

Transfer Learning of BERT: Using DNABERT for the Classification of Enhancer Sequences

Sawyer Lehman, Nicholas Dibley

May 7, 2023

1 Introduction

Enhancer regions are DNA sequences that can activate or enhance the transcription of nearby genes. These regions can be located upstream, downstream, or within the introns of the gene they regulate, and they can be far away from the transcription start site. Enhancer regions contain specific DNA sequences that can bind to transcription factors and other regulatory proteins, which in turn can influence the activity of the nearby genes. These regions are often identified by profiling specific epigenetic marks and using chromatin immunoprecipitation coupled with Next-generation sequencing (NGS), specifically ChIP-Seq [1]. While ChIP-Seq can be used to identify enhancer regions, there are several downsides to using this technique for this purpose. ChIP-Seq has a limited resolution, meaning that it may not be able to precisely identify the boundaries of enhancer regions. Enhancers can be several kilobases long, and ChIP-Seq may only be able to identify the location of the transcription factor binding site within that region. ChIP-Seq can also produce false positives, where it identifies regions that are not actually enhancers but are bound by the transcription factor due to nonspecific interactions or experimental artifacts. In addition, enhancer regions are often tissue-specific, meaning that they are active only in certain cell types. ChIP-Seq experiments can be biased towards specific cell types, leading to the identification of enhancers that are not relevant in other cell types. Lastly, functional validation is often required to confirm that these regions actually act as enhancers in vivo [2]. This can be a challenging and time-consuming process.

Unlike in some other regions of the genome, enhancer sequences do not contain well-defined consensus motifs. Instead, enhancers can contain a diverse range of motifs, which makes it challenging to identify the most relevant ones based on frequency alone. Additionally, as mentioned before, the activity of enhancers is context-specific and can vary across cell types and developmental stages as well. Many computational tools used to identify motifs in enhancer sequences generate a large number of false positives, such as correlation based algorithms in [3]. These false positives can make it difficult to distinguish between biologically relevant motifs and noise.

1.1 Bioinformatics Research Problem

Due to these limitations in technology and challenges, the next-best alternative is to attempt to make predictions about whether or not a sequence is an enhancer. Certain pre-trained deep learning language models developed for Natural Language Processing (NLP) are capable of zero-shot learning and with fine-tuning, these models can perform tasks that they were not specifically trained for [4]. Furthermore, it has been shown that these language models are able to complete tasks with DNA sequences, such as non-coding regions [5] [6]. The identification and prediction of non-coding regulatory DNA sequences, such as enhancers, that are located many base pairs away from the target gene with a deep learning model would be beneficial. This leads to the question: *Can fine-tuned pre-trained deep learning language models identify enhancers?* We aim to examine the transfer learning abilities of language models intended for non-coding DNA sequences with the goal of classifying enhancer sequences.

2 Machine Learning Methods

The language model that we aim to use for the task of classification of enhancer sequences is the DNABERT model [6]. This model is a variation of BERT [7] that has been pre-trained on non-coding regions of DNA sequences. This model, like many other BERT models, is originally trained for the task of Masked Language Modeling (MLM), which requires the model to predict a masked word, or token, based on the context of all tokens in the sentence surrounding the masked token in both directions. The process for MLM is first, the sequence is tokenized and input into the BERT model, which has many layers of encoders. One of the token inputs in the sequence is masked, and the model predicts the highest probable token in the masked position based on the predefined vocabulary available.

This model can be used for tasks other than MLM with proper fine-tuning. Fine-tuning a pre-trained model allows the original weights to influence the re-training, and the model is slightly updated so that it can perform the desired task. More specifically, the original pre-trained head of the model is discarded, but the body of the model is kept. The weights of the head are then randomly initialized to the specific task that you intend to fine-tune the model for. Finally, the weights of this new head are trained for the specific task with the training data to transfer the knowledge of the pre-trained model to the task.

3 Data

The original source of the dataset we used is curated from the Ensembl database. The data is from [8], which is a repository of genomic data benchmarks such as promoters, enhancers, and open chromatin regions from three organisms: humans, mice, and fruit flies. For the purposes of this report, we are interested in human enhancer regions. The enhancer sequence data is labeled,

and the negative samples are randomly sampled sequences from the organism that are non-overlapping with the enhancer sequences. In total, there are a little over 130,000 sequences. About 100,000 were used for training, with balanced classes: 50,000 sequences of positive samples and 50,000 negative samples. The remaining 30,000 was used for testing, which also has balanced classes: 15,000 positive samples and 15,000 negative samples. The sequences were “tokenized” into 6-mers and encoded so that they could be input into the model.

