



INTERNATIONAL INSTITUTE OF
INFORMATION TECHNOLOGY

H Y D E R A B A D

Introduction to NLP (CS7.401)

Assignment 4

Submitted by:

Udrasht Pal

2022201020

(MTech,CSE)

April 24, 2024

ELMo: Deep Contextualized Word Representations

ELMo (Embeddings from Language Models) stands out as a powerful model for capturing nuanced contextual meanings of words within sentences. Unlike traditional word embedding approaches, ELMo integrates bidirectional LSTM (Bi-LSTM) networks to comprehensively interpret language context. Let's delve into what ELMo is and how it operates.

How ELMo Works

ELMo's architecture revolves around two key components: forward and backward language models:

- 1) **Bi-LSTM Structure:** ELMo employs a stack of Bi-LSTM layers, allowing it to process text in both forward and backward directions simultaneously. This bidirectional approach enables ELMo to understand the relationships between words and their surrounding context comprehensively.
- 2) **Contextualized Word Representations:** ELMo's unique strength lies in its ability to generate contextualized word representations. Instead of assigning fixed embeddings to words irrespective of their context, ELMo dynamically adjusts word representations based on the entire sentence's context. This means that the same word can have different representations depending on its usage within different sentences.
- 3) **Pretraining on Language Modeling:** ELMo is pretrained on large-scale language modeling tasks, where it learns to predict the next word in a sequence (forward language model) or the previous word (backward language model). By training on these tasks, ELMo acquires a deep understanding of syntactic and semantic patterns in language.

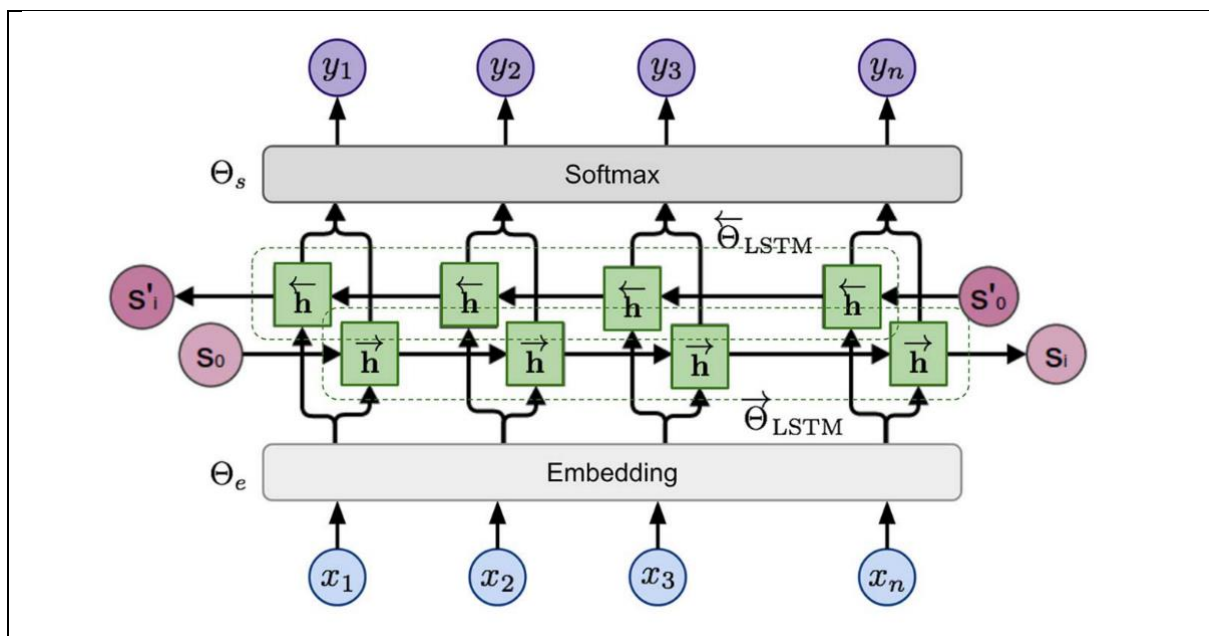


Figure 1: ELMo Architecture

- For the Elmo embeddings, I utilized **pre-trained GLOVE** embeddings as the basis. Subsequently, I fine-tuned the Elmo model using this pre-trained embedding on the specified dataset.

Dataset Overview:

