
Investigating the functionality of Key, Query, Value submodules in Attention Mechanism

Ziyi Xie, Letian Jiang
Center for Data Science
New York University
zx1153@nyu.edu, lj2136@nyu.edu

Abstract

1 Transformer is the state-of-the-art model while the functionality of its attention
2 submodules is still widely discussed. To gain more insights of the attention mech-
3 anism, we investigated the importance of individual submodule component in the
4 Transformer model by two methods: 1) adversarial removal of linguistic informa-
5 tion on the submodules; 2) rewinding the weights of the submodules. By observing
6 the training curve and importance scores, we concluded that V component is more
7 important and attacking or freezing submodules causes overfitting and underfitting.

8 1 Introduction

9 Attention-based models have been one of the latest inventions that improve deep learning model
10 performance in the field of Natural Language Processing. The solely attention-based Transformer
11 model was first proposed in the paper 'Attention is All You Need' [1]. The attention mechanism
12 proposed in this paper contains three major components: Key (K), Query (Q), and Value (V). As
13 shown in figure 1, those three components are used to calculate the Scaled Dot-Product Attention that
14 tells the model how much 'attention' it should pay to certain inputs.

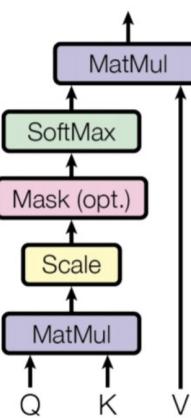


Figure 1: Scaled Dot-Product Attention [1]

15 There have been many other papers that propose improvements over this original transformer model
16 structure to generalize the use of transformers. However, there haven't been many investigations on

17 analyzing the respective roles of Key, Query, and Value components/submodules. Therefore, in this
18 paper, we will use two prevalent methods to analyze how limiting Transformer's K, Q, V access to
19 linguistic information at different layers in training can affect the overall efficiency of training and
20 performance of the model. The first method involves comparing model performance before and after
21 removing linguistic information from the submodules. The second method involves directly changing
22 submodule parameters in several different ways.

23 **2 Setup**

24 **2.1 Model**

25 **2.1.1 Transformer**

26 For the transformer model, we use the similar structure from the paper 'Attention Is All You Need'
27 [1], containing a 5 layers and 5 heads encoder and a single layer linear decoder. The task of this
28 transformer model is to assign a probability distribution for the word following a sequence of words
29 in a sentence, or simply next word prediction.

30 **2.2 Data**

31 The dataset we will be using is WikiText-2 [4]. This (silver) dataset contains tokens collected from the
32 set of verified Good and Featured articles on Wikipedia. We then use Spacy to parse every sentence
33 to produce a dataset that contains the Part-of-Speech (POS) and sentence Dependency (DEP) labels
34 for the WikiText-2 dataset. All the datasets will contain a train set and a validation set. There are
35 102,499 sentences in the train set, and 21,441 sentences in the validation set. During training, the
36 batch size is set to be 20, and the evaluation batch size is set to 10. So there are about 2.4M masked
37 word prediction pairs.

38 **2.2.1 Intermediate Layer Representation Access**

39 Since we are studying how changing K, Q, V could influence the transformer model, we need to know
40 how to extract those intermediate representations. Following the road map in figure 2, we are able to
41 extract K, Q, V representations in each transformer encoder layer and the corresponding matrices
42 W_K, W_Q, W_V that produce the representations

43 **2.3 Evaluation**

44 For evaluation, we mainly focus on looking at the loss and accuracy. The loss is calculated as
45 cross-entropy losses, and the accuracy is calculated by comparing the most likely next word predicted
46 by the model versus the gold standard word.

47 **3 Methodology**

48 **3.1 Experiment 1: Adversarial Attacking Method**

49 In this experiment, we want to investigate how removing syntactic information from intermediate
50 representation inside a transformer affects the training of the transformer. In particular, we will be
51 removing the POS or DEP information from the K, Q, V representations at different layers of the
52 transformer model. To remove such information, we use an adversarial model[3] to remove the
53 information. The adversarial model contains 5 fully connected linear layers with leaky relu and batch
54 normalization in between. The task of this adversarial model is to predict the POS or DEP of the next
55 word given intermediate representations from the transformer model.