Some preliminary cleaning was done to the data, shown in Listing 2 in 8. This was mainly removing any sequences in the training data that were also in the testing data to ensure unbiased evaluations of the models explained in 4.1 and 5. Additionally, a smaller subset, $n=1000$, of the training data with 500 positive and 500 negative samples was made to do a small hyperparameter search explained in 4.1.

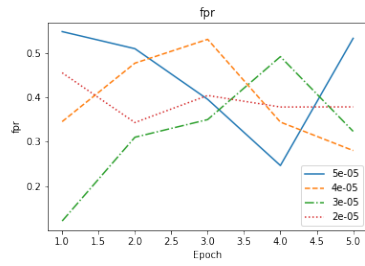
4 Experimental Design

Our overall goal in the project was to test the transferability of DNABERT to a different task. To do this, we first tested the model before fine-tuning the classifier to serve as a baseline performance. Then, a small hyperparameter search was done to see what the best strategy would be to fine-tune the model. Finally, the fine-tuned model was tested and compared to the performance of the original model.

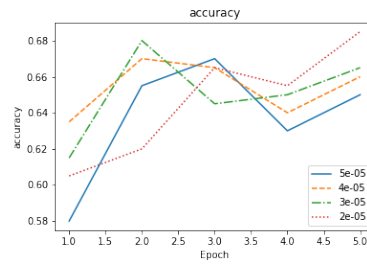
4.1 Hyperparameter search

The first step in fine-tuning the DNABERT model is to choose certain hyperparameters, code shown in Listing 4 in Section 8. To do this, a subset was taken from the training dataset, explained in Section 3. From this subset, 80 percent was used for training and the other 20 percent was used for validation. The learning rates that were considered were $5e-5$, $4e-5$, $3e-5$, and $2e-5$, which are the same Adam optimizer learning rates found in [7]. Each round of fine-tuning with the different learning rates was done for 5 epochs and the performances on the validation set for each epoch were recorded. Another parameter that affects the model training is the batch size, but due to limitations, we could only use a batch size of 8.

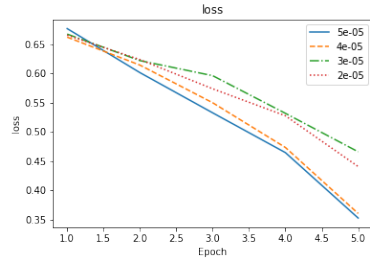
Based on this small hyperparameter search, we can see that the model trained with a learning rate of $3e-5$ at 2 epochs had the best accuracy (Figure 1b) and false positive rate (Figure 1a). Although this model with these hyperparameters did not perform the best in terms of the training loss (Figure 1c) and F1 score (Figure 1d), we know that oftentimes the loss on a validation set is somewhat inflated. In addition, it is sometimes beneficial to have a lower F1 score during validation since if this were too high, it could indicate that the model is overfitting.



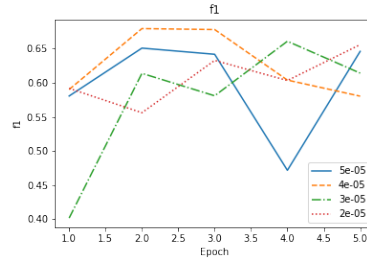
(a) False Positive Rate per Epoch



(b) Accuracy per Epoch



(c) Training Loss per Epoch



(d) F1 Score per Epoch

Figure 1: Validation metrics for DNABERT model fine-tuned at learning rates of $5e-5$ (blue solid line), $4e-5$ (orange dashed), $3e-5$ (green dot-dashed), and $2e-5$ (red dotted) over 5 epochs.

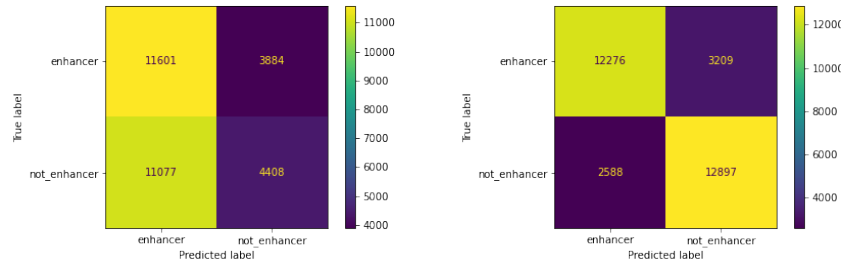
5 Experimental Results

The model was first evaluated before fine-tuning with the testing set, to see the baseline performance. The performances of both models are shown in Table 1. The original model with no updated weights to the classifier was neither accurate nor precise, with an accuracy of 51.7% and a precision of 51.1%. The recall, which is also known as the true positive rate or sensitivity, is relatively high at 74.9%. Upon examining the confusion matrix, the model is predicting a lot of false positives, which contributes to the sensitivity of the model (Figure 2a).

