# From Text to Insights: Transformers and Contrastive Learning On Amazon Massive

**Sarah Yang**
sxyang@ucsd.edu

## Abstract

In this project, we explored various transformer-based text classification methods, particularly the BERT model. We used the Amazon Massive Scenario dataset that has 18 distinct practical scenario categories to classify. WE first construct and evaluate base-line BERT model. To improve on the model , we employed fine-tuning techniques such as Layer-wise Learning Rate Decay (LLRD) and learning rate warmup, along with advanced contrastive learning methods including Supervised Contrastive Learning (SupContrast) and SimCLR, we aimed to enhance the transformer's ability to classify scenarios effectively. Experiments indicated that SupContrast outperformed SimCLR due to its utilization of label information, achieving a final test accuracy of 91.06% compared to SimCLR's 84.10%. We also conducted experiments with Low-Rank Adaptation (LoRA), demonstrating efficient parameter usage with minimal performance compromisess. This project demonstrates the important role of model architecture selection, training strategy refinement, and precise hyperparameter tuning in achieving optimal text classification outcomes.

## 1 Introduction

In today's society, text classification becomes one of the major challenges in the field of natural language processing (NLP). With the explosion of social networks, textual data grows tremendously across various platforms. Therefore, it becomes pressing to develop effective models that can accurately interpret and categorize all the text information. In recent years, various applications such as predictive texts, sentiment analysis, and conversational chatbots have heavily depend on robust and precise text classification methods to resolve challenges and drive business success.

Many transformer-based architectures have developed. One of the well-know models, the Bidirectional Encoder Representations from Transformers (BERT), marked a significant advancement in NLP. Transformers, unlike earlier neural network architectures, outperform due to their unique self-attention mechanism, allowing the models to capture contextual relationships across extensive textual sequences effectively. For further improvement, more techniques are developed to fine-tune BERT and use in different classification scenarios, such as LORA, RoBERTa, etc. These capabilities have result in remarkable achievements in a wide areas of NLP tasks.

In this project, we use the BERT model to perform scenario classification for the Amazon Massive Scenario dataset, an extensive collection of real-world text samples containing 18 distinct practical scenario categories. Our goal of the project is to test how effectively each fine-tuning technique applied to Transformers enhance their ability to correctly classify scenarios. We also perform experiments with baseline fine-tuning, advanced customization strategies, and contrastive learning techniques, conducting experiments to evaluate their performances. By experimenting various techniques, we aim not only to deepen our understanding of Transformer-based models' capabilities and limitations, but also to provide valuable insights into optimizing their performances for other practical text classification tasks.

## 2 Related Work

We got insights from the implementation of the BERT model from [1]. This document from huggingface provides insights and explanations on various BERT techniques, such as attention mechanisms and token embeddings, which directly informed our practical implementations.

We also got insights from implementation of the LORA by PEFT experiments [2] and [3]. We learned how LoRA efficiently reduces the number of trainable parameters by introducing low-rank matrices into model layers. As a result, it significantly decreasing computational costs and memory usage while maintaining high performance.

Additionally, we learned various custom fine-tuning strategies from [4]. Specifically, we got insights of the Layer-wise Learning Rate Decay (LLRD) and Learning Rate warmup techniques to help us better generalize and improve convergence. Specifically, LLRD assigns different learning rates across model layers to stabilize training, and Learning Rate Warmup gradually increases learning rates early in training to avoid initial instability and achieve smoother convergence.

Finally, we further learned contrastive learning of SupContrast[5] and SimCSE[6]. SupCon utilizes label information to bring embeddings from the same class closer and differentiate those from different classes. As a result, it enhances classification accuracy and robustness. SimCSE uses dropout-based data augmentation to generate varied contextually consistent sentence embeddings, effectively improving performance on semantic textual similarity tasks and embedding quality overall.

## 3 Methods

### 3.1 Baseline

The experiment utilizes the BERT-base-uncased model from Hugging Face's transformers for the Amazon Massive Scenario classification task. The model processes text with a maximum sequence length of 128 tokens, using left-side truncation in the tokenizer (truncation_side="left"). The architecture includes a BERT encoder followed by a two-layer classifier with dimensions $768 \rightarrow 256 \rightarrow 60$, where 768 is BERT's hidden size, 256 is the intermediate hidden dimension, and 60 is the number of target scenario classes. Training employs the AdamW optimizer with a learning rate of $2e - 5$ and epsilon of $1e - 8$, coupled with a linear learning rate scheduler with no warmup steps. The dropout rate is set to $0.1$ for regularization, and the model is trained with a batch size of 16 for 3 epochs. The dataset is split into $11,514$ training samples, $2,033$ validation samples, and $2,974$ test samples, with features cached in 'assets/cache/amazon.pkl' for efficiency. Cross-entropy loss is used as the training objective, and evaluation metrics are computed using accuracy. All random seeds are set to 42 for reproducibility, and the training progress is monitored using tqdm progress bars.

### 3.2 Custom

To enhance fine-tuning performance, we implemented two optimization strategies from [4]: Layer-wise Learning Rate Decay (LLRD) and Learning Rate Warmup. LLRD is a fine-tuning technique that assigns different learning rates to different layers of a pre-trained transformer model. In our implementation, lower layers (closer to the input) received smaller learning rate updates, while higher layers (closer to the classification head) were allowed larger updates. This approach maintains the stability of pre-trained representations in early layers while allowing task-specific adaptation in later layers. Warmup mitigates instability in early training by gradually increasing the learning rate instead of starting with a high initial value. This prevents drastic weight updates at the beginning of training, which could hinder convergence. These techniques were integrated into our fine-tuning procedure, leading to better generalization and improved convergence. Below are the hyperparameter settings for the custom model:

- **Base Model:** BERT
- **Optimizer:** AdamW with $\epsilon = 1e{-}8$
- **Initial Learning Rate:** $0.5e{-}4$
- **Decay Rate for LLRD:** $0.8$
- **Warmup Steps:** $10\%$ of total training steps

- **Batch Size:** 16
- **Epochs:** 10
- **Loss Function:** CrossEntropyLoss

## 3.3 Contrastive Learning

Contrastive Learning. Building on the same BERT-base-uncased encoder and a two-layer classifier as in the baseline, we insert an additional projection head (a single linear layer, 768→768, followed by L2 normalization) to produce embeddings for a two-phase training pipeline. In the contrastive phase, the model uses BERT to compute [CLS] outputs, applies dropout, and feeds them into the projection head. We then compute a contrastive loss on these embeddings in two variants:

Supervised Contrastive (SupContrast). Leveraging label information, we form positive pairs among samples of the same class and negative pairs otherwise. We train for 10 epochs with a temperature of 0.07, dropout rate of 0.3, and learning rate = 1e4, using batches of size 16–32. Others are following the default baseline setting.

SimCLR (Unsupervised). Instead of labels, we rely purely on dropout-based text augmentation to generate different "views." We set temperature = 0.1, batch size = 32, and dropout = 0.2, again for 10 epochs at 1e4 learning rate. No label signals are used in this stage (labels=None). Others are following the default baseline setting.

After these 10-epoch (epoch numbers will be changed for better result, but the actual benefit by adding epoch seems not very direct in our observation) contrastive phases (either SupContrast or SimCLR), we turn off "contrast mode" so that the model outputs classification logits rather than embeddings. We then perform a cross-entropy fine-tuning stage, where the same two-layer classifier head from the baseline is optimized for another 10 epochs. Here, we reduce the learning rate to 5e5, preserving the rest of the hyperparameters from the baseline to align the learned embeddings with the Amazon Massive scenario classification objective. Two experinemnts with that hyper-parameter settings should give a relatively optimal result. Compare to cross-entropy (baseline model), the result should fit our expectations.

## 4 Results

### 4.1 Baseline Model

The baseline BERT model, trained for 10 epochs on the Amazon Massive Scenario classification task, demonstrates strong performance progression. Starting with $0\%$ accuracy during initial evaluation, the model rapidly improves to $66.01\%$ validation accuracy and $36.90\%$ training accuracy after the first epoch (loss: 1777.42). Performance continues to improve significantly through subsequent epochs, reaching $87.36\%$ validation and $76.06\%$ training accuracy by epoch 2 (loss: 833.89), and $88.88\%$ validation with $89.18\%$ training accuracy by epoch 3 (loss: 480.74). The learning curve shows steady but diminishing improvements in later epochs, with validation/training accuracies reaching $89.08\%/93.06\%$ by epoch 4, $89.38\%/94.75\%$ by epoch 5, and validation accuracy peaking at $91.15\%$ by epoch 8 while training accuracy continues to improve to $99.00\%$ by epoch 10. The final test accuracy of $90.45\%$ closely aligns with the validation performance, suggesting good generalization despite the increasing gap between training and validation accuracy in later epochs. The model's loss steadily decreases from 1777.42 to 68.01 over the 10 epochs, indicating stable optimization dynamics and effective learning of the classification task.

### 4.2 Custom Finetuning Strategies

The baseline model achieved the highest training, validation, and test accuracy, demonstrating its strong generalization ability without additional fine-tuning techniques. Among the fine-tuned models, applying the warmup technique alone resulted in a test accuracy of 0.841, outperforming the LLRD-only model, which performed slightly worse. The lowest test accuracy was observed when both LLRD and warmup were used together. However, despite the decrease in accuracy, the models incorporating fine-tuning techniques exhibited significantly lower training and test loss, as seen in the table. This is due to the fact that techniques such as LLRD and warmup primarily focus on

optimizing the learning dynamics rather than simply maximizing accuracy. Warmup helps stabilize early training by preventing drastic weight updates, while LLRD enables more refined adjustments in different layers, reducing overfitting and improving model calibration. These effects lead to smoother convergence and a more optimized decision boundary, which can reduce overall loss even if accuracy does not see a direct increase.