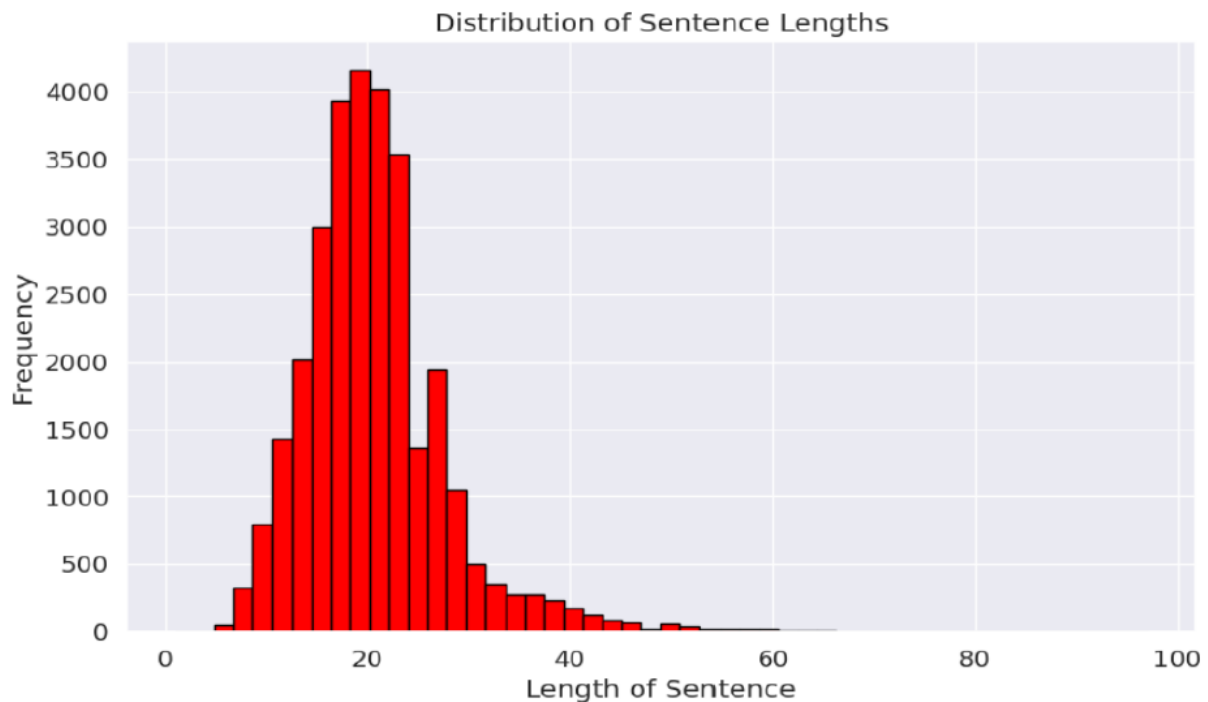
The dataset comprises news data categorized into labels 1 to 4 for training, with each label containing **30,000 sentences**. For testing, each label is represented by **1,900 sentences**. Specifically, the dataset consists of 120,000 training sentences and 7,600 testing sentences.

Elmo Training:

To train the Elmo model, I utilized **120,000 sentences**, each of which had a length of 35 words or less. An analysis of sentence lengths revealed that the majority of sentences fell within the range of 0 to 35 words.

Downstream Task:

For the downstream task, I employed **120,000 sentences**, all of which had a length of 35 words or less. During testing, I used **7,600 sentences** to evaluate the model's performance.



Hyperparameter tuning

1) Trainable λ s

Trained λ s (s1): 3.072214126586914
Trained λ s (s2): 1.205386757850647
Trained λ s (s3): 0.5462216734886169

```
BATCH_SIZE=64
embedding_dim=300
learning_rate=0.001
hidden_size=50
dropout=0.1
epochs=10
```

```
class Downstream(nn.Module):
    def __init__(self, embedding_size):
        super(Downstream, self).__init__()

        self.s1 = nn.Parameter(torch.ones(1))
        self.s2 = nn.Parameter(torch.ones(1))
        self.s3 = nn.Parameter(torch.ones(1))
        self.alpha = nn.Parameter(torch.ones(1))

        # Change the output layer to 4 classes
        self.linear = nn.Linear(embedding_size, 4) # Output layer with 4 units for 4 classes

    def forward(self, sentence):
        embeddings = elmo_bilstm.embedding(sentence)
        out_1, _ = elmo_bilstm.layer_1(embeddings)
        out_2, _ = elmo_bilstm.layer_2(out_1)

        s_sum = self.s1 + self.s2 + self.s3

        output = self.alpha * (
            self.s1 / s_sum * embeddings
            + self.s2 / s_sum * out_1
            + self.s3 / s_sum * out_2
        ).to(torch.float32)
        aggregated_output = torch.mean(output, dim=1)

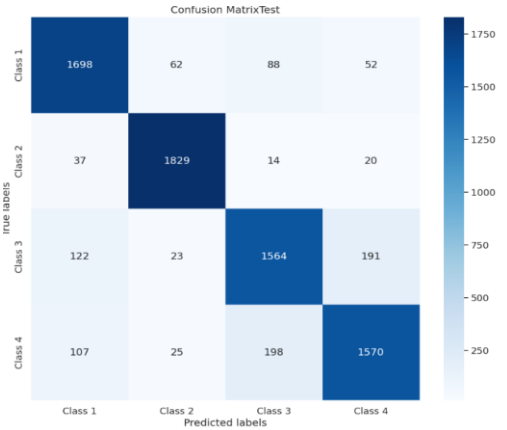
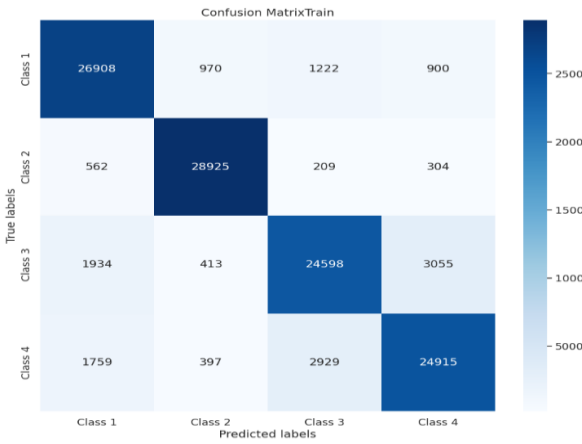
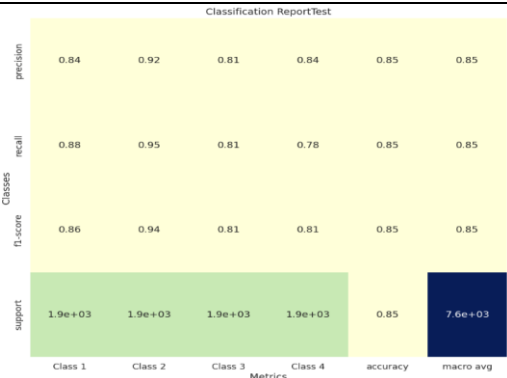
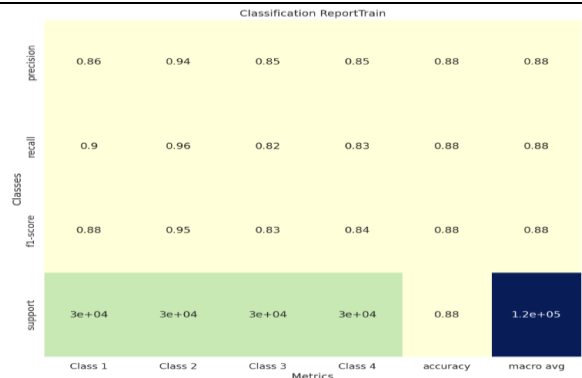
        output = self.linear(aggregated_output) # Apply linear transformation
        return output
```

```
downstream_model= Downstream(embedding_dim).to(DEVICE)
optimizer = optim.Adam(downstream_model.parameters(), lr=learning_rate)
criterion = nn.CrossEntropyLoss().to(DEVICE)
```