After fine-tuning the model with the full training dataset, with a learning rate of $3e-5$ at 2 epochs, the same testing set was used to evaluate the performance. The accuracy and precision were improved greatly to 81.3% and 82.6% respectively. The recall and specificity also improved to 79.3% and 83.3% respectively. Examining the confusion matrix revealed that the model is able to predict true negatives more than true positives, which contributes to the specificity of this model (Figure 2b).

metric	original	fine-tuned
Accuracy	0.517	0.813
Precision	0.511	0.826
Recall	0.749	0.793
Specificity	0.285	0.833

Table 1: Evaluation results of DNABERT model before (original) and after (fine-tuned) fine-tuning.



(a) Confusion matrix of predictions **before** fine-tuning. (b) Confusion matrix of predictions **after** fine-tuning.

Figure 2: Confusion matrices of predictions.

6 Conclusion

Based on the results in Section 5, we see that the fine-tuned version of the DNABERT model is capable of transferring the pre-trained knowledge to the task well. We can tell that the fine-tuning made a difference because the performance on the testing dataset before the model head was fine-tuned was essentially random guessing, which is to be expected since the weights of the classification head added to the model were initialized with randomized weights as mentioned in Section 2. The model before fine-tuning was much more sensitive, which follows the character of producing a lot of false positives much like the other models for enhancer classification. After fine-tuning, the model was able to predict more true negatives, making it a more specific model.

7 Discussion

Overall, the benefit of using a pre-trained model such as BERT is the ability to transfer across tasks or even domains. This is particularly useful if the dataset is too small or noisy for a fully-fledged deep-learning model on its own. This characteristic is also beneficial because it mitigates the need for feature engineering and domain-specific knowledge that is required when creating an entirely new deep-learning model. A challenge that could arise with fine-tuning is if the data is too different from the original pre-training data, which would lead to a poorly underfitting model. Another challenge when fine-tuning, or training a deep learning model in general, is that it is computationally expensive. While attempting to fine-tune the model, we suffered greatly from the lack of GPU resources available to us, hence why we were limited to a smaller batch size. Lastly, the overall lack of interpretability of the decisions made by the model creates a problem when trying to understand what features are important in your data.

References

- [1] S. Blinka, M. H. Reimer, K. Pulakanti, L. Pinello, G.-C. Yuan, and S. Rao, “Identification of transcribed enhancers by genome-wide chromatin immunoprecipitation sequencing,” *Enhancer RNAs: Methods and Protocols*, pp. 91–109, 2017.
- [2] A. Visel, S. Minovitsky, I. Dubchak, and L. A. Pennacchio, “Vista enhancer browser—a database of tissue-specific human enhancers,” *Nucleic acids research*, vol. 35, no. suppl.1, pp. D88–D92, 2007.
- [3] J. M. Hariprakash and F. Ferrari, “Computational biology solutions to identify enhancers-target gene pairs,” *Computational and structural biotechnology journal*, vol. 17, pp. 821–831, 2019.
- [4] J. Wei, M. Bosma, V. Y. Zhao, K. Guu, A. W. Yu, B. Lester, N. Du, A. M. Dai, and Q. V. Le, “Finetuned language models are zero-shot learners,” *arXiv preprint arXiv:2109.01652*, 2021.
- [5] G. Benegas, S. S. Batra, and Y. S. Song, “Dna language models are powerful zero-shot predictors of non-coding variant effects,” *bioRxiv*, pp. 2022–08, 2022.
- [6] Y. Ji, Z. Zhou, H. Liu, and R. V. Davuluri, “Dnabert: pre-trained bidirectional encoder representations from transformers model for dna-language in genome,” *Bioinformatics*, vol. 37, no. 15, pp. 2112–2120, 2021.
- [7] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” *arXiv preprint arXiv:1810.04805*, 2018.
- [8] K. Gresova, V. Martinek, D. Cechak, P. Simecek, and P. Alexiou, “Genomic benchmarks: A collection of datasets for genomic sequence classification,” *bioRxiv*, pp. 2022–06, 2022.