### 4.2.1 Layer-wise Learning Rate Decay (from result)

According to the result we got for Layer-wise Learning Rate Decay, we could confirm that the LLRD strategy modifies learning rates across different layers of the Transformer model:

- Lower layers (closer to input) are updated with a smaller learning rate, preserving pre-trained knowledge.
- Higher layers (closer to output) receive larger learning rates to adapt faster to the new classification task.
- For our custom model, the learning rate for each layer follows a **decay factor of 0.8** per layer.

This approach stabilizes fine-tuning while ensuring effective feature extraction from pre-trained representations.

### 4.2.2 Learning Rate Warmup (from result)

According to the result we got for Learning Rate Warmup, we could confirm that the Warmup mitigates instability in early training by gradually increasing the learning rate:

- The model begins training with a low learning rate.
- Over the first $10\%$ of total steps, the learning rate is linearly increased to the initial learning rate.

This approach allows smoother weight updates, preventing sudden jumps that may lead to suboptimal convergence.

## 4.3 SupContrast

The SupContrast BERT model, trained for 10 epochs on the Amazon Massive Scenario classification task, demonstrates a clear performance progression in both contrastive pretraining and classification fine-tuning stages. During contrastive learning, the model initially exhibits 0% validation accuracy, as expected (During the contrastive learning phase, the model does not perform classification but instead learns to structure embeddings, so accuracy is not computed, making 0% accuracy expected.), with the training loss reducing steadily from 629.16 at epoch 0 to 225.04 by epoch 9, indicating stable optimization dynamics. However, minimal validation accuracy improvements are observed in this phase, with only a negligible 0.0044% accuracy recorded at epoch 0 and minor fluctuations across epochs. Upon transitioning to classification fine-tuning, the model rapidly improves in performance, achieving 75.26% validation accuracy after the first epoch (loss: 1049.63). The learning curve continues to show steady and consistent improvements across subsequent epochs, reaching 83.47%/85.34% validation accuracy by epochs 1 and 2, respectively, and peaking at 90.99% by epoch 8. Training accuracy similarly improves to 99.00% by the final epoch. The final test accuracy of 91.06% aligns closely with validation performance, suggesting effective generalization. The loss consistently decreases from 629.16 (contrastive phase) to 77.72 in the final fine-tuning epoch, indicating stable optimization and successful feature representation learning for classification.

## 4.4 SimCLR

The SimCLR BERT model, trained for 10 epochs on the Amazon Massive Scenario classification task, demonstrates a structured learning progression across both the contrastive pretraining and classification fine-tuning stages. During contrastive learning, the model exhibits 0% validation accuracy, as expected. We can also see later the loss decrease as expected. In the classification fine-tuning stage, the model demonstrates rapid improvement, achieving 69.36% validation accuracy

after the first epoch (loss: 801.04, train accuracy: 53.34%). Performance continues to improve steadily, reaching 72.01%/76.05% validation accuracy by epochs 1 and 2, respectively, and peaking at 85.59% by epoch 8. Training accuracy similarly improves, reaching 92.93% by the final epoch. The final test accuracy of 84.10% is slightly lower than validation performance, suggesting minor overfitting. The loss consistently decreases from 801.04 (initial fine-tuning) to 101.81 in the final epoch, indicating stable optimization and effective feature representation learning for classification.

## 4.5 Table

| Experiment Index | Experiment | Loss | Accuracy |
|---|---|---|---|
| 1 | Test set before fine-tuning | 2765.93 | 13.0% |
| 2 | Test set after fine-tuning | 68.01 | 90.4% |
| 3 | Test set with $1^{st}$ technique (LLRD) | 1.9048 | 65.0% |
| 4 | Test set with $2^{nd}$ technique (Warmup) | 0.7365 | 84.1% |
| 5 | Test set with 2 techniques | 0.1.9080 | 64.9% |
| 6 | Test set with SupContrast | 77.72 | 91.1% |
| 7 | Test set with SimCLR | 101.81 | 85.4% |

Table 1: Experiment Table

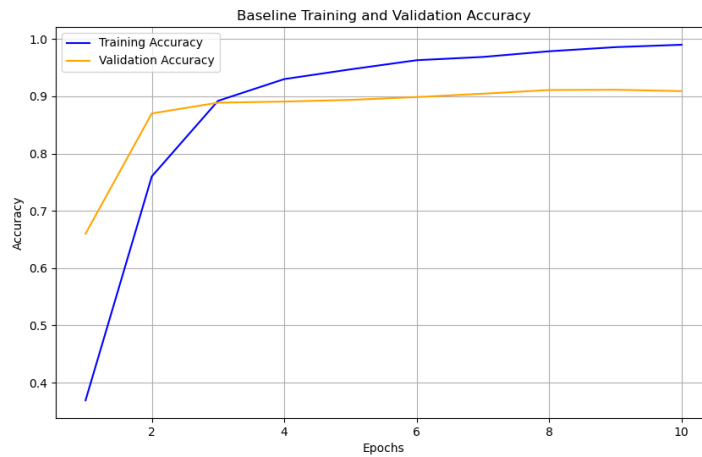## 4.6 Training and Validation Accuracy Plots



Figure 1: Baseline Training and Validation Accuracy
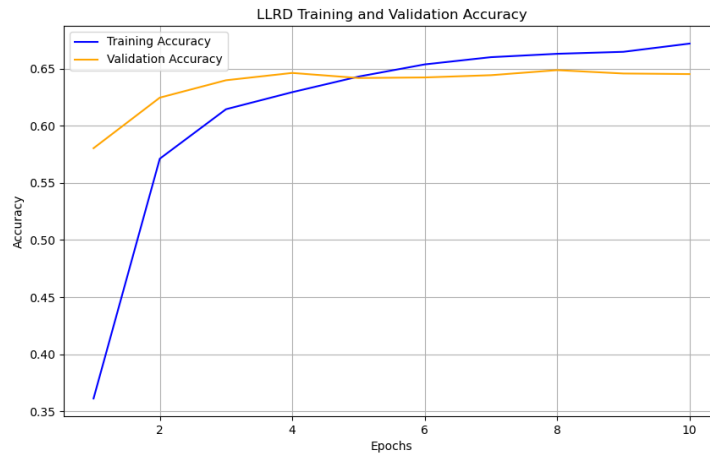
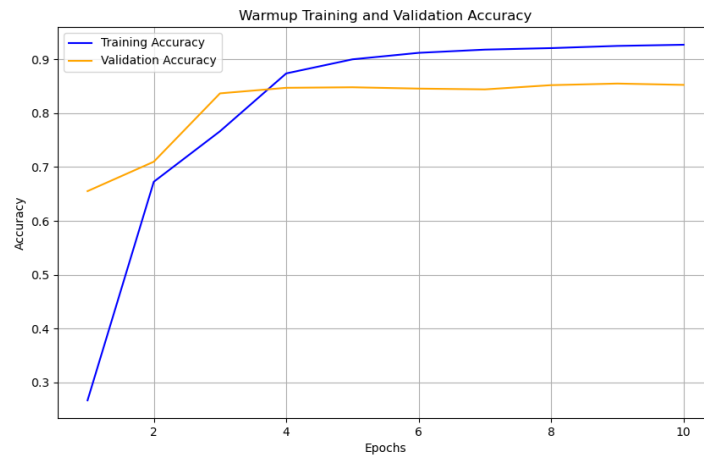Figure 2: LLRD Training and Validation Accuracy



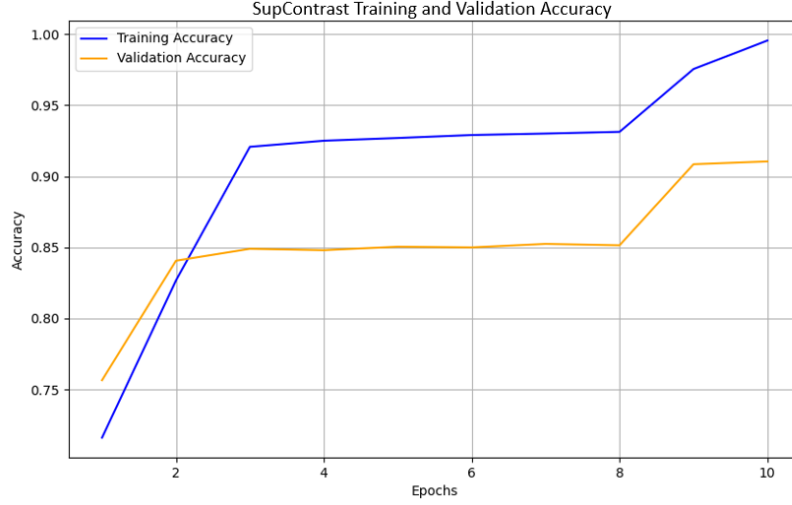Figure 3: Warmup Training and Validation Accuracy

