

---

# Betti numbers of attention graphs is all you really need

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We apply methods of topological analysis to the attention graphs, calculated on  
2 the attention heads of the BERT model (Devlin et al. (2019)). Our research shows  
3 that the classifier built upon basic persistent topological features (namely, Betti  
4 numbers) of the trained neural network can achieve classification results on par with  
5 the conventional classification method. We show the relevance of such topological  
6 text representation on three text classification benchmarks. For the best of our  
7 knowledge, it is the first attempt to analyze the topology of an attention-based  
8 neural network, widely used for Natural Language Processing.

## 9 1 Introduction

10 Modern Neural Networks embed data into a high-dimensional space. Moreover, each layer and even a  
11 layer part can be considered as separate embedding, where the information about interconnections of  
12 these separate embeddings is encoded by some weighted directed graph. In particular, one can apply  
13 various methods to investigate such graphs for attention heads in multi-headed attention models, such  
14 as BERT.

15 Conventionally, the BERT model is used for sentence classification by adding a softmax-based  
16 classification layer upon the output embedding. Instead, we propose to use a linear classifier built  
17 solely upon the persistent topological features (namely, the first two Betti numbers) without using any  
18 information about the order of tokens or to which particular token each weight relates. We have found  
19 that it provides a classification quality on par with the conventional classification method in numerous  
20 tasks. Moreover, on some tasks such as linguistic acceptability and spam detection, our topological  
21 classifiers outperform the usual BERT-based classification. We conclude that the topology of the  
22 attention graphs of the trained BERT model contains enough information for solving considered  
23 classification tasks. The second outcome is, that the proposed text representation, based only on the  
24 first two Betti numbers of the attention graph, can solve the task having lower dimensionality than  
25 BERT embedding.

26 The attention graphs are built as follows. Each attention head in the Transformer architecture  
27 calculates weights of each token in the sentence with respect to every other token, and the next level  
28 representation is constructed using these weights. The attention graph for each head is a complete  
29 digraph (with loops which appear when the token "pays attention" to itself) whose vertices are the  
30 tokens and the attention weights are the weights of the edges.

31 In the classifiers, we use the BERT-based classification model, which is initialized with pre-trained  
32 BERT weights and then is fine-tuned for a given two-class classification task. After fine-tuning, we  
33 extract the persistent features of each head of this model for each data sample and then train a logistic  
34 regression classifier upon these features.

35 Note that our results also confirm that different attention heads contain different amounts of informa-  
 36 tion. These results are well aligned to previous works on BERT (Michel, Levy, and Neubig (2019),  
 37 Clark et al. (2019)).

## 38 2 Related work

39 There are several recent insights obtained by topological analysis of the neural representations of  
 40 realistic datasets. The results of Naitzat, Zhitnikov, and Lim (2020) demonstrate that a deep neural  
 41 network with *ReLU* activation function tends to simplify the topology of the data from layer to layer,  
 42 with the smallest Betti numbers on the output representations. Topological features are also shown  
 43 to be efficient for predicting the generalization ability of the network, its efficiency and stability  
 44 to adversarial examples (Corneanu, Escalera, and Martinez (2020), Corneanu et al. (2019), Rieck  
 45 et al. (2018)). An overview of persistent topology methods, both practical and theoretical, with the  
 46 focus on Artificial Neural Networks analysis, can be found in Otter et al. (2017), Chazal and Michel  
 47 (2017). At the same time, while many researchers are focused on applications of topology to modern  
 48 AI algorithms, there are efforts in the mathematical society to further expand the set of applicable  
 49 methods (Bergomi et al. (2019), Chowdhury et al. (2019), Manin and Marcolli (2020)).

## 50 3 Background

### 51 3.1 Topological background

52 In our approach, we use the following two numerical attributes of an arbitrary graph  $G = (E, V)$ : the  
 53 number  $\beta_0$  of connected components and the number  $\beta_1$  of independent cycles of  $G$ . If one considers  
 54 the graph  $G$  as a simplicial complex, these numbers are equal to its Betti numbers. Note that the Betti  
 55 numbers  $\beta_0$  and  $\beta_1$  of a graph filtration keep the whole information about the persistent homology  
 56 barcodes, see Appendix A for details.