8 Appendix

```
1 # Imports:
2 import pandas as pd
3 import numpy as np
4 import random
5 from tqdm import trange
6 from tabulate import tabulate
7 import matplotlib.pyplot as plt
8 from sklearn.model_selection import train_test_split
9 from sklearn import metrics
10 import torch
11 import torch.nn as nn
12 from torch.utils.data import DataLoader
13 from datasets import load_dataset
14 from transformers import AutoTokenizer,
15     AutoModelForSequenceClassification, AdamW
16 from torch.utils.data import TensorDataset, DataLoader,
17     RandomSampler, SequentialSampler
18
19 # Functions:
20 def data_to_df(data, subset):
21     """huggingface raw data to df,
22     subset can be 'train' or 'test'
23     """
24     lst_data = []
25     for i in range(data[subset].num_rows):
26         seq = data[subset][i]['seq']
27         label = data[subset][i]['label']
28         lst_data.append([seq, label])
29
30     df = pd.DataFrame(lst_data, columns=['seq', 'label'])
31     return df
32
33 def preprocessing(input_text, tokenizer):
34     """
35     Returns <class transformers.tokenization_utils_base.BatchEncoding
36     > with the following fields:
37     - input_ids: list of token ids
38     - token_type_ids: list of token type ids
39     - attention_mask: list of indices (0,1) specifying which tokens
40       should considered by the model (return_attention_mask = True).
41     """
42     return tokenizer.encode_plus(
43         input_text,
44         add_special_tokens = True,
45         max_length = 128,
46         padding='max_length',
47         return_attention_mask = True,
48         truncation=True,
49         return_tensors = 'pt'
50     )
51
52 def print_rand_sentence_encoding(text):
53     """
54     Displays tokens, token IDs and attention mask of a random text
55     sample
56     """
```



```

51     '''
52     index = random.randint(0, len(text) - 1)
53     tokens = tokenizer.tokenize(tokenizer.decode(token_id[index]))
54     token_ids = [i.numpy() for i in token_id[index]]
55     attention = [i.numpy() for i in attention_masks[index]]
56
57     table = np.array([tokens, token_ids, attention]).T
58     print(tabulate(table,
59                   headers = ['Tokens', 'Token IDs', 'Attention
Mask'],
60                           tablefmt = 'fancy_grid')
61           )
62
63
64 def b_tp(preds, labels):
65     '''
66     Returns True Positives (TP): count of correct predictions of
actual class 1
67     '''
68     return sum([preds == labels and preds == 1 for preds, labels in
zip(preds, labels)])
69
70 def b_fp(preds, labels):
71     '''
72     Returns False Positives (FP): count of wrong predictions of
actual class 1
73     '''
74     return sum([preds != labels and preds == 1 for preds, labels in
zip(preds, labels)])
75
76 def b_tn(preds, labels):
77     '''
78     Returns True Negatives (TN): count of correct predictions of
actual class 0
79     '''
80     return sum([preds == labels and preds == 0 for preds, labels in
zip(preds, labels)])
81
82 def b_fn(preds, labels):
83     '''
84     Returns False Negatives (FN): count of wrong predictions of
actual class 0
85     '''
86     return sum([preds != labels and preds == 0 for preds, labels in
zip(preds, labels)])
87
88 def b_metrics(preds, labels):
89     '''
90     Returns the following metrics:
91     - accuracy      = (TP + TN) / N
92     - precision     = TP / (TP + FP)
93     - recall        = TP / (TP + FN)
94     - false positive rate = FP / (FP + TN)
95     - f1 = (2 * TP) / ((2 * TP) + FP + FN)
96     '''
97     preds = np.argmax(preds, axis = 1).flatten()
98     labels = labels.flatten()

```

```

99     tp = b_tp(preds, labels)
100     tn = b_tn(preds, labels)
101     fp = b_fp(preds, labels)
102     fn = b_fn(preds, labels)
103
104     b_accuracy = (tp + tn) / (tp + tn + fp + fn)
105     b_precision = tp / (tp + fp) if (tp + fp) > 0 else 'nan'
106     b_recall = tp / (tp + fn) if (tp + fn) > 0 else 'nan'
107     #b_specificity = tn / (tn + fp) if (tn + fp) > 0 else 'nan'
108     b_fpr = fp / (fp + tn) if (fp + tn) > 0 else 'nan'
109     b_f1 = (2*tp) / ((2*tp) + fp + fn) if ((2*tp) + fp + fn) > 0
        else 'nan'
110
111     return b_accuracy, b_precision, b_recall, b_fpr, b_f1
112
113 def train_dna_bert(model, train_dataloader, validation_dataloader,
        epochs):
114     # Recommended number of epochs: 2, 3, 4. See: https://arxiv.org/pdf/1810.04805.pdf
115     epochs = epochs
116
117     performance = []
118
119     for _ in trange(epochs, desc = 'Epoch'):
120
121         # ===== Training =====
122
123         # Set model to training mode
124         model.train()
125
126         # Tracking variables
127         tr_loss = 0
128         nb_tr_examples, nb_tr_steps = 0, 0
129
130         for step, batch in enumerate(train_dataloader):
131             batch = tuple(t.to(device) for t in batch)
132             b_input_ids, b_input_mask, b_labels = batch
133             optimizer.zero_grad()
134             # Forward pass
135             train_output = model(b_input_ids,
136                                 token_type_ids = None,
137                                 attention_mask = b_input_mask,
138                                 labels = b_labels)
139
140             # Backward pass
141             train_output.loss.backward()
142             optimizer.step()
143             # Update tracking variables
144             tr_loss += train_output.loss.item()
145             nb_tr_examples += b_input_ids.size(0)
146             nb_tr_steps += 1
147
148         # ===== Validation =====
149
150         model.eval()
151
152         val_accuracy = []
153         val_precision = []