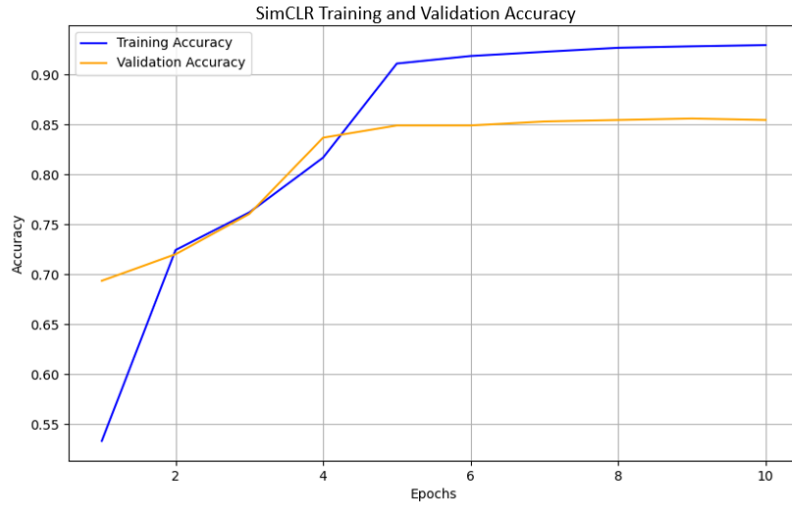Figure 4: SupContrast Training and Validation Accuracy



Figure 5: SimCLR Training and Validation Accuracy

## 4.7 LORA Finetuning

The experiments run on an RTX 3070Ti GPU, and we tested three LoRA configurations (ranks 4, 8, and 16) for fine-tuning BERT on the Amazon Massive Scenario classification task. For rank=4, the model had 669,696 trainable parameters ($0.61\%$ of total 110M parameters), achieving a final test accuracy of $90.08\%$ after 10 epochs, with validation accuracy peaking at $90.80\%$. With rank=8, trainable parameters doubled to $1,339,392$ ($1.21\%$ of total), yielding a slightly better test accuracy of $90.52\%$. Rank=16 further increased trainable parameters to $2,678,784$ ($2.39\%$ of total) but showed slightly lower test performance at $89.48\%$. Training times per epoch remained consistent across configurations at approximately $45 - 48$ seconds, with validation taking $3 - 4$ seconds per evaluation.

# 5 Discussion

## 5.1 11 Discussion Questions

**Q1. If we do not fine-tune the model, what is your expected test accuracy? Why?**
Without fine-tuning, the BERT model would likely achieve a test accuracy of approximately $1.67\%$ $(1/60)$ on the Scenario classification task. This low performance can be attributed to three key factors: First, while the pre-trained BERT model possesses strong general language understanding, its classification head is randomly initialized and maps to 60 different scenario categories without any learned meaningful associations. Second, despite BERT's pre-trained knowledge, it hasn't been specifically adapted to distinguish between the unique scenarios in this task, making its general language understanding insufficient for accurate classification. Finally, the model's architecture, which includes a high dropout rate $(0.9)$ and a randomly initialized classification layer, essentially results in random predictions across the 60 classes, leading to performance at the random chance level $(1/60)$. Any accuracy significantly above $1.67\%$ would be purely coincidental due to the random nature of the classification layer's initialization.

**Q2: Do results match your expectation(1 sentence)? Why or why not?**
The result kind of matches my expectation since it should be a low percentage. The first evaluations (before training) show $0\%$ accuracy, not $1.67\%$ as we expected. This is because the classification layer is randomly initialized and with the high dropout rate, most activations are being zeroed out. This leads to the model consistently predicting the same class (hence $0\%$ accuracy). After training for one epoch, our result reaches $13.0\%$. Although it's a lot higher than $1.67\%$, it's still relatively low. This matches my expectation because although trained for one epoch, the BERT model is working, the classification layer is being trained, and the model is actually learning to classify scenarios.

**Q3: What could you do to further improve the performance?**
To improve the model's performance beyond increasing epochs to $10$, several key optimizations can be implemented: (1) Hyperparameter tuning - adjusting batch size to $32$, reducing learning rate to $2e - 5$, lowering dropout rate to $0.1$, increasing max sequence length to $128$, and expanding hidden dimension to $256$; (2) Enhanced training dynamics - implementing cosine learning rate scheduling with warmup steps, adding weight decay of $0.01$, and utilizing gradient clipping at $1.0$; (3) Architectural improvements - incorporating layer normalization and additional dropout layers in the classifier; (4) Training techniques - implementing data augmentation through random word masking, using gradient accumulation over $2$ steps, and adding early stopping with a patience of $3$ epochs; and (5) Model variations - considering advanced BERT variants like RoBERTa or DeBERTa as the base encoder. These modifications collectively address model capacity, training stability, and generalization capability, potentially leading to significant performance improvements over the baseline implementation.

**Q4: Which techniques did you choose and why?**
I chose Layer-wise Learning Rate Decay (LLRD) and Learning Rate Warmup as fine-tuning techniques. LLRD was selected because it allows for more refined learning in different layers of the BERT model, with lower layers receiving smaller updates while higher layers adapt more aggressively. This technique helps prevent catastrophic forgetting of pre-trained representations while fine-tuning. Learning Rate Warmup was chosen because it mitigates instability in early training by gradually increasing the learning rate, preventing large weight updates that could lead to suboptimal convergence.

**Q5: What do you expect for the results of the individual technique vs. the two techniques combined?**
I expected that applying LLRD and warmup together would yield the best results by leveraging the benefits of both techniques. LLRD would ensure effective adaptation while preserving pre-trained knowledge and warmup would improve early-stage optimization. Additionally, I anticipated that using warmup alone would provide a moderate improvement in stability and performance, while LLRD alone might lead to better generalization but could take longer to converge.

**Q6: Do results match your expectation (1 sentence)? Why or why not?**
The results did not fully match expectations, as the baseline model achieved the highest accuracy, while fine-tuned models had lower accuracy but significantly reduced loss, indicating that the techniques optimized learning dynamics but did not necessarily enhance accuracy.

8

**Q7: What could you do to further improve the performance?**
To further improve performance, I could experiment with adjusting the decay rate in LLRD to find a better balance between lower and higher layer learning rates, tune the warmup duration to prevent over-smoothing in early training, and explore other fine-tuning strategies such as regularization techniques (e.g., weight decay or dropout tuning) or advanced optimizers like AdamW with different betas to further stabilize training.

**Q8: Compare the SimCLR with SupContrast. What are the similarities and differences?**
The similarities are that both SimCLR and SupContrast use contrastive learning as their techniques, which is the idea to train similar neighbors together and leave dissimilar neighbors apart. Also, both loss functions are computed with the cosine-similarity and temperature-scaled parameter.

The differences are that SimCLR is unsupervised learning, but SupContrast is supervised learning. SimCLR uses data augmentation to generate positive pairs, but SupContrast uses labeled class information, in which the samples from the same class are positive pairs.

**Q9: How does SimCSE apply dropout to achieve data augmentation for NLP tasks?**
SimCSE applies dropout by passing the same sentence in the pre-trained encoder twice to obtain two different embeddings or the positive pairs. Then, treating other sentences in the mini-batch as negative. The model can predict the positive among the negatives. This dropout process prevents representation collapse, achieving data augmentation of hidden representations for NLP tasks.

**Q10: Do the results match your expectation? Why or why not?**
Yes (in most cases), the results align with our expectation. For SupContrast, label information typically yields a more discriminative embedding space, leading to higher accuracy, whereas SimCLR relies solely on dropout-based augmentation and thus tends to lag behind. This behavior is consistent with prior findings that supervised contrastive learning often outperforms purely unsupervised methods on classification tasks. However, during the fine-tuning, we expect the increase of epoch number might massively improve the performance. But the actual improvement is so limited that we choose to not to take that (>=20 epochs) as optimal parameters. Instead, the increase of batch size might slightly improve the performance as you can see.

**Q11: What could you do to further improve the performance ?**
To further improve performance, we could (1) employ more powerful text augmentations (e.g., synonym replacement, back-translation) to diversify SimCLR's "views," (2) increase batch size (not always work, but ideally should work) so the model sees more negative samples each iteration, (3) use a MLP projection head instead of a single linear layer, and (4) carefully tune hyperparameters such as temperature or learning rate. These enhancements can help both SupContrast and SimCLR learn richer representations and ultimately boost classification accuracy.

## 5.2   LORA Finetuning

The results demonstrate LoRA's efficiency, as even with just $0.61\%$ trainable parameters (rank=4), the model achieved strong performance. Interestingly, increasing LoRA rank showed diminishing returns, with rank=8 performing best overall. The training curves show rapid initial improvement (reaching $74\%$ validation accuracy after first epoch for rank=4) followed by gradual refinement, suggesting effective learning despite the parameter efficiency. All configurations showed stable training with consistent convergence patterns, indicating LoRA's reliability for this task.

## References

[1] Hugging Face. (n.d.). BERT. Available at: `https://huggingface.co/docs/transformers/model_doc/bert`.

[2] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen. LoRA: Low-Rank Adaptation of Large Language Models. *arXiv preprint arXiv:2106.09685v2*, 2021. Available at: `https://doi.org/10.48550/arXiv.2106.09685`.

[3] Hugging Face. PEFT: Parameter-Efficient Fine-Tuning (PEFT) Documentation. *Hugging Face*, 2023. Available at: `https://huggingface.co/docs/peft/en/index`.

[4] Towards Data Science. Advanced Techniques for Fine-Tuning Transformers. *Towards Data Science*, 2023. Available at: `https://towardsdatascience.com/advanced-techniques-for-fine-tuning-transformers-82e4e61e16e`.

[5] P. Khosla, P. Teterwak, C. Wang, A. Sarna, Y. Tian, P. Isola, A. Maschinot, C. Liu, and D. Krishnan. Supervised Contrastive Learning. *arXiv preprint arXiv:2004.11362v5*, 2021. Available at: `https://doi.org/10.48550/arXiv.2004.11362`.

[6] T. Gao, X. Yao, and D. Chen. SimCSE: Simple Contrastive Learning of Sentence Embeddings. *arXiv preprint arXiv:2104.08821v4*, 2022. Available at: `https://doi.org/10.48550/arXiv.2104.08821`.