56 To better explain how these two models work together to achieve the information removal , lets denote

- 57 • F as the transformer model

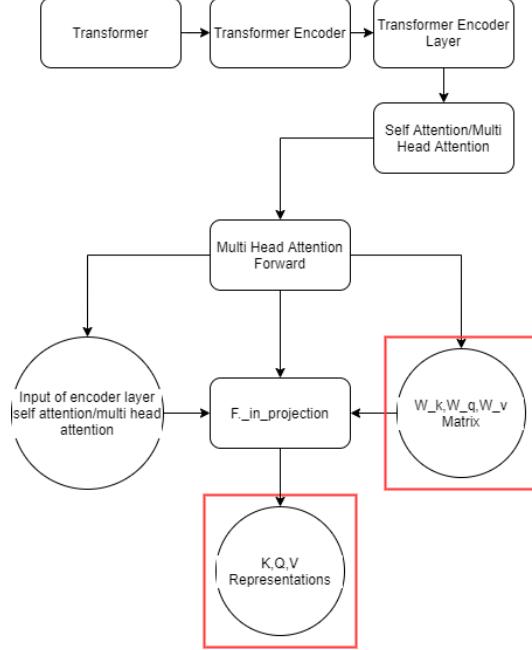


Figure 2: The method used to extract K, Q, V representations and the associated matrix operations to calculate the representations (marked in red boxes)

- 58 • A as the Adversarial model
 59 • x as the input of the transformer model (sentences)
 60 • x' as intermediate representation value in the transformer model
 61 • y as the correct label for the next word prediction task
 62 • y' as the correct label for the POS or DEP of the next word prediction task
- 63 Now we can define some cross entropy loss functions:
 64 $loss_{adversarial} = Loss(A(x', y'))$,
 65 $loss_{transformer} = Loss(F(x, y))$,
 66 $loss_{total} = loss_{transformer} - \lambda * loss_{adversarial}$, where λ is some scalar constant that represents
 67 how strong the presence of the adversarial network is.
- 68 The two models will be trained alternatively. So within each training epoch, two backpropogations
 69 are done. When backpropagation is done on $loss_{total}$, the transformer model is optimized to decrease
 70 the loss on the next word prediction task and change its intermediate representations to increase
 71 the loss of the adversarial model because of the minus sign. When backpropagation is done on
 72 $loss_{adversarial}$, the adversarial model is optimized to decrease its loss on predicting the POS or DEP
 73 of the next word.
- 74 So there is a back and forth process between the training of the transformer model and the adversarial
 75 model: the adversarial model is trying to learn the POS or DEP information of the next word from the
 76 transformer intermediate representations, while the transformer is trying to change its intermediate
 77 representations to prevent the adversarial model from effective learning. In this way, the adversarial
 78 model is able to prevent the intermediate layer representations in the transformer to contain POS or
 79 DEP information.
- 80 The choice of x' is the K, Q, V representations at layer 1 through layer 5 of the transformer model.
 81 We say that we are attacking the K representations at layer 1 of the transformer if the input of the
 82 adversarial model is the K representations at layer 1 of the transformer.

83 For each choice of intermediate representations, we train both models for 100 epochs and look at
84 the loss of the transformer model for evidence about the importance of POS or DEP information in
85 different sub-modules (K, Q, V) of a transformer model during training. We also train a baseline
86 transformer model without the adversarial model as a comparison.

87 **3.2 Experiment 2: Submodule Weights Editing**

88 In this experiment, we investigate the module/submodule importance by rewinding submodule
89 weights and freezing submodule weights. More specifically, we explore the model performance drop
90 as we modify a specific submodule (K, Q, V) by: 1) rewinding the final submodule weights back to
91 previous submodule weights (before convergence); 2) rewinding the model weights of the submodules
92 back to the convex combination of the initial weights and the final weights; 3) freezing the submodule
93 weights during the training process. In all the methods, the models weights are modified at one layer
94 and one attention component at a time. The three methods are described in the following subsections,
95 and the experiment results are in Section 5.