### 57 3.2 BERT model

58 BERT (Devlin et al. (2019)) is the pre-trained model, which achieves state of the art results for many  
 59 NLP tasks. The model is based on Transformer architecture, introduced in Vaswani et al. (2017). The  
 60 BERT model is pre-trained on the large amount of data with Masked Language Modelling and Next  
 61 Sentence Prediction objectives. For downstream tasks the task-specific classifier is attached to the  
 62 BERT output layer and the model is fine-tuned. In our experiments we use the uncased BERT-base  
 63 version, which consists of 12 layers, with 12 attention heads in each. The input of each attention head  
 64 is a matrix  $X$  consisting of the  $d$ -dimensional representations (row-wise) of  $m$  tokens of the sentence,  
 65 so that  $X$  is of size  $m \times d$ . The output of the head is the updated matrix of the representations  $X^{\text{out}}$ ,  
 66 that is,

$$X^{\text{out}} = W^{\text{attn}}(XW^V)$$

$$\text{with } W^{\text{attn}} = \text{softmax}\left(\frac{(XW^Q)(XW^K)^T}{\sqrt{d}}\right), \quad (1)$$

67 where  $W^Q, W^K, W^V$  are trained projection matrices of size  $d \times d$  and  $W^{\text{attn}}$  is the  $m \times m$  matrix  
 68 of attention weights (cf. (Vaswani et al., 2017, Sec. 3.2)). One can interpret each element  $w_{ij}^{\text{attn}}$  as  
 69 a weight of  $j$ -th input’s influence on  $i$ -th output; larger weights mean stronger connection between  
 70 corresponding tokens.

## 71 4 Our method

72 Let us be given some dataset  $S = \{s_i\}_{i=1}^N$  of  $N$  natural language texts encoded with  $m$  tokens each  
 73 and pre-trained attention-based model  $M$ . First of all, we fix some set of thresholds  $T = \{t_i\}_{i=1}^k, 0 <$   
 74  $t_1 < t_2 < \dots < t_k < 1$  and chose a subset of heads of the model  $H_M$ , on which we will perform  
 75 calculations.

76 Then we feed each text sample  $s = s_i$  to the input of the model  $M$  and obtain the matrix  $W^{\text{attn}} =$   
 77  $(w_{i,j}^{\text{attn}})$  on each head  $h \in H_M$ . This matrix defines a weighted complete digraph with loops  $\Gamma_s^h$  with  
 78  $m$  vertices, where  $w_{ij}^{\text{attn}}$  is the weight of the edge  $j \rightarrow i$ .

79 After it, for each graph  $\Gamma_s^h$  and for each threshold level  $t_i \in T$  we build an unweighted directed graph  
80  $\Gamma_s^h(t_i)$  as follows. The set of vertices of  $\Gamma_s^h(t_i)$  is the same as the one for the graph  $\Gamma_s^h$ , moreover, an  
81 edge of  $\Gamma_s^h$  belongs to the new graph  $\Gamma_s^h(t_i)$  if and only if its weight in  $\Gamma_s^h$  is at least  $t_i$ . This way we  
82 assign a sequence of graphs  $\Gamma_s^h(t_i), t_i \in T$  to each text sample for each head of the model.

83 For each unweighted directed graph  $\Gamma_s^h(t_i)$  we also consider the corresponding undirected graph  
84  $\overline{\Gamma_s^h(t_i)}$  by setting an undirected edge  $v_1v_2$  for each pair of vertices  $v_1$  and  $v_2$  which are connected by  
85 an edge in at least one direction in the graph  $\Gamma_s^h(t_i)$ . Then we count  $\beta_0, \beta_1$  of undirected graph  $\overline{\Gamma_s^h(t_i)}$ .  
86 More precisely the process of features calculation for each data sample is described in Algorithm 1.