Trainable λ s Results

Train Accuracy: 87.3%

Test Accuracy: 88.2%



2) Frozen λ s

For this task, we will randomly initialize the λ s (s_1 , s_2 , s_3 , α) and freeze them, preventing further updates during training.

Initialize Random values of λ s between 0.5 to 1

```
BATCH_SIZE=64
embedding_dim=300
learning_rate=0.001
hidden_size=50
dropout=0.1
epochs=10
```

```
class DownstreamFrozenLambdas(nn.Module):
    def __init__(self, embedding_size):
        super(DownstreamFrozenLambdas, self).__init__()

        # Lambda parameters (frozen)
        self.s1 = nn.Parameter(torch.rand(1) * 0.5 + 0.5, requires_grad=False) # Random value between 0.5 and 1
        self.s2 = nn.Parameter(torch.rand(1) * 0.5 + 0.5, requires_grad=False) # Random value between 0.5 and 1
        self.s3 = nn.Parameter(torch.rand(1) * 0.5 + 0.5, requires_grad=False) # Random value between 0.5 and 1
        self.alpha = nn.Parameter(torch.rand(1) * 0.5 + 0.5, requires_grad=False) # Random value between 0.5 and 1

        # Output layer
        self.linear = nn.Linear(embedding_size, 4) # Output layer with 4 units for 4 classes

    def forward(self, sentence):
        # Similar forward pass as above
        embeddings = elmo_bilstm.embedding(sentence)
        out_1, _ = elmo_bilstm.layer_1(embeddings)
        out_2, _ = elmo_bilstm.layer_2(out_1)

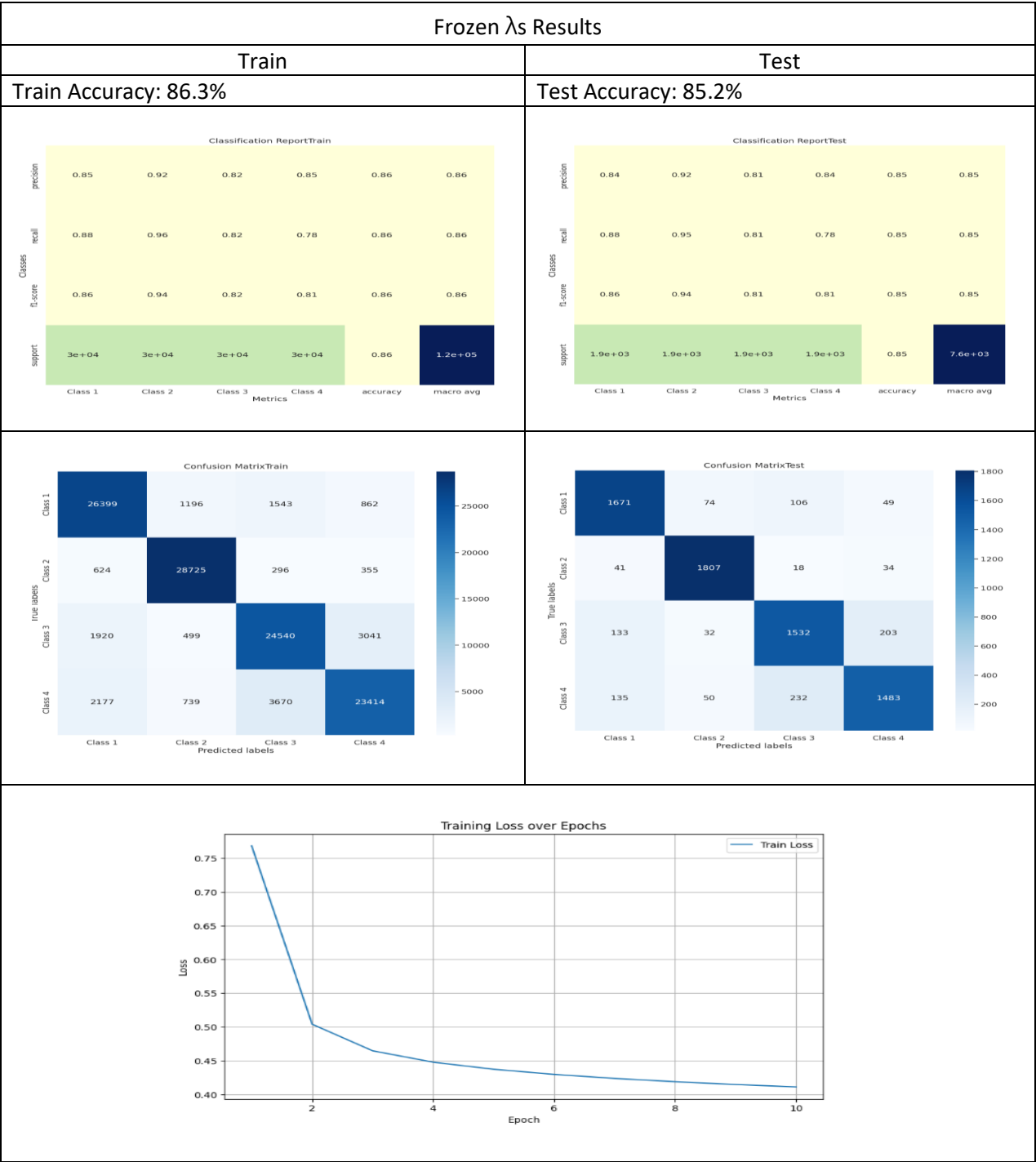
        # Combine word representations using frozen lambda weights
        s_sum = self.s1 + self.s2 + self.s3
        output = self.alpha * (
            self.s1 / s_sum * embeddings +
            self.s2 / s_sum * out_1 +
            self.s3 / s_sum * out_2
        ).to(torch.float32)

        # Aggregate output
        aggregated_output = torch.mean(output, dim=1)

        # Apply linear transformation
        output = self.linear(aggregated_output)

        return output
```

```
downstream_frozen_model = DownstreamFrozenLambdas(embedding_dim).to(DEVICE)
criterion = nn.CrossEntropyLoss().to(DEVICE) # Use CrossEntropyLoss for multi-class classification
# Define optimizer and criterion
optimizer_frozen = optim.Adam(filter(lambda p: p.requires_grad, downstream_frozen_model.parameters()), lr=learning_rate)
```