96 **3.2.1 Convex Combination**

97 Observing each submodule in a model demonstrate different robustness characteristics to parameter
98 perturbation, Zhang et al [5] reported the *module criticality* phenomenon. In another paper,
99 Chatterji et al. [6] formalized the *criticality* score metric and proved the metric is effective in
100 revealing the network generalization. For each submodule *att*, they proposed to rewind submodule
101 weights back to a convex combination, which is defined as $\theta_{att}^\alpha = (1 - \alpha)\theta_{att}^0 + \alpha \cdot \theta_{att}^f$, $\alpha \in [0, 1]$,
102 where θ_{att}^0 denotes the initial weights and θ_{att}^f denotes the final weights.

103 The criticality of a submodule component *att* is definded as:

$$\text{Criti}_{att} = \min \alpha_{att} \text{ s.t. } \text{Loss} \left(\text{Model with } \theta_{att}^f \right) - \text{Loss} \left(\text{Model with } \theta_{att}^{rew} \right) < \epsilon \quad (1)$$

104 where θ_{att}^f denotes the final model and θ_{att}^{rew} denotes the weights after rewinding, and ϵ is a threshold.
105 The criticality score describes the minimum α that the weights can get to initialization while the
106 performance drop is smaller than the threshold ϵ . More intuitively, a submodule can be said critical if
107 the model performance drops significantly after being modified: a high criticality score indicates that
108 the model rely heavily on the specific submodule.

109 **3.2.2 Rewinding to Mid-stage Weights**

110 The previous method inspire us to rewind the submodule weights of the final model back to any
111 model weights during the training. Suppose the model is trained to convergence after n epochs, then
112 we rewind $\theta_{att}^f = \theta_{att}^i$, $i \in [0, n]$ where θ_{att}^i is the model weight at the training epoch i .

113 Though we cannot calculate criticality score, we can observe the submodule importance during the
114 training.

115 **3.3 Freezing submodules during the training**

116 Different from previous methods that part of model weights (submodules) are modified at evaluation,
117 another method is to freeze the submodule weights during the training process by setting gradients to
118 zeros.

119 In the method, the frozen model training curves are compared with the baseline model trained
120 normally.

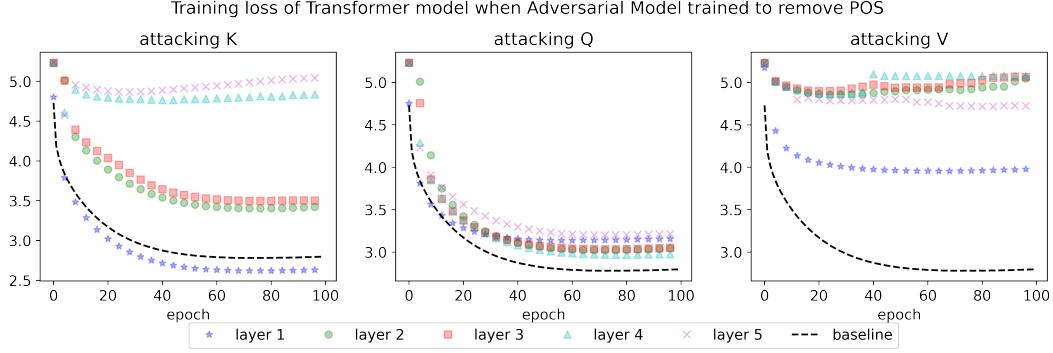


Figure 3: Training loss of the transformer model when using the adversarial model to attack POS information in K, Q, V representations at different layers

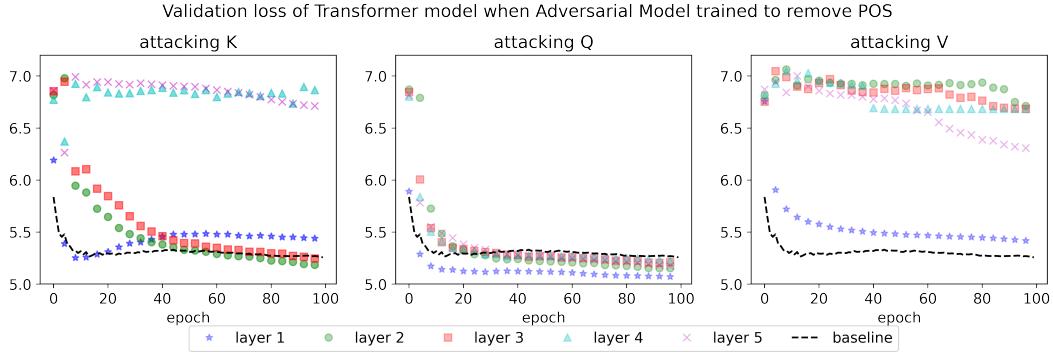


Figure 4: Validation loss of the transformer model when using the adversarial model to attack POS information in K, Q, V representations at different layers