```

```

153     val_recall = []
154     val_fpr = []
155     val_f1 = []
156
157     for batch in validation_dataloader:
158         batch = tuple(t.to(device) for t in batch)
159         b_input_ids, b_input_mask, b_labels = batch
160         with torch.no_grad():
161             # Forward pass
162             eval_output = model(b_input_ids,
163                                 token_type_ids = None,
164                                 attention_mask = b_input_mask
165                                 )
166
167             logits = eval_output.logits.detach().cpu().numpy()
168             label_ids = b_labels.to('cpu').numpy()
169
170             b_accuracy, b_precision, b_recall, b_fpr, b_f1 =
171             b_metrics(logits, label_ids)
172             val_accuracy.append(b_accuracy)
173             if b_precision != 'nan': val_precision.append(
174             b_precision)
175             if b_recall != 'nan': val_recall.append(b_recall)
176             if b_fpr != 'nan': val_fpr.append(b_fpr)
177             if b_f1 != 'nan': val_f1.append(b_f1)
178
179             print('\n\t - Train loss: {:.4f}'.format(tr_loss /
180             nb_tr_steps))
181             print('\t - Validation Accuracy: {:.4f}'.format(sum(
182             val_accuracy)/len(val_accuracy)))
183             print('\t - Validation Precision: {:.4f}'.format(sum(
184             val_precision)/len(val_precision)) if len(val_precision)>0 else
185             '\t - Validation Precision: NaN')
186             print('\t - Validation Recall: {:.4f}'.format(sum(
187             val_recall)/len(val_recall)) if len(val_recall)>0 else '\t -
188             Validation Recall: NaN')
189             print('\t - Validation FPR: {:.4f}'.format(sum(val_fpr)/len(
190             val_fpr)) if len(val_fpr)>0 else '\t - Validation FPR: NaN')
191             print('\t - Validation F1: {:.4f}\n'.format(sum(val_f1)/len(
192             val_f1))) if len(val_f1)>0 else '\t - Validation F1: NaN'
193
194             performance.append({'loss':tr_loss / nb_tr_steps,
195                                 'accuracy':sum(val_accuracy)/len(
196             val_accuracy),
197                                 'precision':sum(val_precision)/len(
198             val_precision) if len(val_precision)>0 else 0,
199                                 'recall':sum(val_recall)/len(val_recall
200             ) if len(val_recall)>0 else 0,
201                                 'fpr':sum(val_fpr)/len(val_fpr) if len(
202             val_fpr)>0 else 0,
203                                 'f1':sum(val_f1)/len(val_f1) if len(
204             val_f1)>0 else 0
205                                 }
206             )
207
208     return performance

```

```

195 def make_performance_list(x):
196     loss = []
197     acc = []
198     prec = []
199     recall = []
200     fpr = []
201     f1 = []
202
203     for epoch in x:
204         loss.append(epoch['loss'])
205         acc.append(epoch['accuracy'])
206         prec.append(epoch['precision'])
207         recall.append(epoch['recall'])
208         fpr.append(epoch['fpr'])
209         f1.append(epoch['f1'])
210
211     return loss, acc, prec, recall, fpr, f1
212
213 def make_plot(title, x1, x2, x3, x4, y):
214     plt.plot(y, x1, label = '5e-05', linestyle='--')
215     plt.plot(y, x2, label = '4e-05', linestyle='--')
216     plt.plot(y, x3, label = '3e-05', linestyle='-.')
217     plt.plot(y, x4, label = '2e-05', linestyle=':')
218     plt.legend()
219     plt.title(title)
220     plt.xlabel('Epoch')
221     plt.ylabel(title)
222     file_name = './visualizations/param_search_'+title+'.png'
223     plt.savefig(file_name)
224     plt.show()
225
226 def test_metrics(df):
227     pos = df[df['true_label']=='enhancer']
228     tp = len(pos[pos['prediction']=='enhancer'])
229     fn = len(pos[pos['prediction']=='not_enhancer'])
230     neg = df[df['true_label']=='not_enhancer']
231     fp = len(neg[neg['prediction']=='enhancer'])
232     tn = len(neg[neg['prediction']=='not_enhancer'])
233
234     print('accuracy: ', ((tp+tn)/(tp+tn+fp+fn)))
235     print('precision: ', (tp/(tp+fp)))
236     print('recall: ', (tp/(tp+fn)))
237     print('specificity: ', (tn/(tn+fp)))
238     print('f1 score: ', ((2*tp)/((2*tp)+fp+fn)))