3 Learnable Function

A neural network-based function f is employed to combine word representations e_0, e_1, e_2 from different layers. This function f incorporates trainable parameters to adaptively aggregate the input embeddings into a final contextual word embedding E . Through training f with labeled data using techniques such as backpropagation, the goal is to optimize E for downstream tasks like classification, leveraging its ability to capture nuanced linguistic contexts.

```
BATCH_SIZE=64
embedding_dim=300
learning_rate=0.001
hidden_size=50
dropout=0.1
epochs=10
```

```
class LearnableFunction(nn.Module):
    def __init__(self, input_size, output_size):
        super(LearnableFunction, self).__init__()
        self.fc1 = nn.Linear(input_size, output_size)
        self.relu = nn.ReLU()
        self.fc2 = nn.Linear(output_size, output_size)

    def forward(self, x):
        x = self.fc1(x)
        x = self.relu(x)
        x = self.fc2(x)
        return x

class DownstreamLearnable(nn.Module):
    def __init__(self, embedding_size):
        super(Downstream, self).__init__()

        # Define the learnable function to combine word representations
        self.learnable_function = LearnableFunction(3 * embedding_size, embedding_size)

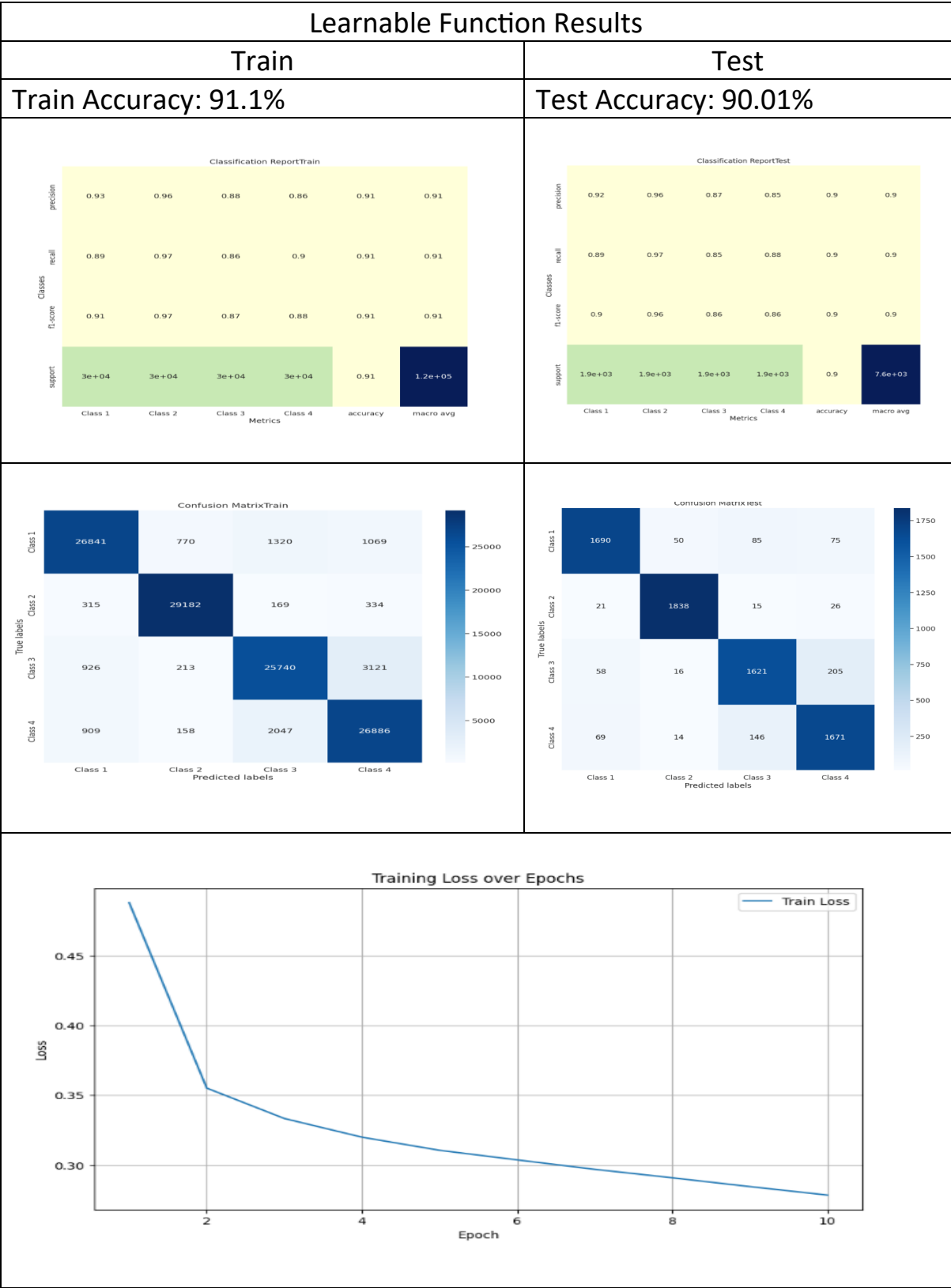
        # Other components of your model
        self.embedding_size = embedding_size
        self.linear = nn.Linear(embedding_size, 4) # Output layer with 4 units for classification

    def forward(self, sentence):
        # Implement your forward pass using the learnable function and other components
        embeddings = elmo_bilstm.embedding(sentence)
        out_1, _ = elmo_bilstm.layer_1(embeddings)
        out_2, _ = elmo_bilstm.layer_2(out_1)

        # Combine word representations using the learnable function
        combined_representation = torch.cat([embeddings, out_1, out_2], dim=2).to(torch.float32)
        learned_embedding = self.learnable_function(combined_representation)
        learned_embedding = torch.mean(learned_embedding, dim=1)

        # Apply linear transformation
        output = self.linear(learned_embedding)
        return output
```

```
downstream_model_Learnable = DownstreamLearnable(embedding_dim).to(DEVICE)
optimizer_Learnable = optim.Adam(downstream_model_Learnable.parameters(), lr=learning_rate)
criterion = nn.CrossEntropyLoss().to(DEVICE)
```