---

**Algorithm 1** Topological features calculation

---

**Require:** Text sample  $s$

**Require:** Set of chosen attention heads  $H_M$  of attention-based model  $M$

**Require:** Thresholds array  $T$

**Ensure:** Features array  $Features$

```

2:   procedure FEATURES_CALCULATION( $s, H_M, T$ )
3:     for all  $h \in H_M$  do
4:       Calculate attention graph  $\Gamma_s^h = (V, E, W_{h,s}^{attn})$  on sample  $s$  on head  $h$ 
5:       for all  $t \in T$  do
6:          $E_s^h(t) \leftarrow \{e \in E(\Gamma_s^h) : W_{h,s}^{attn}(e) \geq t\}$        $\triangleright$  Filtration:
7:          $\Gamma_s^h(t) \leftarrow (V, E_s^h(t))$                                  $\triangleright$  Ignoring weights of remaining edges
8:          $\overline{E_s^h(t)} \leftarrow \{\{i, j\} : (i, j) \in E_s^h(t)\}$        $\triangleright$  Ignoring edges directions
9:          $\overline{\Gamma_s^h(t)} \leftarrow (V, \overline{E_s^h(t)})$ 
10:        Calculate  $\beta_0(\overline{\Gamma_s^h(t)}), \beta_1(\overline{\Gamma_s^h(t)})$      $\triangleright$  Calculating Betti numbers of undirected graph
11:      end for
12:    end for
13:     $Features \leftarrow \left[ \beta_0(\overline{\Gamma_s^h(t)}), \beta_1(\overline{\Gamma_s^h(t)}) \right]_{t \in T}^{h \in H_M}$ 
14:  return  $Features$ 
15: end procedure

```

---

87 After this features calculation, we train the logistic regression on features, obtained for each sample of  
88 the train subset of the dataset and then make predictions on features of samples from the test subset.

## 89 5 Experiments

### 90 5.1 Datasets

91 We performed our experiments on the following datasets, labeled for different classification tasks.

92 **The Corpus of Linguistic Acceptability** ("CoLA") dataset (Warstadt, Singh, and Bowman (2018))  
93 contains 10,657 sentences, labeled by acceptability (grammaticality) and divided into public (open)  
94 and test (hidden) parts. The public part of dataset contains 9,594 sentences and is divided, in turn,  
95 into training and development ("CoLA<sub>dev</sub>") sets. The test set ("CoLA<sub>test</sub>") contains 1,063 sentences  
96 with labels, hidden from the developer.

97 **Large Movie Review Dataset v1.0** ("IMDB") (Maas et al. (2011)) contains 50,000 movie reviews,  
98 labeled by sentiment: "positive" or "negative". Labeled reviews are divided into two equal subsets,  
99 purposed for training and for testing. We applied additional lengths restriction to the samples of this  
100 dataset to obtain attention graphs of a reasonable size. Namely, we kept all reviews of size less than  
101 128 tokens after tokenization with standard BERT uncased tokenizer ("Imdb <sup>$\leq 128$</sup> "), and pruned away  
102 all others. After it, 5505 reviews remained in total. Then we divided the subset, suggested for testing  
103 purposes, into equal development and test sets.

104 **The SMS Spam Collection v.1** ("SPAM") (Almeida, Hidalgo, and Yamakami (2011)) is a public set  
105 of SMS (text) labeled messages that have been collected for mobile phone spam research. It contains

106 5,574 real and non-encoded messages, tagged as legitimate (ham) or spam. For our purposes, we  
 107 divided it into train, development and test ("SPAM<sub>test</sub>") sets randomly in proportion 80 : 10 : 10.

108 We used "development" subsets for tuning logistic regression hyperparameters: maximum amount of  
 109 iterations and  $l_2$ -regularization coefficient. "Test" subsets were used for final validation.