```

Listing 1: Imports and functions

```

1 # Preparing data:
2 # loading data from huggingface
3 raw_data = load_dataset("katarinagresova/
    Genomic_Benchmarks_human_enhancers_ensembl")
4 # converting huggingface raw data to df
5 train = data_to_df(raw_data, 'train')
6 test = data_to_df(raw_data, 'test')
7 # getting overlapping samples in the train and test
8 match_lst = []
9 for x in train['seq'].values.tolist():
10     for y in test['seq'].values.tolist():

```

```

11         if x==y:
12             idx = train.index[train['seq'] == x].tolist()
13             match_lst.extend(idx)
14 # removing the matching samples from the training set
15 train = train.drop(match_lst)
16 # creating a training subset with 500 positive and 500 negative
   samples
17 pos_train = train[train['label'] == 1].sample(n=500, random_state
   =42) # getting random sample of 500 positive samples
18 neg_train = train[train['label'] == 0].sample(n=500, random_state
   =42) # getting random sample of 500 negative samples
19 train_subset = pos_train.append(neg_train, ignore_index=True).
   sample(frac=1, random_state=42) # combining and shuffling
20 # writing data to csv files
21 train.to_csv('./data/train.csv', index=False)
22 test.to_csv('./data/test.csv', index=False)
23 train_subset.to_csv('./data/train_subset.csv', index=False)

```

Listing 2: Preparing Data.

```

1 # setting device to use gpu
2 device = torch.device('cuda' if torch.cuda.is_available() else 'cpu'
   ')
3 # loading DNABERT model and putting model to gpu
4 id2label = {0: "not_enhancer", 1: "enhancer"}
5 label2id = {"not_enhancer": 0, "enhancer": 1}
6 tokenizer = AutoTokenizer.from_pretrained("zhihan1996/DNA_bert_6",
   trust_remote_code=True)
7 model = AutoModelForSequenceClassification.from_pretrained("
   zhihan1996/DNA_bert_6", num_labels=2,
8
   id2label
   =id2label, label2id=label2id,
9
   trust_remote_code=True
10
   )
11 model.cuda()
12 # loading test dataset
13 test = pd.read_csv('./data/test.csv')
14 pos = test[test['label']==1]
15 neg = test[test['label']==0]
16 neg = neg.sample(n=len(pos), random_state=42)
17 test = pd.concat([pos, neg], ignore_index=True)
18 test_text = test.seq.values.tolist()
19 # formatting sequences to 6-mers
20 test_text = [' '.join([seq[i:i+6] for i in range(0, len(seq), 6)])
   for seq in test_text]
21 test_labels = test.label.values.tolist()
22 test_output = []
23 for seq, lab in zip(test_text, test_labels):
24     test_ids = []
25     test_attention_mask = []
26     encoding = preprocessing(seq, tokenizer)
27
28     test_ids.append(encoding['input_ids'])
29     test_attention_mask.append(encoding['attention_mask'])
30     test_ids = torch.cat(test_ids, dim = 0)
31     test_attention_mask = torch.cat(test_attention_mask, dim = 0)
32

```

```

33     with torch.no_grad():
34         output = model(test_ids.to(device), token_type_ids = None,
35                        attention_mask = test_attention_mask.to(device))
36
37         prediction = 'enhancer' if np.argmax(output.logits.cpu().numpy
38        ()).flatten().item() == 1 else 'not_enhancer'
39
40         if lab == 1:
41             true_lab = 'enhancer'
42         else:
43             true_lab = 'not_enhancer'
44
45         test_output.append([seq, true_lab, prediction])
46
47 results = pd.DataFrame(test_output, columns = ['sequence', '
48         true_label', 'prediction'])
49 results.to_csv('zero_shot_predictions.csv', index=False)

