Table 1: Additional experimenst on GPT-J and GPT2-XL to illustrate our method's generalization ability. The mean and standard deviation are reported for 3 repetitions with different ICL examples.

	GPT-J					GPT2-XL				
Method	ICL	CC	DC	PC	UniBias	ICL	CC	DC	PC	UniBias
SST-2	90.67 <sub>1.06</sub>	90.75 <sub>1.88</sub>	91.90 <sub>1.12</sub>	89.14 <sub>1.68</sub>	<b>92.01</b> <sub>1.04</sub>	56.69 <sub>3.39</sub>	81.35 <sub>4.36</sub>	87.00 <sub>2.66</sub>	<b>88.69</b> <sub>1.03</sub>	<b>88.69</b> <sub>1.10</sub>
WiC	50.21 <sub>0.99</sub>	51.83 <sub>0.53</sub>	50.89 <sub>0.90</sub>	51.67 <sub>0.59</sub>	<b>52.77</b> <sub>0.52</sub>	49.060.25	$50.88_{0.92}$	50.31 <sub>1.33</sub>	48.69 <sub>1.79</sub>	<b>51.22</b> <sub>0.69</sub>
COPA	47.33 0.47	$47.00_{0.35}$	$48.66_{0.94}$	52.99 <sub>5.88</sub>	<b>55.33</b> <sub>2.62</sub>	49.00 1.63	46.331.70	49.33 <sub>1.25</sub>	<b>52.33</b> <sub>2.36</sub>	52.00 <sub>3.74</sub>
MR	89.261.65	85.31 <sub>2.36</sub>	90.29 <sub>0.45</sub>	89.660.27	<b>90.39</b> <sub>0.36</sub>	51.11 <sub>1.76</sub>	$72.00_{4.44}$	82.631.80	83.07 <sub>1.56</sub>	<b>83.53</b> <sub>1.02</sub>
RTE	50.543.19	53.99 <sub>0.68</sub>	54.13 <sub>1.19</sub>	$49.58_{4.85}$	<b>56.20</b> <sub>1.80</sub>	52.37 <sub>0.34</sub>	$52.91_{0.45}$	53.15 <sub>0.29</sub>	$53.30_{1.28}$	<b>56.16</b> <sub>1.23</sub>
Avg.	65.60	65.78	67.17	66.61	69.34	51.65	60.69	64.48	65.32	66.32

Table 2: Additional experiments on streamlining grid search by adopting a fixed set of thresholds. The **best** and <u>second-best</u> results are marked by bold and undrline, respectively. Experiments are conducted using Llama2-7b.

	SST2	MMLU	COPA	RTE	MR	Trec	Avg.
ICL	87.226.03	41.732.25	67.602.30	66.217.30	89.371.83	72.9212.42	70.84
CC	92.24 <sub>3.39</sub>	43.72 <sub>0.97</sub>	$67.80_{2.17}$	64.33 <sub>3.68</sub>	$91.77_{1.42}$	76.443.21	72.72
DC	94.15 <sub>1.22</sub>	$43.57_{1.38}$	$60.40_{2.79}$	65.49 <sub>2.09</sub>	92.35 <sub>0.23</sub>	77.163.94	72.19
PC	93.90 <sub>1.54</sub>	34.123.41	67.803.70	62.594.71	91.39 <sub>1.65</sub>	74.925.78	70.79
UniBias	<b>94.54</b> <sub>0.62</sub>	<b>44.83</b> <sub>0.24</sub>	<b>69.00</b> <sub>2.74</sub>	<b>67.65</b> <sub>6.44</sub>	92.19 <sub>0.37</sub>	<b>80.80</b> <sub>3.17</sub>	74.84
UniBias with fixed thresholds	<u>94.42</u> 0.80	<u>44.47</u> 0.93	<u>68.20</u> 1.79	<u>67.51</u> 4.97	92.35 <sub>0.17</sub>	<u>80.32</u> 4.40	74.54

Table 3: Shared biased attention heads and their frequency of occurrence across 12 datasets in the Llama2-7b model. Each entry represents a specific (layer index, head index) combination.

(19, 10)	(19, 14)	(16, 29)	(19, 21)	(25, 21)	(16, 11)	(18, 31)	(18, 1)
6	5	4	4	3	3	3	3

Table 4: Additional experiments on eliminating common biased components. Attention heads list in Table 3 are removed from the original Llama2-7b model, and the modified model is then evaluated across multiple tasks.

	SST2	MMLU	COPA	RTE	MR	Trec	Avg.
ICL Unibias	87.22 <sub>6.03</sub> 94.54 <sub>0.62</sub>	41.73 <sub>2.25</sub> <b>44.83</b> <sub>0.24</sub>	67.60 <sub>2.30</sub> 69.00 <sub>2.74</sub>	66.21 <sub>7.30</sub> 67.65 <sub>6.44</sub>	$\begin{array}{c} 89.37_{1.83} \\ 92.19_{0.37} \end{array}$	72.92 <sub>12.42</sub> <b>80.80</b> <sub>3.17</sub>	70.84 <b>74.84</b>
Eliminating Common Biased Components	94.32 <sub>0.60</sub>	44.201.14	68.00 <sub>2.87</sub>	$67.37_{4.60}$	<b>92.43</b> <sub>0.09</sub>	77.60 <sub>4.75</sub>	73.98



1

Figure 1: Performance of Unibias using unlabeled samples as support set. It is compared against standard ICL and the original Unibias method.

Identified Biased Attention Heads in Llama2-7b



Figure 2: Identified biased attention heads across 12 datasets with 5 repetitions for each dataset.