## 110 5.2 Results

	CoLA <sub>dev</sub>	CoLA <sub>test</sub>	Imdb <sub>test</sub> <sup>≤128</sup>	SPAM <sub>test</sub>
BERT	0.559 (82.0%)	0.492	0.833 (91.7%)	0.941 (98.7 %)
$\beta_0, \beta_1, 144$ heads	0.549 (81.1 %)	<b>0.508</b>	0.812 (90.6 %)	<b>0.950 (98.9 %)</b>
$\beta_0, \beta_1, 12$ best heads	0.532 (80.7 %)	0.463	0.805 (90.3 %)	0.878 (97.3 %)
$\beta_0, \beta_1, 3$ best heads	0.452 (77.1 %)	0.456	0.799 (90.0 %)	0.809 (96.1 %)
$\beta_0, \beta_1, 1$ best head	0.427 (76.4 %)	0.385	0.735 (86.9%)	0.606 (92.3 %)
Test examples amount	1043	1033	1415	556

Table 1: The comparison of our classification methods with the conventional BERT-based classifier by Matthew score and accuracy (in brackets, %).

As an efficiency measure of a linear classifier, we use Matthew score (Matthew coefficient), which is calculated by formula

$$MCC = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}},$$

111 where we denote by  $FP, FN, TN$ , and  $TP$  the amount of false positive, false negative, true negative,  
 112 and true positive predictions of our classifier, respectively. We also note classifications accuracy in  
 113 brackets for those datasets, where test labels are available in open access.

114 For these experiments we fine-tuned BERT on each of datasets separately and used the set of six  
 115 weight thresholds for calculating Betti numbers. In Table 1 we emphasized in bold the results which  
 116 surpassed the result of the fine-tuned BERT classifier.

117 For the first experiment we used the features calculated on all 144 heads. For consequent experiments,  
 118 we checked Matthew score of classification upon features, built from the graph on each head on the  
 119 train set, and ranged heads in descending order according to it. Then we picked 12, 3 or 1 heads with  
 120 the best Matthew score and used them for calculation of classification features (Betti numbers) on  
 121 development and test sets.

122 It's noticeable that topological features of attention graphs on particular heads have different linear  
 123 separability. For more information about this see Appendix B.

## 124 6 Conclusion and further research

125 We have shown that the topology of attention graphs contains enough information for classifying  
 126 texts by three different attributes: linguistic acceptability, sentiment, and being SPAM or not. Thus,  
 127 we see here some degree of universality for distinguishing different text properties.

128 Moreover, the result of our linear classifier, trained on topological features, surpassed the result  
 129 of the conventional BERT-based classifier on the hidden test subset of the Corpus of Linguistic  
 130 Acceptability dataset and is a little better on the SMS Spam Collection v. 1 dataset. This allows us to  
 131 suppose that these features may contain even more generalized task-relevant information than the  
 132 BERT output embedding. Plans of our future research include checking this daring statement with  
 133 other topological features and other threshold collections. Particularly, in our current work we didn't  
 134 use the information about directions of graph edges, which could be utilized with directed graph  
 135 invariants, such as number of simple directed cycles and number of strongly connected components  
 136 of a digraph.

137 Another possible direction for future work is to use the information about differences between linear  
 138 separability scores on different heads to determine which heads are more or less important for each  
 139 particular task. Which can potentially be used as a base for new strategies of efficiently decreasing  
 140 the model size.