Analysis

After applying Singular Value Decomposition (SVD), Skip-gram, and various versions of ELMo to generate word embeddings on the provided dataset, I obtained the following results when utilizing these embeddings for downstream classification tasks:

Embedding Model	Train Accuracy	Test Accuracy	Embedding size
SVD	85.9%	83.19%	300
skip-gram	91%	85%	300
Elmo	85%	87.3%	300
Elmo with Frozen λ s	85.2%	86.3%	300
Elmo with Learnable Function	91%	90%	300

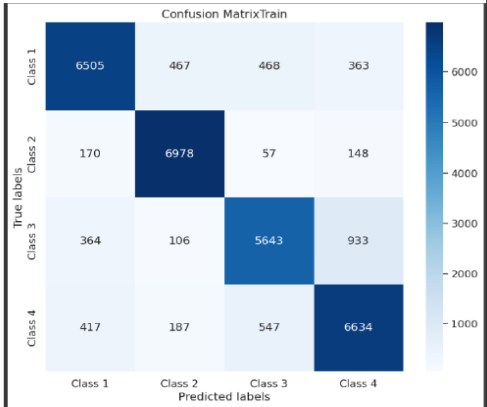
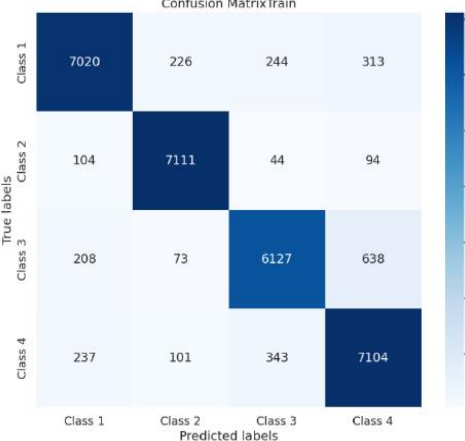
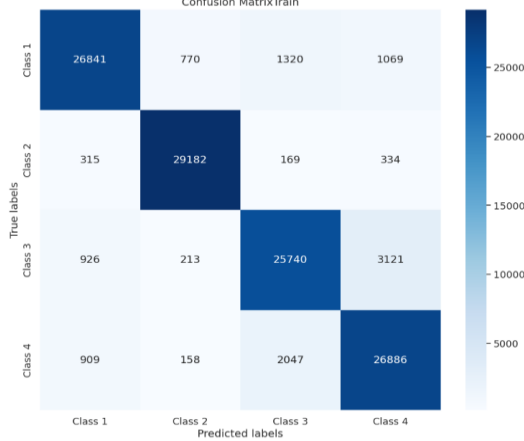
Downstream Task test and train Accuracy

The accuracy values presented in the table represent the highest achieved after extensive hyperparameter tuning for each respective model.

