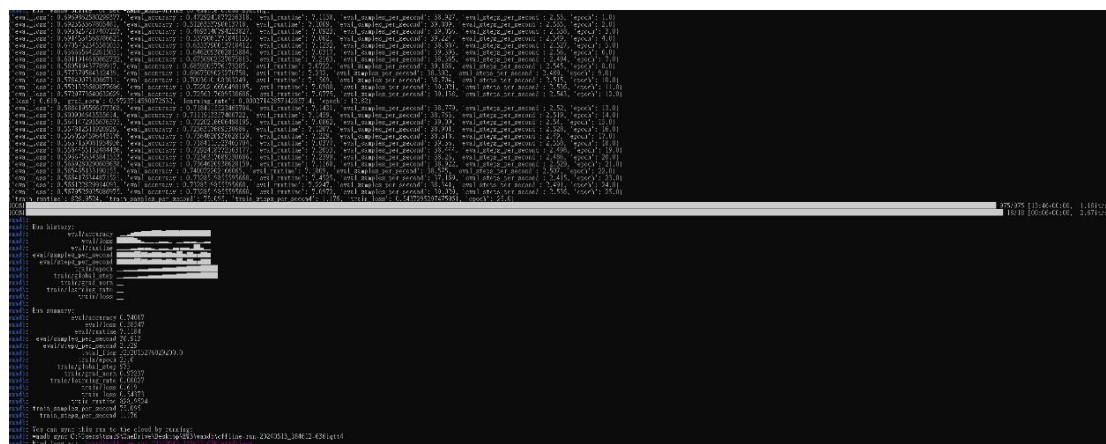## Model Analysis

In this assignment, I selected RoBERTa, a natural language processing model based on BERT. It enhances BERT by removing the Next Sentence Prediction task, as research has shown that it contributes little to performance improvements. Instead, RoBERTa only uses the Masked Language Model task for training. Additionally, RoBERTa is trained for a longer time and with more steps, enabling the model to capture language features more effectively.

In the code, I used the `transformers.AutoModelForSequenceClassification.from_pretrained()` function, setting `num_labels` to 2 for classification, and then loading the model. Afterward, I used the `trainer` function to train the model with the dataset. This function optimizes the model parameters through backpropagation and gradient descent and computes performance metrics on the validation set at the end of each training epoch, such as accuracy and loss. The trained model is saved, and I used the default AdamW optimizer.

- Rte dataset_bitfit: Acc: 0.74007



- Rte dataset_lora: Acc: 0.72563

- Sst2 dataset_bitfit: Acc: 0.8922



- Sst2 dataset_lora: Acc: 0.92317



## PEFT Discussion

For transformer models, Bitfit freezes most model parameters and only updates the bias parameters and the classification layer specific to the task. Examples of bias parameters include:

- Bias in the MLP layer

- Bias involved in calculating query, key, and value in the attention module and merging multiple attention results
- Bias in the Layer Normalization layer

These bias parameters account for around 0.08% of the total parameters in a model like BERT.

By updating only a portion of the bias parameters and task-specific classification layer parameters, the model can better recognize specific downstream tasks, thereby improving accuracy without training a large number of parameters.

The best hyperparameters I found are:

**Rte dataset:**
```
args = transformers.TrainingArguments(
  output_dir='/save',
  num_train_epochs=25,
  learning_rate=5e-4,
  per_device_train_batch_size=16,
  per_device_eval_batch_size=16,
  gradient_accumulation_steps=4,
  warmup_steps=100,
  weight_decay=0.01,
  evaluation_strategy='epoch',
  save_strategy='epoch',
  save_steps=100,
  eval_steps=100,
  load_best_model_at_end=True,
  metric_for_best_model='accuracy'
)

lora_config = LoraConfig(
  task_type=TaskType.SEQ_CLS,
  r=2,
  lora_alpha=8,
  lora_dropout=0.01
)
```
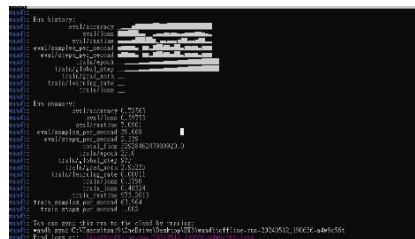Since the Rte dataset is small, I used more epochs to increase its training.

**Sst2 dataset:**
```

```
args = transformers.TrainingArguments(
 output_dir='/save',
 num_train_epochs=5,
 learning_rate=2e-5,
 per_device_train_batch_size=16,
 per_device_eval_batch_size=16,
 gradient_accumulation_steps=4,
 warmup_steps=100,
 weight_decay=0.01,
 evaluation_strategy='epoch',
 save_strategy='epoch',
 save_steps=100,
 eval_steps=100,
 load_best_model_at_end=True,
 metric_for_best_model='accuracy'
)

lora_config = LoraConfig(
 task_type=TaskType.SEQ_CLS,
 r=2,
 lora_alpha=8,
 lora_dropout=0.01
)
```

In full-finetuning, we update all the model parameters, which usually requires a smaller learning rate to avoid disrupting previously learned knowledge. However, in PEFT, we only update part of the model's parameters, often those in the last few layers or certain other layers. These layers may have different feature representations and learning speeds, so the learning rate may need to be adjusted to better suit their updates. If certain layers require larger parameter changes, a higher learning rate might be needed, while for other layers, a smaller learning rate could prevent overfitting. Therefore, the learning rate needs to be adjusted according to the specific layers chosen for partial fine-tuning and the requirements of the task to achieve optimal performance.

## PEFT Comparison

Rte dataset_bitfit:                              Rte dataset_lora:

Sst2 dataset_bitfit:                    sst2 dataset_lora:



We can see that regardless of the dataset, the accuracy measured is similar. I believe this could be because the SST-2 and Rte datasets are relatively simple, and the capabilities of these different models may be sufficient for the task, so there is no noticeable performance difference.

The size of the Lora rank affects model training in the following ways:

A lower rank may result in underfitting, which, while improving generalization, could cause the model to miss important features in the data, reducing its performance. The parameter count will decrease, whereas a higher rank could lead to overfitting and increased computational complexity due to more parameters and larger matrices. This may increase training time, especially for large datasets and complex models, where the model may overfit the training data, leading to poor performance on unseen data, and the parameter count increases.