141 **References**

- 142 Almeida, T. A.; Hidalgo, J. M. G.; and Yamakami, A. 2011. Contributions to the study of SMS  
143 spam filtering: new collection and results. In Hardy, M. R. B.; and Tompa, F. W., eds., *ACM*  
144 *Symposium on Document Engineering*, 259–262. ACM. ISBN 978-1-4503-0863-2. URL <http://dblp.uni-trier.de/db/conf/doceng/doceng2011.html#AlmeidaHY11>.  
145
- 146 Bergomi, M. G.; Frosini, P.; Giorgi, D.; and Quercioli, N. 2019. Towards a topological–geometrical  
147 theory of group equivariant non-expansive operators for data analysis and machine learning. In  
148 *Nature Machine Intelligence* 1.9, 423–433.
- 149 Chazal, F.; and Michel, B. 2017. An introduction to Topological Data Analysis: fundamental and  
150 practical aspects for data scientists. *ArXiv* abs/1710.04019. URL <https://arxiv.org/pdf/1710.04019.pdf>.  
151
- 152 Chowdhury, S.; Gebhart, T.; Huntsma, S.; and Yutin, M. 2019. Path homologies of deep feedforward  
153 networks. In *18th IEEE International Conference On Machine Learning And Applications (ICMLA)*,  
154 1077–1082.
- 155 Clark, K.; Khandelwal, U.; Levy, O.; and Manning, C. D. 2019. What Does BERT Look At? An  
156 Analysis of BERT’s Attention. *CoRR* abs/1906.04341. URL <https://arxiv.org/pdf/1906.04341.pdf>.
- 157 Corneanu, C. A.; Escalera, S.; and Martinez, A. M. 2020. Computing the Testing Error without  
158 a Testing Set. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*  
159 *Recognition*, 2677–2685.
- 160 Corneanu, C. A.; Madadi, M.; Escalera, S.; and Martinez, A. M. 2019. What does it mean to learn  
161 in deep networks? And, how does one detect adversarial attacks? In *Proceedings of the IEEE*  
162 *Conference on Computer Vision and Pattern Recognition*, 4757–4766.
- 163 Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2019. BERT: Pre-training of Deep Bidirectional  
164 Transformers for Language Understanding. In *NAACL-HLT (1)*.
- 165 Maas, A. L.; Daly, R. E.; Pham, P. T.; Huang, D.; Ng, A. Y.; and Potts, C. 2011. Learning Word  
166 Vectors for Sentiment Analysis. In *Proceedings of the 49th Annual Meeting of the Association for*  
167 *Computational Linguistics: Human Language Technologies*, 142–150. Portland, Oregon, USA:  
168 Association for Computational Linguistics. URL <http://www.aclweb.org/anthology/P11-1015>.
- 169 Manin, Y.; and Marcolli, M. 2020. Homotopy Theoretic and Categorical Models of Neural Informa-  
170 tion Networks. *ArXiv* abs/2006.15136. URL <https://arxiv.org/pdf/2006.15136.pdf>.
- 171 Michel, P.; Levy, O.; and Neubig, G. 2019. Are Sixteen Heads Really Better than One? In  
172 Wallach, H.; Larochelle, H.; Beygelzimer, A.; d’Álché Buc, F.; Fox, E.; and Garnett, R., eds.,  
173 *Advances in Neural Information Processing Systems* 32, 14014–14024. Curran Associates, Inc.  
174 URL <http://papers.nips.cc/paper/9551-are-sixteen-heads-really-better-than-one.pdf>.
- 175 Naitzat, G.; Zhitnikov, A.; and Lim, L.-H. 2020. Topology of deep neural networks. *ArXiv*  
176 abs/2004.06093. URL <https://arxiv.org/pdf/2004.06093.pdf>.
- 177 Otter, N.; Porter, M. A.; Tillmann, U.; Grindrod, P.; and Harrington, H. A. 2017. A roadmap for the  
178 computation of persistent homology. In *EPJ Data Science* 6.1, 17.
- 179 Rieck, B.; Togninalli, M.; Bock, C.; Moor, M.; Horn, M.; Gumbsch, T.; and Borgwardt, K. 2018.  
180 Neural Persistence: A Complexity Measure for Deep Neural Networks Using Algebraic Topology.  
181 In *International Conference on Learning Representations*.
- 182 Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, L.; and  
183 Polosukhin, I. 2017. Attention Is All You Need. In *Advances in neural information processing*  
184 *systems*, 5998–6008.
- 185 Warstadt, A.; Singh, A.; and Bowman, S. R. 2018. Neural Network Acceptability Judgments. *arXiv*  
186 preprint arXiv:1805.12471 .