121 4 Results and discussion of Adversarial Attacking Transformer submodules

122 4.1 Results

123 Figure 3 is the training loss of the transformer when an adversarial network is used to attack the POS
 124 information in K, Q, V representations at layer 1 through layer 5, and a baseline transformer loss
 125 when no adversarial network is used. We find that when POS in K representations in the earliest
 126 layer is being attacked, the train loss of the transformer model is improved over the baseline model.
 127 However, if we look at validation loss as in figure 4, the validation loss of such attack increases over
 128 the baseline model. Therefore attacking POS in K from the earliest layer causes the model to overfit.
 129 Attacking POS in K from middle layers increases the transformer model train loss but result in a
 130 lower than baseline validation loss, making the model more robust. Attacking POS of K in late layers
 131 makes both the train and validation loss close to initialization.

132 When POS in Q representations is attacked, the train loss of attacking all five layers cluster closely
 133 above the baseline training loss. The better than baseline validation loss further shows that such
 134 attacks make the model more robust

135 When POS in V are attacked in middle and late layers, we observe close to initialization train and
 136 validation loss. Attacking POS in V from the earliest layer leads to a moderately worse train loss but
 137 very close to baseline validation loss.

138 In figure 5 is the training loss of the transformer model when DEP information is being attacked.
 139 When attacking DEP in K, the train loss is higher at all layers. For middle to late layers, the train loss
 140 is close to initialization, but for early layers, the train loss is still reasonably close to the baseline.

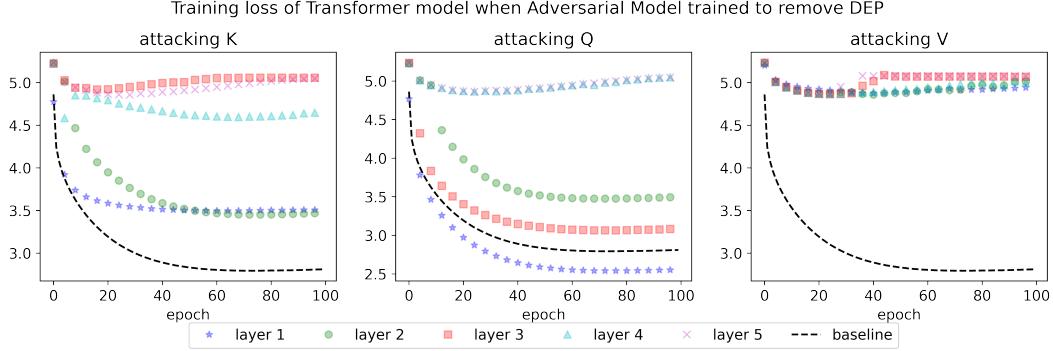


Figure 5: Training loss of the transformer model when using the adversarial model to attack DEP information in K, Q, V representations at different layers

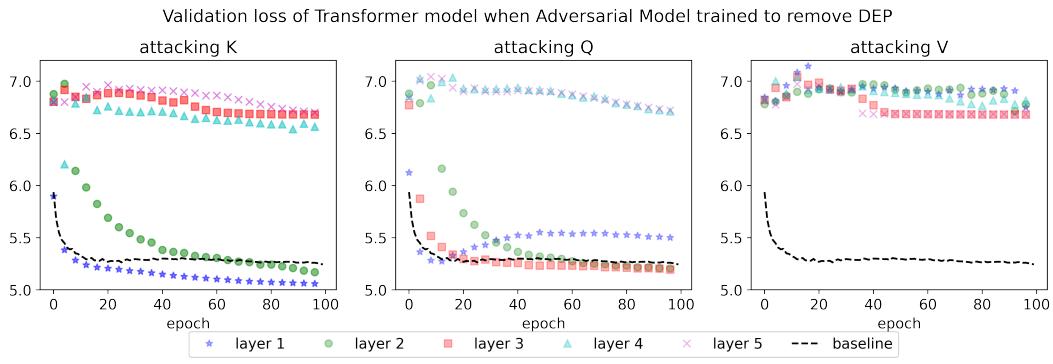


Figure 6: Validation loss of the transformer model when using the adversarial model to attack DEP information in K, Q, V representations at different layers