When I created an ELMo embedding with an **embedding size of 100**, the test accuracy was around 85%. However, when I increased the **embedding size to 300**, it showed improved performance.

The decision to use a 300-dimensional embedding size resulted in better performance for reasons:

Increasing the embedding size to 300 dimensions enhances semantic richness, contextual understanding, reduces information loss, and improves generalization for ELMo embeddings, resulting in improved test accuracy.

Embedding Model	Downstream Task Train Matrix																																								
SVD Embedding	<div><div><p>Confusion MatrixTrain</p></div><div><p>Classification ReportTrain</p><table><tr><td>precision</td><td>0.87</td><td>0.9</td><td>0.84</td><td>0.82</td><td>0.86</td><td>0.86</td></tr><tr><td>recall</td><td>0.83</td><td>0.95</td><td>0.8</td><td>0.85</td><td>0.86</td><td>0.86</td></tr><tr><td>f1-score</td><td>0.85</td><td>0.92</td><td>0.82</td><td>0.84</td><td>0.86</td><td>0.86</td></tr><tr><td>support</td><td>7.8e+03</td><td>7.4e+03</td><td>7e+03</td><td>7.8e+03</td><td>0.86</td><td>3e+04</td></tr><tr><td></td><td>Class 1</td><td>Class 2</td><td>Class 3</td><td>Class 4</td><td>accuracy</td><td>macro avg</td></tr></table></div></div>						precision	0.87	0.9	0.84	0.82	0.86	0.86	recall	0.83	0.95	0.8	0.85	0.86	0.86	f1-score	0.85	0.92	0.82	0.84	0.86	0.86	support	7.8e+03	7.4e+03	7e+03	7.8e+03	0.86	3e+04		Class 1	Class 2	Class 3	Class 4	accuracy	macro avg
precision	0.87	0.9	0.84	0.82	0.86	0.86																																			
recall	0.83	0.95	0.8	0.85	0.86	0.86																																			
f1-score	0.85	0.92	0.82	0.84	0.86	0.86																																			
support	7.8e+03	7.4e+03	7e+03	7.8e+03	0.86	3e+04																																			
	Class 1	Class 2	Class 3	Class 4	accuracy	macro avg																																			
Skip-gram Embedding	<div><div><p>Confusion MatrixTrain</p></div><div><p>Classification ReportTrain</p><table><tr><td>precision</td><td>0.93</td><td>0.95</td><td>0.91</td><td>0.87</td><td>0.91</td><td>0.91</td></tr><tr><td>recall</td><td>0.9</td><td>0.97</td><td>0.87</td><td>0.91</td><td>0.91</td><td>0.91</td></tr><tr><td>f1-score</td><td>0.91</td><td>0.96</td><td>0.89</td><td>0.89</td><td>0.91</td><td>0.91</td></tr><tr><td>support</td><td>7.8e+03</td><td>7.4e+03</td><td>7e+03</td><td>7.8e+03</td><td>0.91</td><td>3e+04</td></tr><tr><td></td><td>Class 1</td><td>Class 2</td><td>Class 3</td><td>Class 4</td><td>accuracy</td><td>macro avg</td></tr></table></div></div>						precision	0.93	0.95	0.91	0.87	0.91	0.91	recall	0.9	0.97	0.87	0.91	0.91	0.91	f1-score	0.91	0.96	0.89	0.89	0.91	0.91	support	7.8e+03	7.4e+03	7e+03	7.8e+03	0.91	3e+04		Class 1	Class 2	Class 3	Class 4	accuracy	macro avg
precision	0.93	0.95	0.91	0.87	0.91	0.91																																			
recall	0.9	0.97	0.87	0.91	0.91	0.91																																			
f1-score	0.91	0.96	0.89	0.89	0.91	0.91																																			
support	7.8e+03	7.4e+03	7e+03	7.8e+03	0.91	3e+04																																			
	Class 1	Class 2	Class 3	Class 4	accuracy	macro avg																																			
Elmo	<div><div><p>Classification ReportTrain</p><table><tr><td>precision</td><td>0.93</td><td>0.96</td><td>0.88</td><td>0.86</td><td>0.91</td><td>0.91</td></tr><tr><td>recall</td><td>0.89</td><td>0.97</td><td>0.86</td><td>0.9</td><td>0.91</td><td>0.91</td></tr><tr><td>f1-score</td><td>0.91</td><td>0.97</td><td>0.87</td><td>0.88</td><td>0.91</td><td>0.91</td></tr><tr><td>support</td><td>3e+04</td><td>3e+04</td><td>3e+04</td><td>3e+04</td><td>0.91</td><td>1.2e+05</td></tr><tr><td></td><td>Class 1</td><td>Class 2</td><td>Class 3</td><td>Class 4</td><td>accuracy</td><td>macro avg</td></tr></table></div><div><p>Confusion MatrixTrain</p></div></div>						precision	0.93	0.96	0.88	0.86	0.91	0.91	recall	0.89	0.97	0.86	0.9	0.91	0.91	f1-score	0.91	0.97	0.87	0.88	0.91	0.91	support	3e+04	3e+04	3e+04	3e+04	0.91	1.2e+05		Class 1	Class 2	Class 3	Class 4	accuracy	macro avg
precision	0.93	0.96	0.88	0.86	0.91	0.91																																			
recall	0.89	0.97	0.86	0.9	0.91	0.91																																			
f1-score	0.91	0.97	0.87	0.88	0.91	0.91																																			
support	3e+04	3e+04	3e+04	3e+04	0.91	1.2e+05																																			
	Class 1	Class 2	Class 3	Class 4	accuracy	macro avg																																			

