying']
2 encoded_text = get_bert_encoding(net, tokens_a)
3 # 词元: '<cls>', 'a', 'crane', 'is', 'flying', '<sep>'
4 encoded_text_cls = encoded_text[:, 0, :]
5 encoded_text_crane = encoded_text[:, 2, :]
6 encoded_text.shape, encoded_text_cls.shape, encoded_text_crane[0][:3]
```

Out[8]:

```
(torch.Size([1, 6, 128]),
 torch.Size([1, 128]),
 tensor([-1.4044,  1.8831, -0.2498], device='cuda:0', grad_fn=<SliceBackward0>))
```

现在考虑一个句子“a crane driver came”和“he just left”。类似地，`encoded pair[:, 0, :]` 是来自预训练BERT

的整个句子对的编码结果。注意，多义词元“crane”的前三个元素与上下文不同时的元素不同。这支持了BERT表示是上下文敏感的。

In [9]:

```
1 tokens_a, tokens_b = ['a', 'crane', 'driver', 'came'], ['he', 'just', 'left']
2 encoded_pair = get_bert_encoding(net, tokens_a, tokens_b)
3 # 词元: '<cls>', 'a', 'crane', 'driver', 'came', '<sep>', 'he', 'just',
4 # 'left', '<sep>'
5 encoded_pair_cls = encoded_pair[:, 0, :]
6 encoded_pair_crane = encoded_pair[:, 2, :]
7 encoded_pair.shape, encoded_pair_cls.shape, encoded_pair_crane[0][:3]
```

Out[9]:

```
(torch.Size([1, 10, 128]),
 torch.Size([1, 128]),
 tensor([-0.0134,  1.2939, -1.4307], device='cuda:0', grad_fn=<SliceBackward0>))
```