141 Looking at the validation loss, the attacking DEP in K of early layers produces better than baseline
 142 validation loss. When attacking DEP information in Q, we see a very similar train loss pattern
 143 compared to attacking POS in K. Attacking DEP in V at all layers makes the model not learn
 144 anything, because both the train loss and validation loss are just initialization values.

145 **4.2 Discussion**

146 We can try to explain the findings using certain interpretations of the structure of the transformer model
 147 in terms of the positioning of the layers and the structure that produce the K, Q, V representations.
 148 Since the model is taking part of sentences to predict a distribution over possible next words, we
 149 can think of earlier layers functions to transfer words to representations, and later layers function to
 150 transfer representations to probability distributions. Then, for K, Q, V, there is some subtlety to how
 151 they are used. As shown at the beginning of this paper, to calculate the attention score, K and Q are
 152 first multiplied, and then V is multiplied to that product.
 153 Attacking the POS or DEP information in K representations of late layers can keep the train loss
 154 up. This means that when the transformer is trying to generate a probabilistic prediction for the next
 155 words, it is very important for it to have POS/DEP information encoded in K representations. Second,
 156 attacking early layers can cause the model to overfit. This means that encoding linguistic information
 157 can help the transformer model not to overfit. On the contrary, attacking K and Q representations from
 158 middle layers can make the transformer model more robust. This shows that linguistic information
 159 could be a burden for the transformer at the middle layers.

160 Another perspective we want to gain from the experiments is the flow of information within the
 161 transformer model. We care about this because the information is generally transferable, the model
 162 could recover certain missing information in part of the model from other parts of the model. Looking
 163 at the training plots, we see attacking V representations in almost all layers can make the transformer
 164 model not learn on the next word prediction task. This might be because that when calculating
 165 attention scores, the V representations are the last values to be included in the calculation. So if
 166 V representations are being attacked, the model can not recover the calculation of proper attention
 167 scores that allows the transformer model to learn. When attacking Q or K, we can see from train loss
 168 that most of the time the model is able to recover from the attack to some extend. We can also see that
 169 there are some similarities between the K loss curves and Q loss curves for removing both linguistic
 170 information. This should be addressed by the form of the scaled dot product attention, where K and
 171 Q are not really distinguishable in the calculation of the attention score 1.

172 5 Results of Rewinding and Freezing Method

173 In the rewinding experiments, we will first describe the training and validation results of the original
 174 task and then investigate three methods: 1) rewinding with convex combination; 2) rewinding on
 175 mid-stage models; 3) freezing at the training stage. In each of the three methods, we modify the
 176 weights of one submodule at a time.

177 5.1 Observing Results

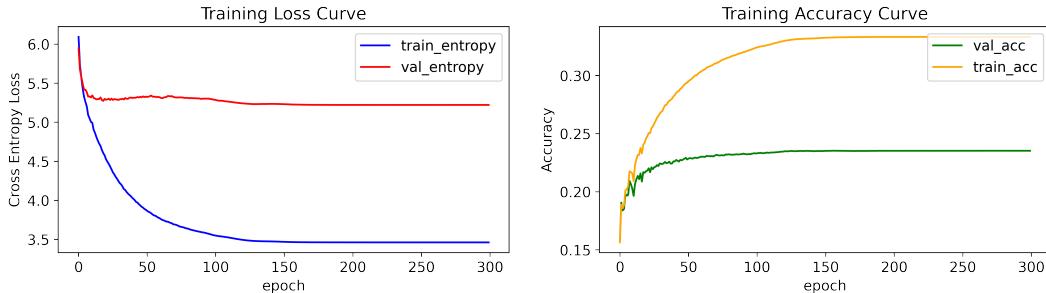


Figure 7: Training and Validation Curve of the Original Task

178 We trained a model to convergence and the training curve is shown in Figure 7. The training loss and
 179 accuracy are 3.462 and 33.32% respectively at convergence, and the validation loss and accuracy are
 180 5.222 and 23.53% at convergence. In the following experiments we set the model at epoch 200 as our
 181 final model and its weight is denoted as θ_{att}^f and the initial untrained model as θ_{att}^0 . This model is
 182 considered the baseline model.