```

Listing 3: Testing model before finetuning.

```

1 device = torch.device('cuda' if torch.cuda.is_available() else 'cpu
2 ')
3 tokenizer = AutoTokenizer.from_pretrained("zhihan1996/DNA_bert_6",
4     trust_remote_code=True)
5 # Loading subset for hyperparameter search
6 train = pd.read_csv('./data/train_subset.csv')
7 train_text = train.seq.values.tolist()
8 # formatting sequences to 6-mers:
9 train_text = [' '.join([seq[i:i+6] for i in range(0, len(seq), 6)])
10               for seq in train_text]
11 train_labels = train.label.values.tolist()
12 # checking to make sure tokenization worked
13 print(rand_sentence_encoding(train_text))
14 # tokenizing train subset
15 token_id = []
16 attention_masks = []
17
18 for sample in train_text:
19     encoding_dict = preprocessing(sample, tokenizer)
20     token_id.append(encoding_dict['input_ids'])
21     attention_masks.append(encoding_dict['attention_mask'])
22
23 token_id = torch.cat(token_id, dim = 0)
24 attention_masks = torch.cat(attention_masks, dim = 0)
25 labels = torch.tensor(train_labels)
26
27 val_ratio = 0.2
28 batch_size = 8
29
30 # splitting the training into train and validation
31 train_idx, val_idx = train_test_split(
32     np.arange(len(labels)),
33     test_size = val_ratio,
34     shuffle = True,
35     stratify = labels)
36
37 train_set = TensorDataset(token_id[train_idx],

```

```

36         attention_masks[train_idx],
37         labels[train_idx])
38
39 val_set = TensorDataset(token_id[val_idx],
40                         attention_masks[val_idx],
41                         labels[val_idx])
42
43 # wrapping in dataloader object
44 train_dataloader = DataLoader(
45     train_set,
46     sampler = RandomSampler(train_set),
47     batch_size = batch_size
48 )
49
50 validation_dataloader = DataLoader(
51     val_set,
52     sampler = SequentialSampler(val_set),
53     batch_size = batch_size
54 )
55
56 id2label = {0: "not_enhancer", 1: "enhancer"}
57 label2id = {"not_enhancer": 0, "enhancer": 1}
58
59 # Recommended learning rates (Adam): 5e-5, 3e-5, 2e-5. See: https://arxiv.org/pdf/1810.04805.pdf
60 learning_rates = [5e-05, 4e-05, 3e-05, 2e-05]
61 performance = {}
62
63 for lr in learning_rates:
64     model = AutoModelForSequenceClassification.from_pretrained("
65         zhihan1996/DNA_bert_6", num_labels=2,
66         id2label
67         =id2label, label2id=label2id,
68         trust_remote_code=True
69     )
70
71     #model = nn.DataParallel(model)
72
73     optimizer = torch.optim.AdamW(model.parameters(), lr = lr, eps
74     = 1e-08)
75
76     model.cuda()
77
78     output = train_dna_bert(model, train_dataloader,
79     validation_dataloader, epochs=5)
80     performance[str(lr)] = output
81     torch.cuda.empty_cache()
82
83 y = [1, 2, 3, 4, 5]
84 x1=performance['5e-05']
85 x2=performance['4e-05']
86 x3=performance['3e-05']
87 x4=performance['2e-05']
88
89 x1_loss, x1_acc, x1_prec, x1_recall, x1_fpr, x1_f1 =
90     make_performance_list(x1)

```

```

86 x2_loss, x2_acc, x2_prec, x2_recall, x2_fpr, x2_f1 =
    make_performance_list(x2)
87 x3_loss, x3_acc, x3_prec, x3_recall, x3_fpr, x3_f1 =
    make_performance_list(x3)
88 x4_loss, x4_acc, x4_prec, x4_recall, x4_fpr, x4_f1 =
    make_performance_list(x4)
89
90 make_plot('loss', x1_loss, x2_loss, x3_loss, x4_loss, y)
91 make_plot('accuracy', x1_acc, x2_acc, x3_acc, x4_acc, y)
92 make_plot('precision', x1_prec, x2_prec, x3_prec, x4_prec, y)
93 make_plot('recall', x1_recall, x2_recall, x3_recall, x4_recall, y)
94 make_plot('fpr', x1_fpr, x2_fpr, x3_fpr, x4_fpr, y)
95 make_plot('f1', x1_f1, x2_f1, x3_f1, x4_f1, y)

```

Listing 4: Hyperparameter search.

```

1  # Fine-tune
2  # -----
3  # further data cleaning:
4  train = pd.read_csv('./data/train.csv')
5  pos = train[train['label']==1]
6  neg = train[train['label']==0]
7  neg = neg.sample(n=len(pos), random_state=42)
8  train = pd.concat([pos, neg], ignore_index=True)
9  train_text = train.seq.values.tolist()
10 # formatting sequences to 6-mers
11 train_text = [' '.join([seq[i:i+6] for i in range(0, len(seq), 6)])
    for seq in train_text]
12 train_labels = train.label.values.tolist()
13
14 token_id = []
15 attention_masks = []
16
17 for sample in train_text:
18     encoding_dict = preprocessing(sample, tokenizer)
19     token_id.append(encoding_dict['input_ids'])
20     attention_masks.append(encoding_dict['attention_mask'])
21
22
23 token_id = torch.cat(token_id, dim = 0)
24 attention_masks = torch.cat(attention_masks, dim = 0)
25 labels = torch.tensor(train_labels)
26
27 val_ratio = 0.2
28 batch_size = 8
29
30 # splitting the training into train and validation
31 train_idx, val_idx = train_test_split(
32     np.arange(len(labels)),
33     test_size = val_ratio,
34     shuffle = True,
35     stratify = labels)
36
37 train_set = TensorDataset(token_id[train_idx],
38                           attention_masks[train_idx],
39                           labels[train_idx])
40
41 val_set = TensorDataset(token_id[val_idx],