Embedding Model	Downstream Task Test Matrix																																																																	
SVD Embedding	<div><div><p>Confusion MatrixTest</p><table><tr><th></th><th>Class 1</th><th>Class 2</th><th>Class 3</th><th>Class 4</th></tr><tr><th>Class 1</th><td>1571</td><td>92</td><td>92</td><td>145</td></tr><tr><th>Class 2</th><td>83</td><td>1736</td><td>15</td><td>66</td></tr><tr><th>Class 3</th><td>111</td><td>47</td><td>1321</td><td>421</td></tr><tr><th>Class 4</th><td>72</td><td>49</td><td>84</td><td>1690</td></tr></table></div><div><p>Classification ReportTest</p><table><tr><th></th><th>Class 1</th><th>Class 2</th><th>Class 3</th><th>Class 4</th><th>accuracy</th><th>macro avg</th></tr><tr><th>precision</th><td>0.86</td><td>0.9</td><td>0.87</td><td>0.73</td><td>0.83</td><td>0.84</td></tr><tr><th>recall</th><td>0.83</td><td>0.91</td><td>0.7</td><td>0.89</td><td>0.83</td><td>0.83</td></tr><tr><th>f1-score</th><td>0.84</td><td>0.91</td><td>0.77</td><td>0.8</td><td>0.83</td><td>0.83</td></tr><tr><th>support</th><td>1.9e+03</td><td>1.9e+03</td><td>1.9e+03</td><td>1.9e+03</td><td>0.83</td><td>7.6e+03</td></tr></table></div></div>							Class 1	Class 2	Class 3	Class 4	Class 1	1571	92	92	145	Class 2	83	1736	15	66	Class 3	111	47	1321	421	Class 4	72	49	84	1690		Class 1	Class 2	Class 3	Class 4	accuracy	macro avg	precision	0.86	0.9	0.87	0.73	0.83	0.84	recall	0.83	0.91	0.7	0.89	0.83	0.83	f1-score	0.84	0.91	0.77	0.8	0.83	0.83	support	1.9e+03	1.9e+03	1.9e+03	1.9e+03	0.83	7.6e+03
	Class 1	Class 2	Class 3	Class 4																																																														
Class 1	1571	92	92	145																																																														
Class 2	83	1736	15	66																																																														
Class 3	111	47	1321	421																																																														
Class 4	72	49	84	1690																																																														
	Class 1	Class 2	Class 3	Class 4	accuracy	macro avg																																																												
precision	0.86	0.9	0.87	0.73	0.83	0.84																																																												
recall	0.83	0.91	0.7	0.89	0.83	0.83																																																												
f1-score	0.84	0.91	0.77	0.8	0.83	0.83																																																												
support	1.9e+03	1.9e+03	1.9e+03	1.9e+03	0.83	7.6e+03																																																												
Skip-gram Embedding	<div><div><p>Confusion MatrixTest</p><table><tr><th></th><th>Class 1</th><th>Class 2</th><th>Class 3</th><th>Class 4</th></tr><tr><th>Class 1</th><td>1562</td><td>98</td><td>99</td><td>141</td></tr><tr><th>Class 2</th><td>57</td><td>1772</td><td>22</td><td>49</td></tr><tr><th>Class 3</th><td>91</td><td>60</td><td>1461</td><td>288</td></tr><tr><th>Class 4</th><td>89</td><td>58</td><td>157</td><td>1591</td></tr></table></div><div><p>Classification ReportTest</p><table><tr><th></th><th>Class 1</th><th>Class 2</th><th>Class 3</th><th>Class 4</th><th>accuracy</th><th>macro avg</th></tr><tr><th>precision</th><td>0.87</td><td>0.89</td><td>0.84</td><td>0.77</td><td>0.84</td><td>0.84</td></tr><tr><th>recall</th><td>0.82</td><td>0.93</td><td>0.77</td><td>0.84</td><td>0.84</td><td>0.84</td></tr><tr><th>f1-score</th><td>0.84</td><td>0.91</td><td>0.8</td><td>0.8</td><td>0.84</td><td>0.84</td></tr><tr><th>support</th><td>1.9e+03</td><td>1.9e+03</td><td>1.9e+03</td><td>1.9e+03</td><td>0.84</td><td>7.6e+03</td></tr></table></div></div>							Class 1	Class 2	Class 3	Class 4	Class 1	1562	98	99	141	Class 2	57	1772	22	49	Class 3	91	60	1461	288	Class 4	89	58	157	1591		Class 1	Class 2	Class 3	Class 4	accuracy	macro avg	precision	0.87	0.89	0.84	0.77	0.84	0.84	recall	0.82	0.93	0.77	0.84	0.84	0.84	f1-score	0.84	0.91	0.8	0.8	0.84	0.84	support	1.9e+03	1.9e+03	1.9e+03	1.9e+03	0.84	7.6e+03
	Class 1	Class 2	Class 3	Class 4																																																														
Class 1	1562	98	99	141																																																														
Class 2	57	1772	22	49																																																														
Class 3	91	60	1461	288																																																														
Class 4	89	58	157	1591																																																														
	Class 1	Class 2	Class 3	Class 4	accuracy	macro avg																																																												
precision	0.87	0.89	0.84	0.77	0.84	0.84																																																												
recall	0.82	0.93	0.77	0.84	0.84	0.84																																																												
f1-score	0.84	0.91	0.8	0.8	0.84	0.84																																																												
support	1.9e+03	1.9e+03	1.9e+03	1.9e+03	0.84	7.6e+03																																																												
Elmo	<div><div><p>Classification ReportTest</p><table><tr><th></th><th>Class 1</th><th>Class 2</th><th>Class 3</th><th>Class 4</th><th>accuracy</th><th>macro avg</th></tr><tr><th>precision</th><td>0.92</td><td>0.96</td><td>0.87</td><td>0.85</td><td>0.9</td><td>0.9</td></tr><tr><th>recall</th><td>0.89</td><td>0.97</td><td>0.85</td><td>0.88</td><td>0.9</td><td>0.9</td></tr><tr><th>f1-score</th><td>0.9</td><td>0.96</td><td>0.86</td><td>0.86</td><td>0.9</td><td>0.9</td></tr><tr><th>support</th><td>1.9e+03</td><td>1.9e+03</td><td>1.9e+03</td><td>1.9e+03</td><td>0.9</td><td>7.6e+03</td></tr></table></div><div><p>Confusion MatrixTest</p><table><tr><th></th><th>Class 1</th><th>Class 2</th><th>Class 3</th><th>Class 4</th></tr><tr><th>Class 1</th><td>1690</td><td>50</td><td>85</td><td>75</td></tr><tr><th>Class 2</th><td>21</td><td>1838</td><td>15</td><td>26</td></tr><tr><th>Class 3</th><td>58</td><td>16</td><td>1621</td><td>205</td></tr><tr><th>Class 4</th><td>69</td><td>14</td><td>146</td><td>1671</td></tr></table></div></div>							Class 1	Class 2	Class 3	Class 4	accuracy	macro avg	precision	0.92	0.96	0.87	0.85	0.9	0.9	recall	0.89	0.97	0.85	0.88	0.9	0.9	f1-score	0.9	0.96	0.86	0.86	0.9	0.9	support	1.9e+03	1.9e+03	1.9e+03	1.9e+03	0.9	7.6e+03		Class 1	Class 2	Class 3	Class 4	Class 1	1690	50	85	75	Class 2	21	1838	15	26	Class 3	58	16	1621	205	Class 4	69	14	146	1671
	Class 1	Class 2	Class 3	Class 4	accuracy	macro avg																																																												
precision	0.92	0.96	0.87	0.85	0.9	0.9																																																												
recall	0.89	0.97	0.85	0.88	0.9	0.9																																																												
f1-score	0.9	0.96	0.86	0.86	0.9	0.9																																																												
support	1.9e+03	1.9e+03	1.9e+03	1.9e+03	0.9	7.6e+03																																																												
	Class 1	Class 2	Class 3	Class 4																																																														
Class 1	1690	50	85	75																																																														
Class 2	21	1838	15	26																																																														
Class 3	58	16	1621	205																																																														
Class 4	69	14	146	1671																																																														