183 5.1.1 Rewinding method 1: Convex Combination

184 In the first rewinding method, we rewind the weights of a specific submodule of a Transformer layer
 185 back to the convex combination: $\theta_{att}^\alpha = (1 - \alpha)\theta_{att}^0 + \alpha \cdot \theta_{att}^f$, where $\alpha \in [0, 1]$, with granularity
 186 0.01.

187 Figure 8 shows the evaluation performance of convex combination on the training set. The color map
 188 describes the loss change for different level of dependence on the initialization. The y-axis shows the
 189 α value in the convex combination and the smaller α stands for a greater weight on the initial state.
 190 The x-axis shows the submodule that we modified. We observe that when we rewind the weights
 191 of submodule v at layer 2,3,4, and 5, there are apparent drops while others only have slight dips in
 192 performance.

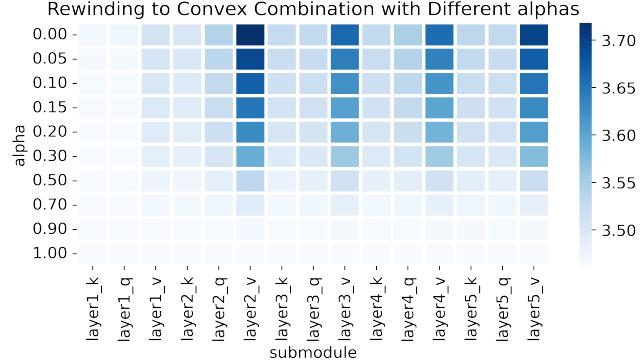


Figure 8: Rewinding Loss of convex combinations

193 Then, we use the metric *criticality* in Equation 1 to quantify the submodule importance. We set the
 194 threshold to be 1% of the final loss, which means we will test how much we may rewind the weight
 195 back to initial state while a performance tabulation is within $\epsilon = 0.0346$.

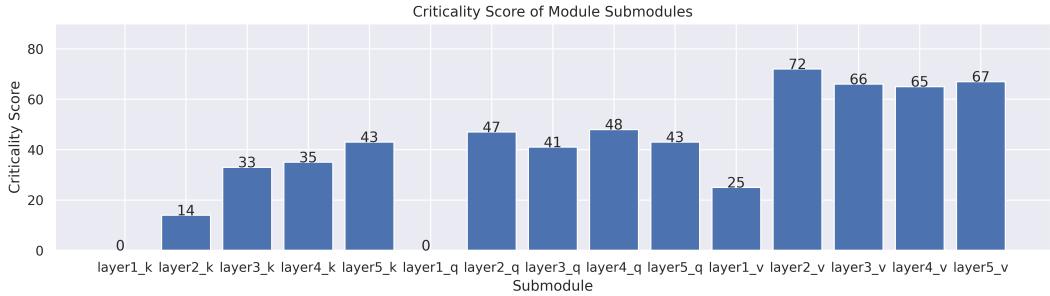


Figure 9: Criticality Score

196 Figure 9 shows the *criticality* score of each submodule. In the last two figures, we observe that

- 197 • Layer 1 has smaller criticality scores among all layers
 198 • The models with v edited have a more significant performance drop while the models with k and q
 199 edited have almost the same amount of performance dip.