187 **Appendix A. Persistent homology and Betti numbers**

188 Recall that a simplicial complex  $K$  is a finite collection of finite sets called *simplices* such that each  
 189 subset of any element of  $K$  also is an element of  $K$ ; such subsets of a simplex are called *faces*. In  
 190 particular, an undirected graph is a simplicial complex where all edges and vertices are its faces.  
 191 The set of all formal  $\mathbb{Z}$ -linear combinations of the  $p$ -dimensional simplices (that is,  $(p - 1)$ -element  
 192 simplices) of  $K$  is denoted  $C_p(K)$ . These linear combinations  $c = \sum_j \gamma_j \sigma_j$  are called  $p$ -chains,  
 193 where the  $\gamma_j \in \mathbb{Z}$  and the  $\sigma_j$  are  $p$ -simplices in  $K$ .

194 The boundary,  $\partial(\sigma_j)$ , is the formal sum of the  $(p - 1)$ -dimensional faces of  $\sigma_j$  and the boundary of  
 195 the chain is obtained by extending  $\partial$  linearly,

$$\partial(c) = \sum_j \gamma_j \partial(\sigma_j),$$

196 with integer coefficients  $\gamma_j$ .

197 The  $p$ -chains that have boundary 0 are called  $p$ -cycles, they form a subgroup  $Z_p(K)$  of  $C_p(K)$ . The  
 198  $p$ -chains that are the boundary of  $(p + 1)$ -chains are called  $p$ -boundaries and form a subgroup  $B_p(K)$   
 199 of  $C_p(K)$ . The quotient group  $H_p(K) = C_p(K)/B_p(K)$  is called the  $p$ -th *homology* of  $K$ . Their  
 200 ranks  $\beta_p = \text{rank } H_p(K)$  of these abelian groups are called *Betti numbers*. The homologies and the  
 201 Betti numbers are classical topological invariants of  $K$ .

202 In particular, a graph  $G = (E, V)$  contains only 0-dimensional and 1-dimensional faces. It follows  
 203 that its topological form is essentially described by the numbers  $\beta_0$  and  $\beta_1$  which are the only nonzero  
 204 Betti numbers. Here  $\beta_0$  is the number of connected components of  $G$ , and  $\beta_1$  is the number of  
 205 independent cycles of the graph (which is equal to  $|E| - |V| + \beta_0$ ).

206 A *subcomplex* of  $K$  is a subset of simplices that is closed under the face relation. A *filtration* of  $K$  is  
 207 a nested sequence of subcomplexes that starts with the empty complex and ends with the complete  
 208 complex,

$$\emptyset \subset K_1 \subset K_2 \subset K_3 \subset \dots \subset K_m = K.$$

209 In particular, to any weighted undirected graph  $G = (V, E)$  and an increasing sequence  $0 = t_0 \leq$   
 210  $t_1 \leq \dots \leq t_m$  such that  $t_m$  is greater or equal to the maximal edge weight in  $G$ , one can associate a  
 211 filtration

$$\emptyset \subset G_{t_0} \subset \dots \subset G_{t_m} = G, \tag{2}$$

212 where  $G_{t_i} = (V, E_{t_i})$  and  $E_{t_i}$  consists of all edges of  $E$  with weight more or equal to  $t_i$ .

213 The  $p$ -th persistent homology of  $K$  is the pair of sets of vector spaces  $\{H_p(K_i) | 0 \leq i \leq l\}$  and  
 214 maps  $\{f_{i,j} : H_p(K_i) \rightarrow H_p(K_j) | 1 \leq i < j \leq l\}$ , where the maps are induced by the inclusion maps  
 215  $K_i \rightarrow K_j$ .

216 Each persistent homology class  $\alpha$  in this sequence is “born” at some  $K_i$  and “dies” at some  $K_j$ . One  
 217 can visualize this as an interval  $[i, j]$ . The collection of all such intervals is called the *barcode* of the  
 218 filtration. It is the most useful invariant of the filtration. Note that the information about the persistent  
 219 homology classes is generally essential to calculate the barcode, whereas the information about the  
 220 Betti numbers only is insufficient.