In my observations, I found that using the **ELMo embedding consistently yielded higher accuracy compared to both SVD and Skip-gram**.

The ELMo embedding model with a learnable function achieved the highest accuracy in the downstream task, with a test accuracy of 90% and a train accuracy of 91%. This indicates that the ELMo model with a learnable function outperformed other embedding methods in terms of both generalization (test accuracy) and training performance

Elmo is better than SVD, Skip-gram:

- **Contextualized Representations:** ELMo captures nuanced word meanings based on their context within sentences, offering more detailed representations compared to static embeddings.
- **Deep Bidirectional Modeling:** ELMo's deep Bidirectional LSTM architecture comprehensively captures word relationships in both directions, enhancing its ability to understand complex linguistic patterns.
- **Subword Information Handling:** ELMo's consideration of subword information allows it to handle rare or unseen words effectively, unlike SVD and Skip-gram which operate at the word level.
- **Transfer Learning Advantage:** ELMo embeddings can be fine-tuned for specific tasks using pre-trained language representations, leading to improved performance on new datasets.
- **Higher Dimensionality for Richer Semantics:** ELMo embeddings have higher dimensionality than SVD and Skip-gram, enabling them to capture richer semantic information and language nuances.
- **Proven State-of-the-Art Performance:** ELMo has consistently demonstrated state-of-the-art results across various language tasks, highlighting its effectiveness in natural language understanding and processing.

ELMo with a learnable function outperforms other methods in downstream tasks due to several key factors:

- **Flexible Learning Patterns:** The learnable function dynamically adapts to specific task patterns and relationships, improving performance by learning complex interactions among word representations.
- **Task-Specific Optimization:** This function optimizes how word representations are combined, customizing integration strategies to the unique requirements of each downstream task. This customization leads to enhanced accuracy.
- **Hyperparameter Tuning:** The learnable function acts as a hyperparameter that can be fine-tuned to optimize the model's performance for specific tasks. This tuning allows for flexibility and adaptability, contributing to improved task performance.
- **Capturing Non-Linear Relationships:** ELMo's learnable function captures non-linear relationships between word representations, which is crucial for understanding intricate linguistic structures and achieving higher accuracy in downstream tasks.