```



```

42         attention_masks[val_idx],
43         labels[val_idx])
44
45 # wrapping in dataloader object
46 train_dataloader = DataLoader(
47     train_set,
48     sampler = RandomSampler(train_set),
49     batch_size = batch_size
50 )
51
52 validation_dataloader = DataLoader(
53     val_set,
54     sampler = SequentialSampler(val_set),
55     batch_size = batch_size
56 )
57
58 id2label = {0: "not_enhancer", 1: "enhancer"}
59 label2id = {"not_enhancer": 0, "enhancer": 1}
60 model = AutoModelForSequenceClassification.from_pretrained("
    zhihan1996/DNA_bert_6", num_labels=2,
61                                     id2label
    =id2label, label2id=label2id,
62
    trust_remote_code=True
63                                     )
64
65 optimizer = torch.optim.AdamW(model.parameters(), lr = 3e-05, eps =
    1e-08)
66 model.cuda()
67 final_output = train_dna_bert(model, train_dataloader,
    validation_dataloader, epochs=2)
68 final_df = pd.DataFrame(final_output)
69 final_df.to_csv('./fine_tune_validation_results.csv')
70
71 # Testing
72 # -----
73 # loading test dataset
74 test = pd.read_csv('./data/test.csv')
75 pos = test[test['label']==1]
76 neg = test[test['label']==0]
77 neg = neg.sample(n=len(pos), random_state=42)
78 test = pd.concat([pos, neg], ignore_index=True)
79 test_text = test.seq.values.tolist()
80 # formatting sequences to 6-mers
81 test_text = [' '.join([seq[i:i+6] for i in range(0, len(seq), 6)])
    for seq in test_text]
82 test_labels = test.label.values.tolist()
83 test_output = []
84 for seq, lab in zip(test_text, test_labels):
85     test_ids = []
86     test_attention_mask = []
87     encoding = preprocessing(seq, tokenizer)
88
89     test_ids.append(encoding['input_ids'])
90     test_attention_mask.append(encoding['attention_mask'])
91     test_ids = torch.cat(test_ids, dim = 0)
92     test_attention_mask = torch.cat(test_attention_mask, dim = 0)

```

```

93
94     with torch.no_grad():
95         output = model(test_ids.to(device), token_type_ids = None,
96                        attention_mask = test_attention_mask.to(device))
97
98         prediction = 'enhancer' if np.argmax(output.logits.cpu().numpy
99         ()).flatten().item() == 1 else 'not_enhancer'
100
101         if lab == 1:
102             true_lab = 'enhancer'
103         else:
104             true_lab = 'not_enhancer'
105
106         test_output.append([seq, true_lab, prediction])
107
108 results = pd.DataFrame(test_output, columns = ['sequence', '
109         true_label', 'prediction'])
110 results.to_csv('./finetune_model_predictions.csv', index=False)

```

Listing 5: Fine-tuning model and testing.

```

1 # Zero-shot predictions
2 # -----
3 zero = pd.read_csv('./zero_shot_predictions.csv', index_col=0)
4 test_metrics(zero)
5 actual = zero['true_label'].values
6 predicted = zero['prediction'].values
7 confusion_matrix = metrics.confusion_matrix(actual, predicted)
8 cm_display = metrics.ConfusionMatrixDisplay(confusion_matrix =
9         confusion_matrix, display_labels = ['enhancer', 'not_enhancer'
10         ])
11 cm_display.plot()
12 plt.savefig('./visualizations/zero_shot_cm.png')
13 plt.show()
14
15 # Fine-tune predictions
16 # -----
17 fine = pd.read_csv('./finetune_model_predictions.csv')
18 test_metrics(fine)
19 actual = fine['true_label'].values
20 predicted = fine['prediction'].values
21 confusion_matrix = metrics.confusion_matrix(actual, predicted)
22 cm_display = metrics.ConfusionMatrixDisplay(confusion_matrix =
23         confusion_matrix, display_labels = ['enhancer', 'not_enhancer'
24         ])
25 cm_display.plot()
26 plt.savefig('./visualizations/final_test_cm.png')
27 plt.show()

```

Listing 6: Getting final metrics from zero-shot and fine-tune models.