200 • Upper encoder-attention layers are more critical than lower encoder-attention layers

201 **5.1.2 Rewinding method 2: Rewind the performance to the mid-stage**

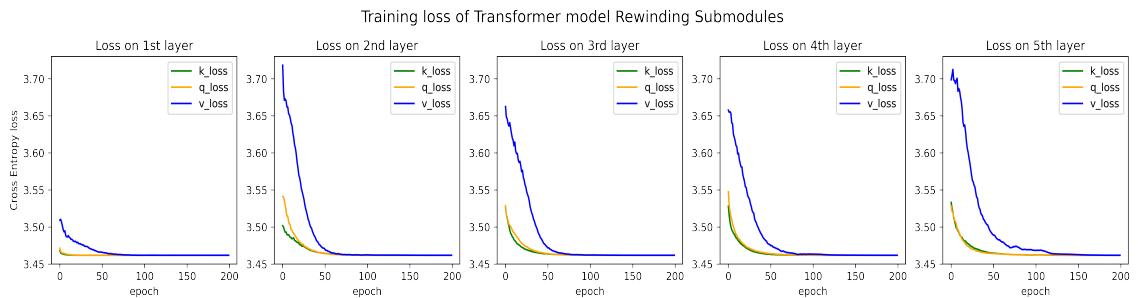


Figure 10: Rewinding Submodules

202 In this method, we edit the submodule weights by replacing the final model weights of a submodule
 203 with the mid-stage weights of that submodule: $\theta_{att}^f = \theta_{att}^i, i \in [1, 200], att \in \{k, q, v\}$ where i
 204 denotes the epoch and att denotes the attention component we edit on.

205 Figure 10 shows the loss curve if we rewind the model weights back to mid-stage weights. We
 206 observe that:

- 207 • Rewinding weights on V brings higher performance change while modifying K and Q has no huge
 208 difference
- 209 • Layer 1 is the most robust encoder layer against rewinding

210 5.1.3 Rewinding method 3: Freezing submodules during the training

211 In this method, we no longer edit the model weights during the evaluation but freeze some submodule
 212 weights throughout the training. Our goal is to see whether freezing submodules affects the model
 213 training process.

214 We freeze the submodules by setting the gradients to be zero during the training and Figure 11 and
 215 Figure 12 shows the training curves and validation curves respectively.

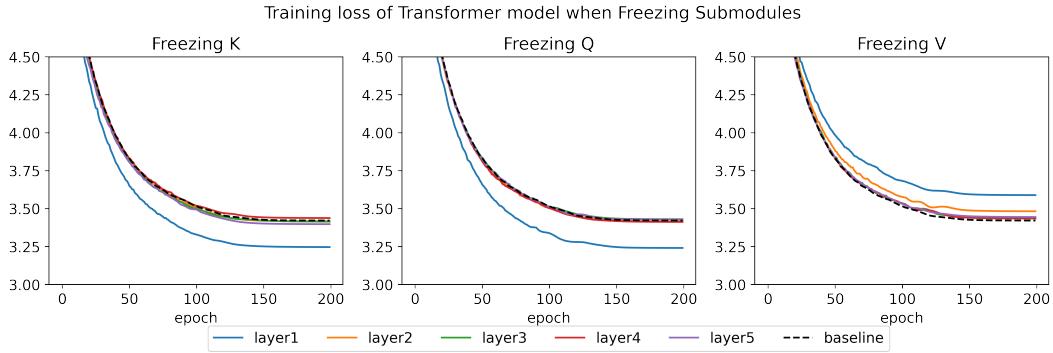


Figure 11: Performance of Freezing Submodules on Train Set

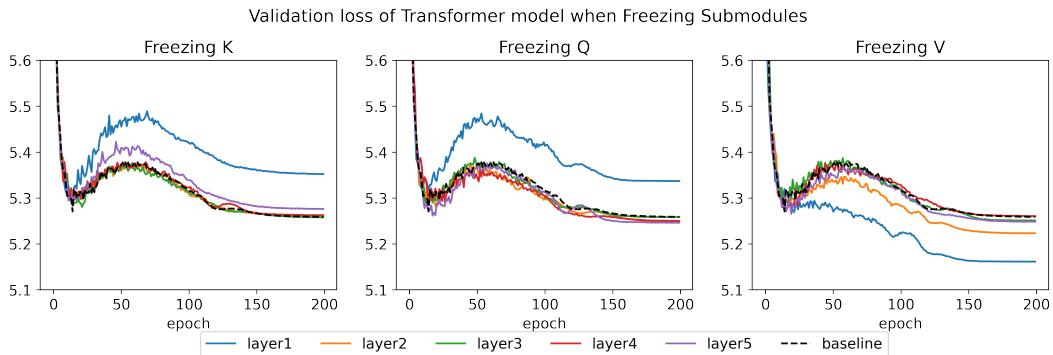


Figure 12: Performance of Freezing Submodules on Valid Set

216 We compare the training curve of the baseline model (in dashed black line) and the model with frozen
 217 submodules. There are some observations:

- 218 • In general, the models with first layer edited have more significant performance changes. Freezing
 219 K and Q submodules in layer 1 drops the training curve but increases the V validation loss.
- 220 • Freezing layer 1 causes over-fitting on K and Q but freezing V causes underfitting
- 221 • Freezing upper layer submodules has little impact on the results

222 **5.2 Discussion**

223 In this part we investigated the submodule importance in Transformer through three experiments and
224 the results show that each submodules are not equally important in different settings.

225 Self attention components of K, Q, and V in layer 1 are important during the training since it causes
226 overfitting if freezing K and Q while underfitting if freezing V. On the contrary, they are the least
227 important during the evaluation since they have very low criticality scores in the rewinding experiment.
228 The overlapped training curves show that attention components in the upper layers are not important
229 since information can be recovered through training, even they are frozen.

230 We also agree that V is more important in the attention dot product as for its high criticality scores,
231 which accords with the adversarial experiment results.

232 **6 Conclusion**

233 In this work we investigated the importance and functionalities of individual submodule component in
234 Transformer model. The experiment results reveal that some submodules in specific layers are more
235 important in the next word prediction task. The V component is more vulnerable under adversarial
236 attacking and has higher criticality scores. Also, attacking or freezing submodules in lower layers
237 can cause performance perturbation and over-fitting in component K and Q but under-fitting in V.
238 Both methods have shown that removing information from K and Q can make the model have less
239 validation loss, therefore more robust. Lastly, the Transformer model shows that it is able to recover
240 from information lost through attacking or freezing in most of the layers and submodules.

241 **7 Future Work**

242 To attain more generalizable results, the same experiment should be done on more different linguistic
243 information. To have more convincing result, it is better to run the attack on each layer and attributes
244 (K,Q,V) multiple times to generate a confidence interval. Because the convergence of the training
245 might change as per observations during the experiment.

246 It would also be interesting to run the experiments on attacking or freezing combinations of K,Q,V or
247 other parts of the transformer model (feed-forward part, decoder attentions, different attention heads,
248 etc). If time commits, it would be also interesting to inhabit both K and Q simultaneously, or test
249 the hypothesis if they are able to compensate for each other. Similarly, we may also inhabit a whole
250 module or inhabit the attention component (one of K, Q, and V) in all the layers.

251 **8 Acknowledgements and Lessons learned**

252 First of all, we would like to sincerely thank Naomi Saphra for guidance, inspiration and help made in
253 this project as an advisor. This is an interesting and fantastic project we have ever done since it's very
254 close to NLP research. We also met some difficulties understanding previous research works in this
255 area, determining the project direction, and understanding the Transformer details. Unlike applying
256 NLP models to some specific tasks, a research-oriented project does not guarantee any meaningful
257 results. Fortunately, by conducting literature review and many experiments, we finally discover some
258 interesting and surprising properties of Transformer models.

259 **9 Contribution**

260 In this Capstone project Ziyi (ZX) contributes mostly to the adversarial experiment and Letian (LJ)
261 contributes mostly to the rewinding&freezing experiment.

262 **References**

- 263 [1] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser,
264 & Illia Polosukhin. (2017). Attention Is All You Need.
- 265 [2] Logan, S. (2019). Using Knowledge Graphs for Fact-Aware Language Modeling. In Proceedings of the 57th
266 Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). Association for
267 Computational Linguistics.
- 268 [3] Yaroslav Ganin, & Victor Lempitsky. (2015). Unsupervised Domain Adaptation by Backpropagation.
- 269 [4] Stephen Merity, Caiming Xiong, James Bradbury, & Richard Socher. (2016). Pointer Sentinel Mixture
270 Models.
- 271 [5] Zhang, C., Bengio, S., & Singer, Y. (2019). Are all layers created equal?. arXiv preprint arXiv:1902.01996.
- 272 [6] Chatterji, N. S., Neyshabur, B., Sedghi, H. (2019). The intriguing role of module criticality in the
273 generalization of deep networks. arXiv preprint arXiv:1912.00528.