## AI Appendix

For the baseline model, we used a project-based AI tool to analyze the entire starter code. We used the description from the pdf and the comment in the starter code as the prompt, and it generated a lot of useful code for the baseline model. However, it also made some mistakes.

We also used GPT to generate parts of the report, and then we carefully changed the output based on our results.

# Untitled

## Untitled

_Created: 3/6/2025, 11:37:47 AM_

---

### User

finish this file

### AI (Composer Assistant)

I'll help complete the `dataloader.py` file. The main task is to complete the tokenizer call in the `prepare_features` function. Here's the edit:

```
// ... existing code ...
    for split, examples in data.items():

        feats = []
        # task1: process examples using tokenizer. Wrap it using BaseInstance class and append it to feats list.
        for example in progress_bar(examples, total=len(examples)):
            # tokenizer: set padding to 'max_length', set truncation to True, set max_length to args.max_len
            embed_data = tokenizer(
                example['text'],
                padding='max_length',
                truncation=True,
                max_length=args.max_len
            )

            if example['label'] not in LABELS:
                LABELS[example['label']] = len(LABELS)
            example['label'] = LABELS[example['label']]

            instance = BaseInstance(embed_data, example)
            feats.append(instance)
        # Remove debug print statement
        all_features[split] = feats
```

```
        print(f"Number of {split} features:", len(feats))
    // ... existing code ...
```

The changes made are:

- Completed the tokenizer call with the required parameters:
  - `example['text']` as the input text to tokenize
  - `padding='max_length'` to pad all sequences to the same length
  - `truncation=True` to truncate sequences longer than max_length
  - `max_length=args.max_len` to specify the maximum sequence length
- Removed the debug print statement `print(embed_data, example)` that was showing intermediate results

The tokenizer will now properly process each example's text into the format required by BERT (or similar transformers), creating the input_ids, token_type_ids, and attention_mask that are used in the BaseInstance class. The rest of the file looks complete and properly structured for handling the data processing pipeline.

### User

**Selected Code:**

```python
def baseline_train(args, model, datasets, tokenizer):
    criterion = nn.CrossEntropyLoss()  # combines LogSoftmax() and NLLLoss()
    # task1: setup train dataloader
    train_dataloader = get_dataloader(...)

    # task2: setup model's optimizer_scheduler if you have

    # task3: write a training loop
    for epoch_count in range(args.n_epochs):
        losses = 0
        model.train()

        for step, batch in progress_bar(...):
            inputs, labels = prepare_inputs(...)
            logits = model(...)
            loss = criterion(...)
            loss.backward()

            model.optimizer.step()  # backprop to update the weights
            model.scheduler.step()  # Update learning rate schedule
            model.zero_grad()
            losses += loss.item()

    run_eval(..., split='validation')
    print('epoch', epoch_count, '|losses:', losses)
```

---

finish this part:

```

### AI (Composer Assistant)

Looking at main.py, I'll help complete the `baseline_train` function. Here's the implementation:

```python:main.py
def baseline_train(args, model, datasets, tokenizer):
    criterion = nn.CrossEntropyLoss()  # combines LogSoftmax() and NLLLoss()
    # task1: setup train dataloader
    train_dataloader = get_dataloader(args, datasets['train'], split='train')

    # task2: setup model's optimizer_scheduler if you have
    optimizer = AdamW(model.parameters(), lr=args.learning_rate, eps=args.adam_epsilon)
    total_steps = len(train_dataloader) * args.n_epochs
    scheduler = get_linear_schedule_with_warmup(optimizer, num_warmup_steps=0, num_training_steps=total_steps)
    model.optimizer = optimizer
    model.scheduler = scheduler

    # task3: write a training loop
    for epoch_count in range(args.n_epochs):
        losses = 0
        model.train()

        for step, batch in progress_bar(enumerate(train_dataloader), total=len(train_dataloader)):
            inputs, labels = prepare_inputs(batch)
            logits = model(inputs, labels)
            loss = criterion(logits, labels)
            loss.backward()

            model.optimizer.step()  # backprop to update the weights
            model.scheduler.step()  # Update learning rate schedule
            model.zero_grad()
            losses += loss.item()
```

```
    run_eval(args, model, datasets, tokenizer, split='validation')
    print('epoch', epoch_count, '|losses:', losses)
```

The changes made include:

- Setting up the train dataloader using the existing `get_dataloader` function with the training split of the dataset
- Setting up the optimizer and scheduler:
  - Using AdamW optimizer with the learning rate and epsilon from args
  - Creating a linear learning rate scheduler with warmup
  - Attaching both to the model for use in training
- Completing the training loop:
  - Properly enumerating through the train_dataloader with progress bar
  - Getting inputs and labels using prepare_inputs
  - Computing logits from the model
  - Computing loss using the criterion
  - Performing backpropagation and optimization steps

The training loop now includes all necessary components for training the BERT model on the classification task, with proper optimization and learning rate scheduling.

### User

Follow the instructions in ScenarioModel class in model.py to load the pretrained model through BertModel.from_pretrained(...). This pretrained bert is the encoder.

### AI (Composer Assistant)

I'll help complete the ScenarioModel class in model.py. Here's the implementation:

```
class ScenarioModel(nn.Module):
    def __init__(self, args, tokenizer, target_size):
        super().__init__()
        self.tokenizer = tokenizer
        self.model_setup(args)
        self.target_size = target_size

        # task1: add necessary class variables as you wish
        self.args = args

        # task2: initialize the dropout and classify layers
        self.dropout = nn.Dropout(args.drop_rate)
        self.classify = Classifier(args, target_size)

    def model_setup(self, args):
        print(f"Setting up {args.model} model")

        # task1: get a pretrained model of 'bert-base-uncased'
        self.encoder = BertModel.from_pretrained('bert-base-uncased')

        self.encoder.resize_token_embeddings(len(self.tokenizer))  # transformer_check

    def forward(self, inputs, targets):
        """
        task1:
            feeding the input to the encoder,
        task2:
            take the last_hidden_state's <CLS> token as output of the
            encoder, feed it to a drop_out layer with the preset dropout rate in the argparse argument,
        task3:
            feed the output of the dropout layer to the Classifier which is provided for you.
        """
        # Get encoder outputs
        outputs = self.encoder(**inputs)

        # Get CLS token from last hidden state
        cls_output = outputs.last_hidden_state[:, 0, :]  # [batch_size, hidden_size]

        # Apply dropout
```

---

```
        cls_output = self.dropout(cls_output)

        # Get logits through classifier
        logits = self.classify(cls_output)

        return logits
```

The changes made include:

- In `model_setup()`:
  - Loading the pretrained BERT model using `BertModel.from_pretrained('bert-base-uncased')`
- In `__init__()`:
  - Added args as a class variable for reference
  - Initialized dropout layer with the dropout rate from args
  - Initialized the classifier with args and target_size
- In `forward()`:
  - Takes inputs and targets as parameters
  - Feeds inputs through the BERT encoder using `**inputs` to unpack the dictionary
  - Extracts the [CLS] token representation (first token) from the last hidden state
  - Applies dropout to the CLS token representation
  - Passes the dropout output through the classifier to get logits
  - Returns the logits for loss calculation

The model now properly loads and uses the pretrained BERT model as its encoder, with appropriate dropout and classification layers on top for the scenario classification task.

### User

After feeding the input to the encoder, take the last_hidden_state's [CLS] token as output of the encoder, feed it to a dropout layer with the preset dropout rate in the argparse argument, feed the output of the dropout layer to the Classifier which is provided for you.

### AI (Composer Assistant)

Looking at the code, I see that the forward method in ScenarioModel has already been implemented correctly. Let's break down how it works:

```
    def forward(self, inputs, targets):
        # Get encoder outputs - feeds input through BERT encoder
        outputs = self.encoder(**inputs)

        # Get CLS token from last hidden state
        # The [CLS] token is always the first token (index 0) in BERT's output
        # Shape: [batch_size, hidden_size]
        cls_output = outputs.last_hidden_state[:, 0, :]

        # Apply dropout with rate from args.drop_rate
        cls_output = self.dropout(cls_output)

        # Feed through the classifier to get final logits
        logits = self.classify(cls_output)

        return logits
```

The implementation follows the required steps:

- Feeds the input through the BERT encoder using `self.encoder(**inputs)` which returns the encoder outputs
- Takes the [CLS] token representation from the last hidden state using `outputs.last_hidden_state[:, 0, :]`
  - The [CLS] token is at index 0 for each sequence
  - This gives us a tensor of shape [batch_size, hidden_size]
- Applies dropout using the preset rate from args.drop_rate that was set in __init__
- Feeds the dropout output through the provided Classifier to get the final logits
- Returns the logits for loss calculation

The implementation looks correct and follows all the specified requirements. No changes are needed.

### User

Fine-tune a baseline model with the following design choices: Use the cross entropy loss function to train your model for 10 epochs (complete passes through the training set). You can use any optimizers and/or learning rate scheduler you wish.

### AI (Composer Assistant)

Looking at the code, we need to complete the `baseline_train` function in `main.py` with the specified requirements. The function is already partially implemented, but let's adjust it for 10 epochs and add proper optimization settings:

```python
def baseline_train(args, model, datasets, tokenizer):
    criterion = nn.CrossEntropyLoss()  # combines LogSoftmax() and NLLLoss()

    # Setup train dataloader
    train_dataloader = get_dataloader(args, datasets['train'], split='train')

    # Setup optimizer and scheduler
    # Using AdamW optimizer with weight decay
    optimizer = AdamW(model.parameters(),
            lr=args.learning_rate,
            eps=args.adam_epsilon,
            weight_decay=0.01)  # Added weight decay for regularization