221 Still, in the case of the filtration associated to a weighted graph (2), the basis of  $H_0$  (respectively,  
 222  $H_1$ ) gives the intervals of the form  $[0, t_i]$  (resp.,  $[t_i, t_m)$ ) only. Given a number  $l = t_k$ , the number of  
 223 intervals of length at most  $l$  for  $H_0$  (respectively, the number of intervals of of length at least  $t_m - l$   
 224 for  $H_1$ ) is therefore equal to the the Betti number  $\beta_0(K_l)$  (resp.,  $\beta_1(K_l)$ ). We see that in this case the  
 225 collection of the Betti numbers  $\beta_i(K_{t_j})$  is sufficient to recover the barcode. Thus, we use just Betti  
 226 numbers of the subgraphs  $K_{t_j}$  as the only topological invariants of our graphs.

227 **Appendix B. Classifiers built by a single head**

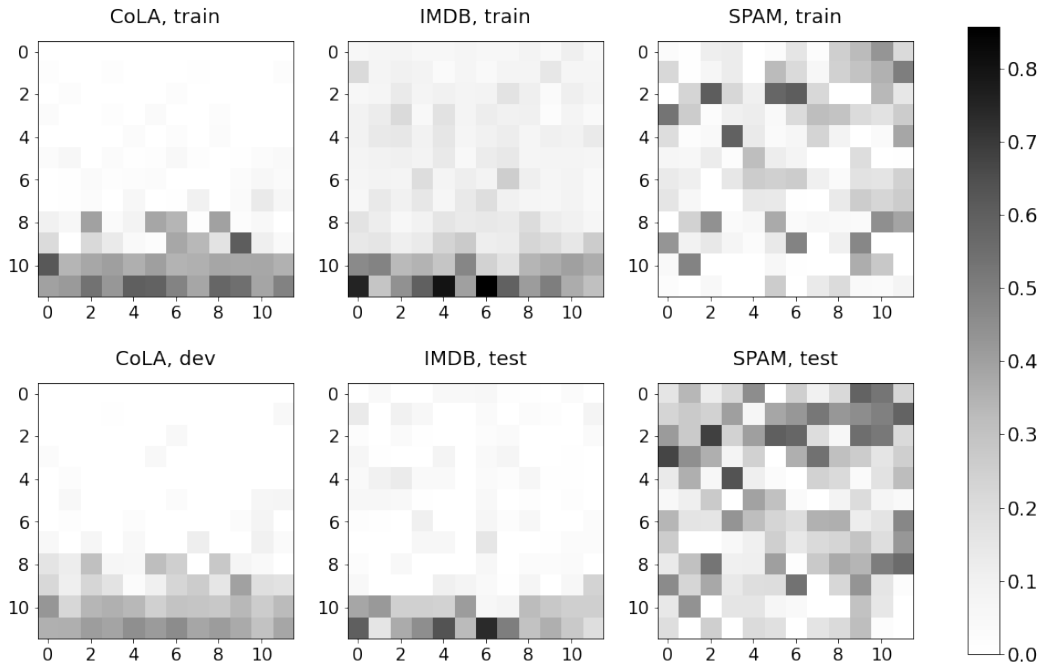


Figure 1: Matthew scores of predictions of linear classifiers, built upon particular attention heads. The number of the layer is displayed on the vertical axis. The number of the head inside the layer is displayed on the horizontal axis.

228 Figure 1 illustrates that the relevance of features, calculated on different heads, varies greatly from  
229 head to head on each task. It also shows that the same head can be more relevant for solving one task  
230 but less relevant for solving other ones. On the other hand, we can see similar patterns on the train  
231 and test/development sets for each task separately (in each column of Figure 1). This means that the  
232 head importance, derived from this score, is generalized to unseen examples and therefore can be  
233 used for feature selection.