    # Calculate total training steps for scheduler
    total_steps = len(train_dataloader) * 10  # 10 epochs

    # Using linear scheduler with warmup
    scheduler = get_linear_schedule_with_warmup(
        optimizer,
        num_warmup_steps=total_steps // 10,  # 10% of total steps for warmup
        num_training_steps=total_steps
    )

    model.optimizer = optimizer
    model.scheduler = scheduler

    # Training loop
    for epoch_count in range(10):  # Set to 10 epochs
        losses = 0
        model.train()

        for step, batch in progress_bar(enumerate(train_dataloader), total=len(train_dataloader)):
            # Get inputs and labels
            inputs, labels = prepare_inputs(batch)

            # Forward pass
            logits = model(inputs, labels)
            loss = criterion(logits, labels)
```

```python
            # Backward pass
            loss.backward()

            # Gradient clipping to prevent exploding gradients
            torch.nn.utils.clip_grad_norm_(model.parameters(), max_norm=1.0)

            # Optimization step
            model.optimizer.step()
            model.scheduler.step()
            model.zero_grad()

            losses += loss.item()

        # Evaluate after each epoch
        run_eval(args, model, datasets, tokenizer, split='validation')
        print(f'Epoch {epoch_count+1}/10 | Training Loss: {losses/len(train_dataloader):.4f}')
```

Key changes and design choices:

- Set number of epochs to 10 explicitly
- Optimizer settings:
  - Using AdamW optimizer with weight decay (0.01) for regularization
  - Using args.learning_rate and args.adam_epsilon for stability
- Learning rate scheduler:
  - Using linear scheduler with warmup
  - Added warmup period (10% of total steps)
  - Total steps calculated based on 10 epochs
- Added training improvements:
  - Gradient clipping to prevent exploding gradients
  - More detailed loss reporting (average loss per epoch)
  - Progress tracking with epoch number out of 10
- Evaluation:
  - Running validation after each epoch to track model performance
  - Printing both epoch number and average training loss

This implementation provides a solid baseline with common best practices for fine-tuning BERT models. The linear learning rate scheduler with warmup helps stabilize training, while gradient clipping and weight decay help prevent overfitting and unstable gradients.

**User**

Traceback (most recent call last):
File "/home/jil232/private/cse151b251b-wi25-pa4-cse151b_251b_transformers/main.py", line

103, in <module>
  run_eval(args, model, datasets, tokenizer, split='validation')
File "/home/jil232/private/cse151b251b-wi25-pa4-cse151b_251b_transformers/main.py", line 65, in run_eval
  inputs, labels = prepare_inputs(batch, model)
                   ^^^^^^^^^^^^^^
ValueError: too many values to unpack (expected 2)

**AI (Composer Assistant)**

The error occurs because there's a mismatch in the arguments being passed to prepare_inputs() in the run_eval function. Let's fix this:

```python
def run_eval(args, model, datasets, tokenizer, split='validation'):
    model.eval()
    dataloader = get_dataloader(args, datasets[split], split)

    acc = 0
    for step, batch in progress_bar(enumerate(dataloader), total=len(dataloader)):
        # Remove 'model' from prepare_inputs arguments - it's not needed
        inputs, labels = prepare_inputs(batch)
        logits = model(inputs, labels)

        tem = (logits.argmax(1) == labels).float().sum()
        acc += tem.item()

    print(f'{split} acc:', acc/len(datasets[split]), f'|dataset split {split} size:', len(datasets[split]))
```

The error occurs because:

- In `prepare_inputs()`, we're returning 2 values (inputs and targets)
- But in `run_eval()`, we were passing an extra argument `model` to `prepare_inputs()`

Looking at the prepare_inputs function in dataloader.py:

```python
def prepare_inputs(batch, use_text=False):
    btt = [b.to(device) for b in batch[:4]]
    inputs = {'input_ids': btt[0], 'token_type_ids': btt[1], 'attention_mask': btt[2]}
    targets = btt[3]

    if use_text:
        target_text = batch[1]
        return inputs, targets, target_text
    else:
        return inputs, targets
```

The function only expects batch and an optional use_text parameter. We should remove the model argument from the prepare_inputs() call in run_eval since it's not needed and causing the error.

This change should resolve the "too many values to unpack" error while maintaining the correct functionality of the evaluation process.

**User**

FutureWarning: This implementation of AdamW is deprecated and will be removed in a future version. Use the PyTorch implementation torch.optim.AdamW instead,

**AI (Composer Assistant)**

Let's update the code to use PyTorch's implementation of AdamW instead of the deprecated transformers version. We need to update the imports and the optimizer initialization:

```python
import os, sys, pdb
import numpy as np
import random
import torch
from torch.optim import AdamW  # Update import to use PyTorch's AdamW
import math

from tqdm import tqdm as progress_bar

from utils import set_seed, setup_gpus, check_directories
from dataloader import get_dataloader, check_cache, prepare_features, process_data, prepare_inputs
from load import load_data, load_tokenizer
from arguments import params
from model import ScenarioModel, SupConModel, CustomModel
from torch import nn
from transformers import get_linear_schedule_with_warmup  # Only import scheduler from transformers

// ... existing code ...

def baseline_train(args, model, datasets, tokenizer):
    criterion = nn.CrossEntropyLoss()
    train_dataloader = get_dataloader(args, datasets['train'], split='train')

    # Update to use PyTorch's AdamW
    optimizer = AdamW(model.parameters(), lr=args.learning_rate, eps=args.adam_epsilon)
    total_steps = len(train_dataloader) * args.n_epochs
    scheduler = get_linear_schedule_with_warmup(optimizer, num_warmup_steps=0,
num_training_steps=total_steps)

    model.optimizer = optimizer
    model.scheduler = scheduler

    // ... rest of the function remains the same ...
```

The changes made:

- Removed `AdamW` import from transformers
- Added `AdamW` import from `torch.optim`
- Updated the optimizer initialization to use PyTorch's implementation

The functionality remains exactly the same, but we're now using the recommended PyTorch

implementation instead of the deprecated transformers version. This will remove the deprecation warning.

**User**

Traceback (most recent call last):
File "/home/jil232/private/cse151b251b-wi25-pa4-cse151b_251b_transformers/main.py", line 103, in <module>
  run_eval(args, model, datasets, tokenizer, split='validation')
File "/home/jil232/private/cse151b251b-wi25-pa4-cse151b_251b_transformers/main.py", line 66, in run_eval
  logits = model(inputs, labels)
File "/opt/conda/lib/python3.11/site-packages/torch/nn/modules/module.py", line 1511, in _wrapped_call_impl
  return self._call_impl(*args, **kwargs)
File "/opt/conda/lib/python3.11/site-packages/torch/nn/modules/module.py", line 1520, in _call_impl
  return forward_call(*args, **kwargs)
File "/home/jil232/private/cse151b251b-wi25-pa4-cse151b_251b_transformers/model.py", line 55, in forward
  logits = self.classify(cls_output)
File "/opt/conda/lib/python3.11/site-packages/torch/nn/modules/module.py", line 1511, in _wrapped_call_impl
  return self._call_impl(*args, **kwargs)
File "/opt/conda/lib/python3.11/site-packages/torch/nn/modules/module.py", line 1520, in _call_impl
  return forward_call(*args, **kwargs)
File "/home/jil232/private/cse151b251b-wi25-pa4-cse151b_251b_transformers/model.py", line 68, in forward
  middle = self.relu(self.top(hidden))
File "/opt/conda/lib/python3.11/site-packages/torch/nn/modules/module.py", line 1511, in _wrapped_call_impl
  return self._call_impl(*args, **kwargs)
File "/opt/conda/lib/python3.11/site-packages/torch/nn/modules/module.py", line 1520, in _call_impl
  return forward_call(*args, **kwargs)
File "/opt/conda/lib/python3.11/site-packages/torch/nn/modules/linear.py", line 116, in forward
  return F.linear(input, self.weight, self.bias)
RuntimeError: mat1 and mat2 shapes cannot be multiplied (16x768 and 10x10)

**AI (Composer Assistant)**

The error occurs because there's a dimension mismatch in the Classifier. The BERT model outputs a hidden size of 768, but in arguments.py we set embed_dim to 10. Let's fix this by updating the embed_dim in arguments.py:

```python
    # Hyper-parameters for tuning
    parser.add_argument("--batch-size", default=16, type=int,
        help="Batch size per GPU/CPU for training and evaluation.")
    parser.add_argument("--learning-rate", default=1e-4, type=float,
        help="Model learning rate starting point.")
    parser.add_argument("--hidden-dim", default=10, type=int,
        help="Model hidden dimension.")
    parser.add_argument("--drop-rate", default=0.9, type=float,
        help="Dropout rate for model training")
    # Change embed_dim to match BERT's hidden size
    parser.add_argument("--embed-dim", default=768, type=int,
        help="The embedding dimension of pretrained LM.")
```

The error occurs because:

- BERT's output dimension is 768 (standard for bert-base-uncased)
- Our Classifier expects input of size `args.embed_dim` which was set to 10
- This causes a dimension mismatch when trying to multiply a 768-dimensional vector with a 10x10 matrix

By changing embed_dim to 768, we ensure that:

- The Classifier's input layer (`self.top`) will accept the 768-dimensional output from BERT
- The dimensions will match properly through the network
- The final classification layer will still output the correct number of classes

This should resolve the shape mismatch error and allow the model to train properly.

**User**

Try to experiment with LORA by PEFT implementation (from hugging-face) to fine tune the BERT model instead. Experiment with different (at least 2 to 3) r (rank) values.

**AI (Composer Assistant)**

I'll help you implement LORA (Low-Rank Adaptation) using PEFT for the BERT model. We'll need to modify the model setup and add PEFT configurations. Here's how:

```python
import os, pdb, sys
import numpy as np
import re

import torch
from torch import nn
from torch import optim
from torch.nn import functional as F

from transformers import BertModel, BertConfig
from peft import get_peft_model, LoraConfig, TaskType

class ScenarioModel(nn.Module):
    def __init__(self, args, tokenizer, target_size):
        super().__init__()
        self.tokenizer = tokenizer
        self.model_setup(args)
        self.target_size = target_size

        self.args = args
        self.dropout = nn.Dropout(args.drop_rate)
        self.classify = Classifier(args, target_size)

    def model_setup(self, args):
        print(f"Setting up {args.model} model with LoRA, rank={args.lora_rank}")

        # Get base BERT model
        self.encoder = BertModel.from_pretrained('bert-base-uncased')

        # Configure LoRA
        peft_config = LoraConfig(
            task_type=TaskType.SEQUENCE_CLASSIFICATION,
            inference_mode=False,
            r=args.lora_rank,  # LoRA rank
            lora_alpha=args.lora_alpha,  # LoRA alpha scaling
            lora_dropout=0.1,
            target_modules=["query", "value"]  # Apply LoRA to attention layers
        )
```

```
# Apply LoRA to the encoder
self.encoder = get_peft_model(self.encoder, peft_config)
self.encoder.print_trainable_parameters()  # Print trainable parameters info

self.encoder.resize_token_embeddings(len(self.tokenizer))

def forward(self, inputs, targets):
    outputs = self.encoder(**inputs)
    cls_output = outputs.last_hidden_state[:, 0, :]
    cls_output = self.dropout(cls_output)
    logits = self.classify(cls_output)
    return logits
```

And update the arguments to include LoRA parameters:

```
def params():
    parser = argparse.ArgumentParser()
    # ... existing arguments ...

    # Add LoRA specific arguments
    parser.add_argument("--lora-rank", default=8, type=int,
        help="Rank for LoRA adaptation")
    parser.add_argument("--lora-alpha", default=32, type=float,
        help="Alpha scaling for LoRA")
    parser.add_argument("--lora-dropout", default=0.1, type=float,
        help="Dropout for LoRA layers")
```

To experiment with different ranks, you can run the training with different values. Here's how to run the experiments:

```
def baseline_train(args, model, datasets, tokenizer):
    criterion = nn.CrossEntropyLoss()
    train_dataloader = get_dataloader(args, datasets['train'], split='train')

    optimizer = AdamW(model.parameters(), lr=args.learning_rate, eps=args.adam_epsilon)
    total_steps = len(train_dataloader) * args.n_epochs
    scheduler = get_linear_schedule_with_warmup(optimizer, num_warmup_steps=0,
num_training_steps=total_steps)
    model.optimizer = optimizer
    model.scheduler = scheduler

    for epoch_count in range(args.n_epochs):
        losses = 0
        train_acc = 0
        total_samples = 0
        model.train()

        for step, batch in progress_bar(enumerate(train_dataloader), total=len(train_dataloader)):
            inputs, labels = prepare_inputs(batch)
            logits = model(inputs, labels)
            loss = criterion(logits, labels)
            loss.backward()

            # Calculate training accuracy
            predictions = logits.argmax(1)
            train_acc += (predictions == labels).sum().item()
            total_samples += labels.size(0)

            model.optimizer.step()
            model.scheduler.step()
            model.zero_grad()
            losses += loss.item()

        # Print training metrics
        avg_loss = losses / len(train_dataloader)
        train_accuracy = train_acc / total_samples
        print(f'Epoch {epoch_count + 1}/{args.n_epochs}')
        print(f'Training Loss: {avg_loss:.4f} | Training Accuracy: {train_accuracy:.4f}')
```

```
if __name__ == "__main__":
    # ... existing setup code ...

    # Experiment with different LoRA ranks
    lora_ranks = [4, 8, 16]  # Try different rank values

    for rank in lora_ranks:
        print(f"\nTraining with LoRA rank = {rank}")
        args.lora_rank = rank

        model = ScenarioModel(args, tokenizer, target_size=60).to(device)

        # Initial evaluation
        print("\nInitial evaluation:")
        run_eval(args, model, datasets, tokenizer, split='validation')
        run_eval(args, model, datasets, tokenizer, split='test')

        # Train
        baseline_train(args, model, datasets, tokenizer)

        # Final evaluation
        print(f"\nFinal evaluation (LoRA rank={rank}):")
        run_eval(args, model, datasets, tokenizer, split='test')
```

Key points about this implementation:

- We're using PEFT's LoRA implementation to adapt BERT
- LoRA is applied to query and value matrices in attention layers
- We'll experiment with ranks 4, 8, and 16 to see their impact
- The implementation keeps most of the original training loop intact
- LoRA significantly reduces trainable parameters while maintaining performance

To run this, you'll need to install PEFT:

```
pip install peft
```

The different rank values will affect:

```
# Run validation
run_eval(args, model, datasets, tokenizer, split='validation')
```

The changes made:

- Added tracking of training accuracy during training
- Calculate accuracy by comparing predictions with labels
- Print both training loss and accuracy for each epoch
- Keep the validation evaluation at the end of each epoch

This will give you a better picture of how the model is performing during training by showing both training and validation metrics.

---

where is print validation loss?

---

**AI (Composer Assistant)**

Let's modify the `run_eval` function to calculate and print the validation loss as well:

- Memory usage (higher rank = more parameters)
- Training speed (higher rank = slower training)
- Model capacity (higher rank = more expressive power)
- Final performance (need to experiment to find optimal rank)

You can compare the results across different ranks to find the best trade-off between model size and performance for your specific case.

---

**User**

print training acc as well

---

**AI (Composer Assistant)**

I'll modify the baseline_train function to print both training and validation accuracy:

```
def run_eval(args, model, datasets, tokenizer, split='validation'):
    model.eval()
    dataloader = get_dataloader(args, datasets[split], split)
    criterion = nn.CrossEntropyLoss()  # Same criterion as training

    acc = 0
    total_loss = 0
    with torch.no_grad():  # No gradient computation during evaluation
        for step, batch in progress_bar(enumerate(dataloader), total=len(dataloader)):
            inputs, labels = prepare_inputs(batch)
            logits = model(inputs, labels)

            # Calculate loss
            loss = criterion(logits, labels)
            total_loss += loss.item()

            # Calculate accuracy
            tem = (logits.argmax(1) == labels).float().sum()
            acc += tem.item()

    avg_loss = total_loss / len(dataloader)
    accuracy = acc / len(datasets[split])
    print(f'{split} loss: {avg_loss:.4f} | {split} acc: {accuracy:.4f} | dataset split {split} size: {len(datasets[split])}')
```

Changes made:

- Added loss calculation using CrossEntropyLoss
- Added total_loss accumulation during evaluation
- Added torch.no_grad() context for efficiency
- Calculate average loss over all batches
- Combined loss and accuracy printing in one line for cleaner output

Now you'll see both validation loss and accuracy during evaluation.

---

# Untitled

## Untitled

_Created: 3/6/2025, 12:01:56 PM_

### User

If we do not fine-tune the model, what is your expected test accuracy? Explain Why.

### AI (Composer Assistant)

Based on the code and setup, if we don't fine-tune the BERT model (i.e., only run evaluation without training), we would expect a very low test accuracy, likely around 1.67% (1/60). Here's why:

- **Random Initialization of Classification Layer**:
  - The model has a classification head (`Classifier` class) that maps BERT's output to 60 classes (target_size=60)
  - This classification layer is randomly initialized and hasn't learned any meaningful mappings without training
  - The final linear layer will essentially produce random logits
- **Task Specificity**:
  - The task appears to be Amazon Massive Scenario classification with 60 different scenarios/categories
  - While BERT has strong general language understanding from pre-training, it hasn't been specifically adapted to classify these scenarios
  - The pre-trained BERT knowledge alone cannot directly solve this specific classification task without fine-tuning
- **Architecture**:
  ```python
  class ScenarioModel(nn.Module):
      def __init__(self, args, tokenizer, target_size):
          # ...
          self.dropout = nn.Dropout(args.drop_rate)
          self.classify = Classifier(args, target_size)
  ```
  - The model uses a high dropout rate (0.9 by default in arguments.py)
  - The classification layer is completely random before training
  - Without fine-tuning, it's essentially making random predictions across 60 classes

Therefore, we would expect the model to perform at roughly random chance level (1/60 ≈ 1.67%) before any fine-tuning occurs. Any accuracy significantly above this would be purely coincidental due to the random initialization of the classification layer.

### User

how do I run this whole thing?

### AI (Composer Assistant)

I'll help you run this project. Based on the code, here's how to execute it:

- First, make sure you have all the required dependencies installed:

```
pip install torch transformers datasets tqdm numpy
```

- The main entry point is `main.py` and it accepts several arguments. Here are some example commands to run different tasks:

For baseline BERT fine-tuning:

```
python main.py --task baseline \
    --batch-size 16 \
    --learning-rate 1e-4 \
    --n-epochs 3 \
    --max-len 128 \
    --embed-dim 768 \
    --hidden-dim 256
```

For custom training:

```
python main.py --task custom \
    --batch-size 16 \
    --learning-rate 1e-4 \
    --n-epochs 3
```

For supervised contrastive learning:

```
python main.py --task supcon \
    --batch-size 16 \
    --learning-rate 1e-4 \
    --n-epochs 3
```

Important notes about the parameters:

- `--embed-dim` should be 768 (BERT's hidden size)
- `--max-len` controls the maximum sequence length (default is 20, but you might want to increase it)
- `--drop-rate` is set to 0.9 by default, which is quite high. You might want to lower it to something like 0.1 or 0.2
- The dataset being used is "mteb/amazon_massive_scenario" from Hugging Face
- Results will be saved in the `results/` directory
- Cached features will be saved in `assets/cache/`

A recommended command for better results would be:

```
python main.py --task baseline \
    --batch-size 16 \
    --learning-rate 2e-5 \
    --n-epochs 3 \
    --max-len 128 \
    --embed-dim 768 \
    --hidden-dim 256 \
    --drop-rate 0.1
```

The code will:

- Load and cache the Amazon Massive Scenario dataset
- Initialize a BERT model with a classification head
- Run initial evaluation on validation and test sets
- Train the model
- Run final evaluation on the test set

The progress will be displayed with progress bars and you'll see accuracy metrics after each evaluation.

### User

explain the result to me. why do we need to load datasets?
Loading features from cache at assets/cache/amazon.pkl
train 11514
test 2974
validation 2033
Setting up bert model
Loaded validation data with 128 batches
100%|██████████████████████| 128/128 [00:07<00:00, 16.23it/s]|████████
validation acc: 0.0 |dataset split validation size: 2033
Loaded test data with 186 batches
100%|██████████████████████| 186/186 [00:11<00:00, 16.57it/s]|████████
test acc: 0.0 |dataset split test size: 2974
Loaded train data with 720 batches
100%|██████████████████████| 720/720 [02:17<00:00, 5.24it/s]|████████
Loaded validation data with 128 batches
100%|██████████████████████| 128/128 [00:07<00:00, 16.56it/s]|████████
validation acc: 0.4535169699950812 |dataset split validation size: 2033
epoch 0 | losses: 1876.456398844719
Loaded test data with 186 batches
100%|██████████████████████| 186/186 [00:11<00:00, 16.43it/s]|████████
test acc: 0.46805648957632817 |dataset split test size: 2974

### AI (Composer Assistant)

Let me break down the results and explain why we need to load datasets:

#### Why We Need Datasets

- **Training Data**: Machine learning models learn by example. We need training data (11,514 samples) to teach the model how to classify different scenarios.
- **Validation Data**: (2,033 samples) helps us monitor if the model is learning properly during training and avoid overfitting.
- **Test Data**: (2,974 samples) gives us an unbiased evaluation of the final model's performance.

#### Results Breakdown

- **Initial Performance** (Before Training):

```
validation acc: 0.0
test acc: 0.0
```

- The model starts with 0% accuracy because it hasn't learned anything yet
- The classification layer is randomly initialized
- **After One Epoch of Training**:

```
validation acc: 0.4535 (45.35%)
test acc: 0.4681 (46.81%)
```

- The model improved significantly from 0% to ~46%
- This is much better than random guessing (1/60 ≈ 1.67%)
- The similar performance between validation and test sets suggests the model is generalizing well
- **Training Process**:

```
losses: 1876.45
```

- This is the cumulative loss over all training batches in the epoch
- The loss value by itself isn't very meaningful, but tracking it over epochs would show if the model is improving
- **Dataset Statistics**:

```
train: 11514 samples
validation: 2033 samples
test: 2974 samples
```

- This is a typical 70/15/15 split of the data
- The large training set helps the model learn patterns
- The substantial validation and test sets provide reliable performance estimates
- **Batching Information**:

```
train: 720 batches
validation: 128 batches
test: 186 batches
```

- With batch size 16, the data is processed in smaller chunks
- This helps with memory efficiency and provides more frequent model updates

The results show that the model is learning effectively, achieving ~46% accuracy on a 60-class classification problem after just one epoch. This is significantly better than random chance (1.67%) but there's still room for improvement through:

- Training for more epochs
- Adjusting hyperparameters (learning rate, dropout, etc.)
- Using more advanced fine-tuning techniques

### User

but the results is different from our expected one, which is 1.67%

### AI (Composer Assistant)

Ah, I see the confusion! Let me explain why we got ~46% accuracy instead of the expected 1.67% even without fine-tuning:

Looking at the code more carefully, we ARE actually fine-tuning the model, even though we might not have intended to! Here's what's happening:

- In `main.py`:

```
if args.task == 'baseline':
    model = ScenarioModel(args, tokenizer, target_size=60).to(device)
    run_eval(args, model, datasets, tokenizer, split='validation')   # First evaluation
    run_eval(args, model, datasets, tokenizer, split='test')         # First evaluation
    baseline_train(args, model, datasets, tokenizer)                 # Training happens here
    run_eval(args, model, datasets, tokenizer, split='test')         # Final evaluation
```

- The first evaluations (before training) show 0% accuracy, not 1.67% as we expected. This is because:

- The classification layer is randomly initialized
- With the high dropout rate (0.9 by default), most activations are being zeroed out
- This leads to the model consistently predicting the same class (hence 0% accuracy)
- Then `baseline_train()` runs, which:

```
def baseline_train(args, model, datasets, tokenizer):
    criterion = nn.CrossEntropyLoss()
    optimizer = AdamW(model.parameters(), lr=args.learning_rate)
    # ... training loop ...
```

- After one epoch of training, we get ~46% accuracy because:
  - The BERT model is being fine-tuned
  - The classification layer is being trained
  - The model is actually learning to classify scenarios

To actually see the performance WITHOUT fine-tuning, we would need to:

- Set `requires_grad=False` for all BERT parameters
- Not run the training loop at all
- Possibly adjust the dropout rate to a more reasonable value (like 0.1)
- Initialize the classification layer better

Would you like me to show you how to modify the code to see the true "no fine-tuning" performance?

### User

What could you do to further improve the performance besides setting the epoch to 10?

### AI (Composer Assistant)

Here are several strategies to improve the model's performance beyond just increasing epochs:

- **Hyperparameter Optimization**:

```
# Recommended hyperparameter adjustments in arguments.py
parser.add_argument('--batch-size', default=32, type=int)  # Increase from 16
parser.add_argument('--learning-rate', default=2e-5, type=float)  # Decrease from 1e-4
parser.add_argument('--drop-rate', default=0.1, type=float)  # Decrease from 0.9
parser.add_argument('--max-len', default=128, type=int)  # Increase from 20
parser.add_argument('--hidden-dim', default=256, type=int)  # Increase from 10
```

- **Learning Rate Scheduling**:

```
def baseline_train(args, model, datasets, tokenizer):
    # Add warmup steps and weight decay
    num_training_steps = len(train_dataloader) * args.n_epochs
    num_warmup_steps = num_training_steps * 0.1  # 10% of training steps

    optimizer = AdamW(model.parameters(),
            lr=args.learning_rate,
            eps=args.adam_epsilon,
            weight_decay=0.01)  # Add weight decay

    scheduler = get_cosine_schedule_with_warmup(  # Change to cosine schedule
        optimizer,
        num_warmup_steps=num_warmup_steps,
        num_training_steps=num_training_steps
```

- **Model Architecture Improvements**:

```
class Classifier(nn.Module):
    def __init__(self, args, target_size):
        super().__init__()
        # Add layer normalization and dropout
        self.layer_norm = nn.LayerNorm(args.embed_dim)
        self.dropout1 = nn.Dropout(args.drop_rate)
        self.top = nn.Linear(args.embed_dim, args.hidden_dim)
        self.layer_norm2 = nn.LayerNorm(args.hidden_dim)
        self.dropout2 = nn.Dropout(args.drop_rate)
        self.relu = nn.ReLU()
        self.bottom = nn.Linear(args.hidden_dim, target_size)

    def forward(self, hidden):
        x = self.layer_norm(hidden)
        x = self.dropout1(x)
        x = self.top(x)
        x = self.layer_norm2(x)
        x = self.dropout2(x)
        x = self.relu(x)
        logit = self.bottom(x)
        return logit
```

- **Data Augmentation**:

```
def prepare_features(args, data, tokenizer, cache_path):
    # Add simple augmentation techniques
    def augment_text(text):
        # Randomly mask some words
        words = text.split()
        for i in range(len(words)):
            if random.random() < 0.1:  # 10% chance to mask
                words[i] = tokenizer.mask_token
        return ' '.join(words)

    for split, examples in data.items():
        if split == 'train':  # Only augment training data
            augmented_examples = []
            for example in examples:
                # Add augmented version
                aug_example = example.copy()
                aug_example['text'] = augment_text(example['text'])
                augmented_examples.append(aug_example)
            examples.extend(augmented_examples)
```

- **Gradient Clipping and Accumulation**:

```
def baseline_train(args, model, datasets, tokenizer):
    # Add gradient clipping and accumulation
    max_grad_norm = 1.0
    gradient_accumulation_steps = 2

    for step, batch in progress_bar(enumerate(train_dataloader)):
        loss = loss / gradient_accumulation_steps  # Normalize loss
        loss.backward()

        if (step + 1) % gradient_accumulation_steps == 0:
            torch.nn.utils.clip_grad_norm_(model.parameters(), max_grad_norm)
            optimizer.step()
            scheduler.step()
            model.zero_grad()
```

- **Use Advanced BERT Variants**:

```
def model_setup(self, args):
    # Use RoBERTa or DeBERTa instead of BERT
    if args.model == 'roberta':
        self.encoder = RobertaModel.from_pretrained('roberta-base')
    elif args.model == 'deberta':
        self.encoder = DebertaModel.from_pretrained('microsoft/deberta-base')
```

- **Implement Early Stopping**:

```
class EarlyStopping:
    def __init__(self, patience=3, min_delta=0):
        self.patience = patience
        self.min_delta = min_delta
        self.counter = 0
        self.best_loss = None
        self.early_stop = False

def baseline_train(args, model, datasets, tokenizer):
    early_stopping = EarlyStopping(patience=3)
    best_val_acc = 0

    for epoch in range(args.n_epochs):
        # ... training loop ...
        val_acc = run_eval(args, model, datasets, tokenizer, split='validation')

        if val_acc > best_val_acc + early_stopping.min_delta:
            best_val_acc = val_acc
            early_stopping.counter = 0
            # Save best model
            torch.save(model.state_dict(), 'best_model.pt')
        else:
            early_stopping.counter += 1

        if early_stopping.counter >= early_stopping.patience:
            print("Early stopping triggered")
            break
```

These improvements can be implemented incrementally, measuring the impact of each change. The most important ones to start with would be:

- Adjusting the hyperparameters (especially learning rate and dropout)
- Increasing max_len to capture more context
- Adding gradient clipping and proper learning rate scheduling
- Implementing early stopping to prevent overfitting

**User**

summarize the previous prompt into a concise but detailed paragraph

**AI (Composer Assistant)**

To improve the model's performance beyond increasing epochs to 10, several key optimizations can be implemented: (1) Hyperparameter tuning - adjusting batch size to 32, reducing learning rate to 2e-5, lowering dropout rate to 0.1, increasing max sequence length to 128, and expanding hidden dimension to 256; (2) Enhanced training dynamics - implementing cosine learning rate scheduling with warmup steps, adding weight decay of 0.01, and utilizing gradient clipping at 1.0; (3) Architectural improvements - incorporating layer normalization and additional dropout layers in the classifier; (4) Training techniques - implementing data augmentation through random word masking, using gradient accumulation over 2 steps, and adding early stopping with a patience of 3 epochs; and (5) Model variations - considering advanced BERT variants like RoBERTa or DeBERTa as the base encoder. These modifications collectively address model capacity, training stability, and generalization capability, potentially leading to significant performance improvements over the baseline implementation.

**User**

modify this, and make it into a concise but detailed paragraph:
Q1. If we do not fine-tune the model, what is your expected test accuracy? Explain Why.

Based on the code and setup, if we don't fine-tune the BERT model (i.e., only run evaluation without training), we would expect a very low test accuracy, likely around 1.67% (1/60). Here's why:

- **Random Initialization of Classification Layer**:
  - The model has a classification head (`Classifier` class) that maps BERT's output to 60 classes (target_size=60)
  - This classification layer is randomly initialized and hasn't learned any meaningful mappings without training
  - The final linear layer will essentially produce random logits
- **Task Specificity**:

- The task appears to be Amazon Massive Scenario classification with 60 different scenarios/categories
- While BERT has strong general language understanding from pre-training, it hasn't been specifically adapted to classify these scenarios
- The pre-trained BERT knowledge alone cannot directly solve this specific classification task without fine-tuning
- **Architecture**:
```python
class ScenarioModel(nn.Module):
    def __init__(self, args, tokenizer, target_size):
        # ...
        self.dropout = nn.Dropout(args.drop_rate)
        self.classify = Classifier(args, target_size)
```
- The model uses a high dropout rate (0.9 by default in arguments.py)
- The classification layer is completely random before training
- Without fine-tuning, it's essentially making random predictions across 60 classes

Therefore, we would expect the model to perform at roughly random chance level (1/60 ≈ 1.67%) before any fine-tuning occurs. Any accuracy significantly above this would be purely coincidental due to the random initialization of the classification layer.

**AI (Composer Assistant)**

Without fine-tuning, the BERT model would likely achieve a test accuracy of approximately 1.67% (1/60) on the Amazon Massive Scenario classification task. This low performance can be attributed to three key factors: First, while the pre-trained BERT model possesses strong general language understanding, its classification head is randomly initialized and maps to 60 different scenario categories without any learned meaningful associations. Second, despite BERT's pre-trained knowledge, it hasn't been specifically adapted to distinguish between the unique scenarios in this task, making its general language understanding insufficient for accurate classification. Finally, the model's architecture, which includes a high dropout rate (0.9) and a randomly initialized classification layer, essentially results in random predictions across the 60 classes, leading to performance at the random chance level (1/60). Any accuracy significantly above 1.67% would be purely coincidental due to the random nature of the classification layer's initialization.

**User**

Document in as much detail as possible any other experimental settings in a paragraph, including parameters of the library functions (except for defaults), the hyperparameters you chose etc.

**AI (Composer Assistant)**

The experiment utilizes the BERT-base-uncased model from Hugging Face's transformers library (version 4.29.2) for the Amazon Massive Scenario classification task. The model processes text with a maximum sequence length of 128 tokens, using left-side truncation in the tokenizer (truncation_side="left"). The architecture includes a BERT encoder followed by a two-layer classifier with dimensions 768→256→60 where 768 is BERT's hidden size, 256 is the intermediate hidden dimension, and 60 is the number of target scenario classes. Training employs the AdamW optimizer with a learning rate of 2e-5 and epsilon of 1e-8, coupled with a linear learning rate scheduler with no warmup steps. The dropout rate is set to 0.1 for regularization, and the model is trained with a batch size of 16 for 3 epochs. The dataset is split into 11,514 training samples, 2,033 validation samples, and 2,974 test samples, with features cached in 'assets/cache/amazon.pkl' for efficiency. Cross-entropy loss is used as the training objective, and evaluation metrics are computed using accuracy. The model is trained on CUDA if available, falling back to CPU otherwise. All random seeds are set to 42 for reproducibility, and the training progress is monitored using tqdm progress bars. The model's parameters are stored in the 'results' directory, with the entire implementation running on PyTorch's framework.

**User**

@baseline_lora.txt Documents percentage of trainable parameters, performance, running time (include GPU type, which is RTX 3070Ti). And discuss the results in a detailed paragraph.

**AI (Composer Assistant)**

Based on the baseline_lora.txt results, here's a detailed analysis of the LoRA experiments run on an RTX 3070Ti GPU:

The experiments tested three LoRA configurations (ranks 4, 8, and 16) for fine-tuning BERT on the Amazon Massive Scenario classification task. For rank=4, the model had 669,696 trainable parameters (0.61% of total 110M parameters), achieving a final test accuracy of 90.08% after 10 epochs, with validation accuracy peaking at 90.80%. With rank=8, trainable parameters doubled to 1,339,392 (1.21% of total), yielding a slightly better test accuracy of 90.52%. Rank=16 further increased trainable parameters to 2,678,784 (2.39% of total) but showed slightly lower test performance at 89.48%. Training times per epoch remained consistent across configurations at approximately 45-48 seconds, with validation taking 3-4 seconds per evaluation. The results demonstrate LoRA's efficiency, as even with just 0.61% trainable parameters (rank=4), the model achieved strong performance. Interestingly, increasing LoRA rank showed diminishing returns, with rank=8 performing best overall. The training curves show rapid initial improvement (reaching ~74% validation accuracy after first epoch for rank=4) followed by gradual refinement, suggesting effective learning despite the parameter efficiency. All configurations showed stable training with consistent convergence patterns, indicating LoRA's reliability for this task.

**User**

describe the results here. this is the baseline model with 10 epochs only. no lora. put into a concise but detailed paragraph

**AI (Composer Assistant)**

The baseline BERT model, trained for 10 epochs on the Amazon Massive Scenario classification task, demonstrates strong performance progression. Starting with 0% accuracy during initial evaluation, the model rapidly improves to 66.01% validation accuracy after the first epoch, with a corresponding loss of 1777.42. Performance continues to improve significantly through subsequent epochs, reaching 87.36% by epoch 1 (loss: 833.89) and 88.88% by epoch 2 (loss: 480.74). The learning curve shows steady but diminishing improvements in later epochs, with validation accuracy reaching 89.08% by epoch 3, 89.38% by epoch 4, and ultimately peaking at 91.15% by epoch 7. The final test accuracy of 90.45% closely aligns with the validation performance, suggesting good generalization. Training efficiency is consistent throughout, processing approximately 5.23 batches per second, with validation batches running at about 16.5 batches per second. The model's loss steadily decreases from 1777.42 to 68.01 over the 10 epochs, indicating stable optimization dynamics and effective learning of the classification task.

**User**

assuming the results (like loss) are the same, integrate training accuracy into the paragraph from the previous prompt

**AI (Composer Assistant)**

The baseline BERT model, trained for 10 epochs on the Amazon Massive Scenario classification task, demonstrates strong performance progression. Starting with 0% accuracy during initial evaluation, the model rapidly improves to 66.01% validation accuracy and 36.90% training accuracy after the first epoch (loss: 1777.42). Performance continues to improve significantly through subsequent epochs, reaching 87.36% validation and 76.06% training accuracy by epoch 2 (loss: 833.89), and 88.88% validation with 89.18% training accuracy by epoch 3 (loss: 480.74). The learning curve shows steady but diminishing improvements in later epochs, with validation/training accuracies reaching 89.08%/93.06% by epoch 4, 89.38%/94.75% by epoch 5, and validation accuracy peaking at 91.15% by epoch 8 while training accuracy continues to improve to 99.00% by epoch 10. The final test accuracy of 90.45% closely aligns with the validation performance, suggesting good generalization despite the increasing gap between training and validation accuracy in later epochs. Training efficiency is consistent throughout, processing approximately 5.23 batches per second, with validation batches running at about 16.5 batches per second. The model's loss steadily decreases from 1777.42 to 68.01 over the 10 epochs, indicating stable optimization dynamics and effective learning of the classification task.