



# VI-OOD: A Unified Representation Learning Framework for Textual Out-of-distribution Detection

Li-Ming Zhan, Bo LIU, Xiao-Ming Wu

Department of Computing

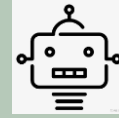
The Hong Kong Polytechnic University

# Problem Definition



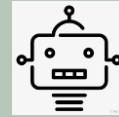
My password was changed.

Don't worry, I can help you find it!



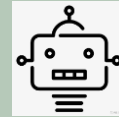
I haven't got my credit card.

Sorry for that. Your card will be arriving soon.



I recently found my card!  
Could you help me re-establish it?

Sorry, I don't understand your question.



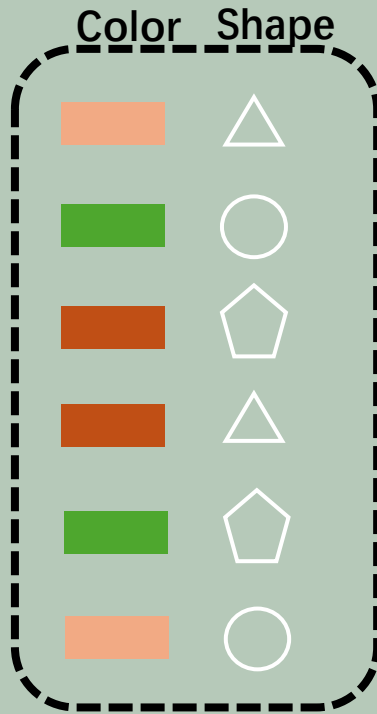
known intent classes

Unknown intent classes

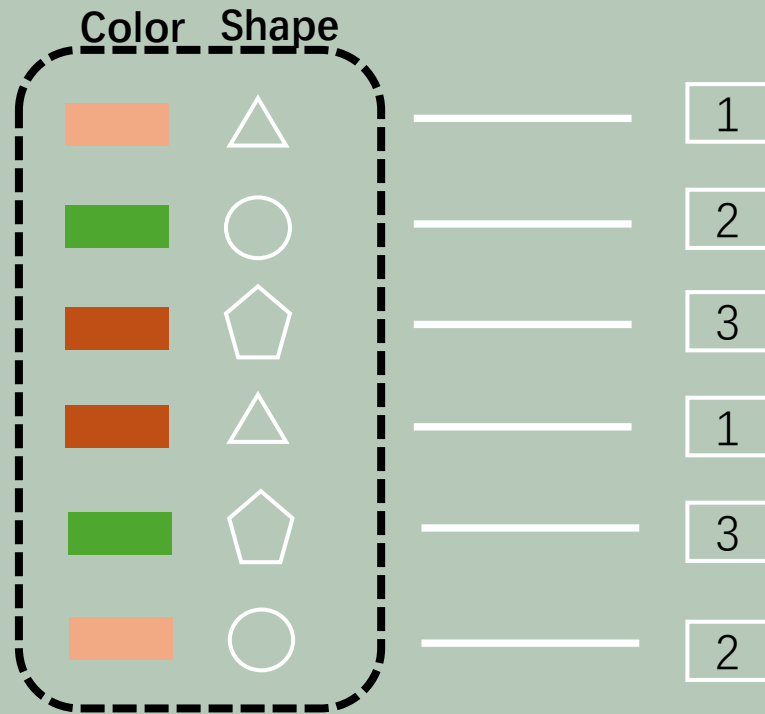
The goal of textual OOD detection is to raise an alarm or become aware of unknown inputs while keeping the known in-distribution (ID) classes correctly classified.

# Motivation

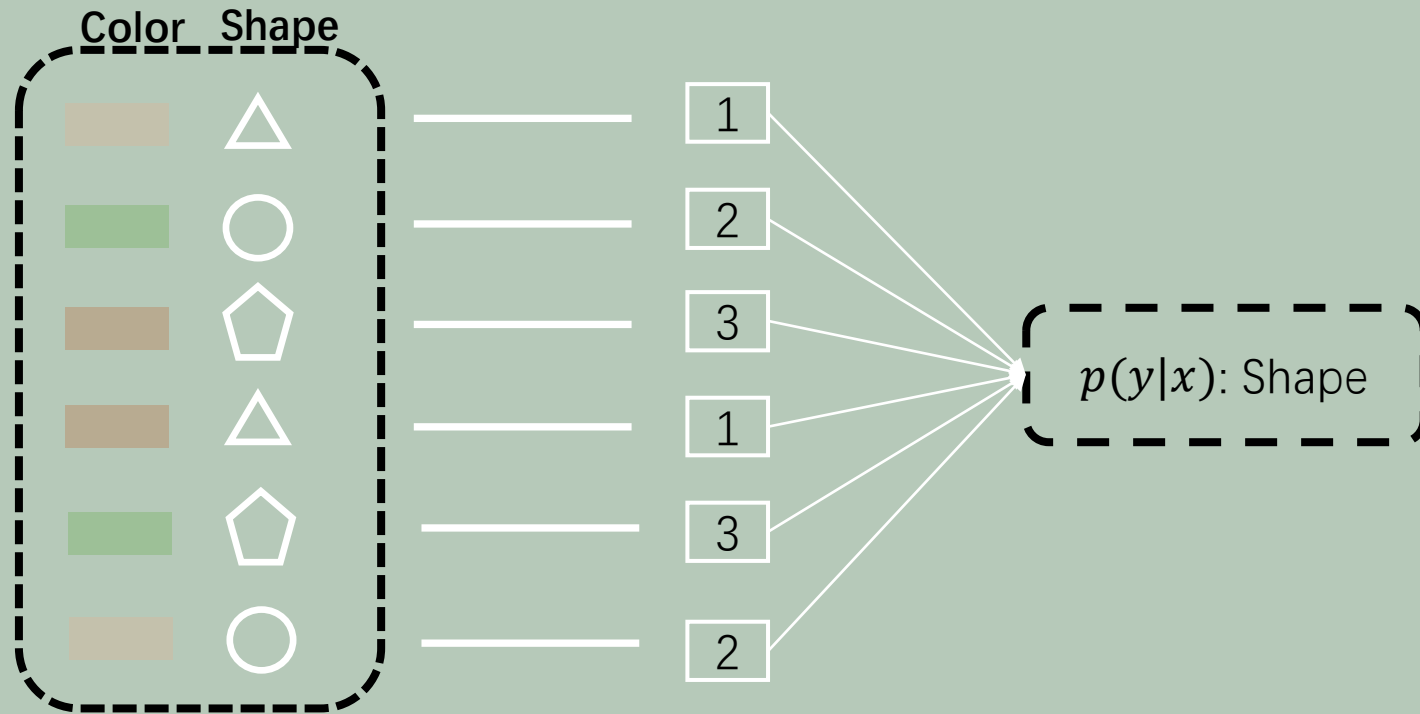
The training objective,  $p(y|x)$ , is biased and does not take into account the ability to handle OOD inputs.



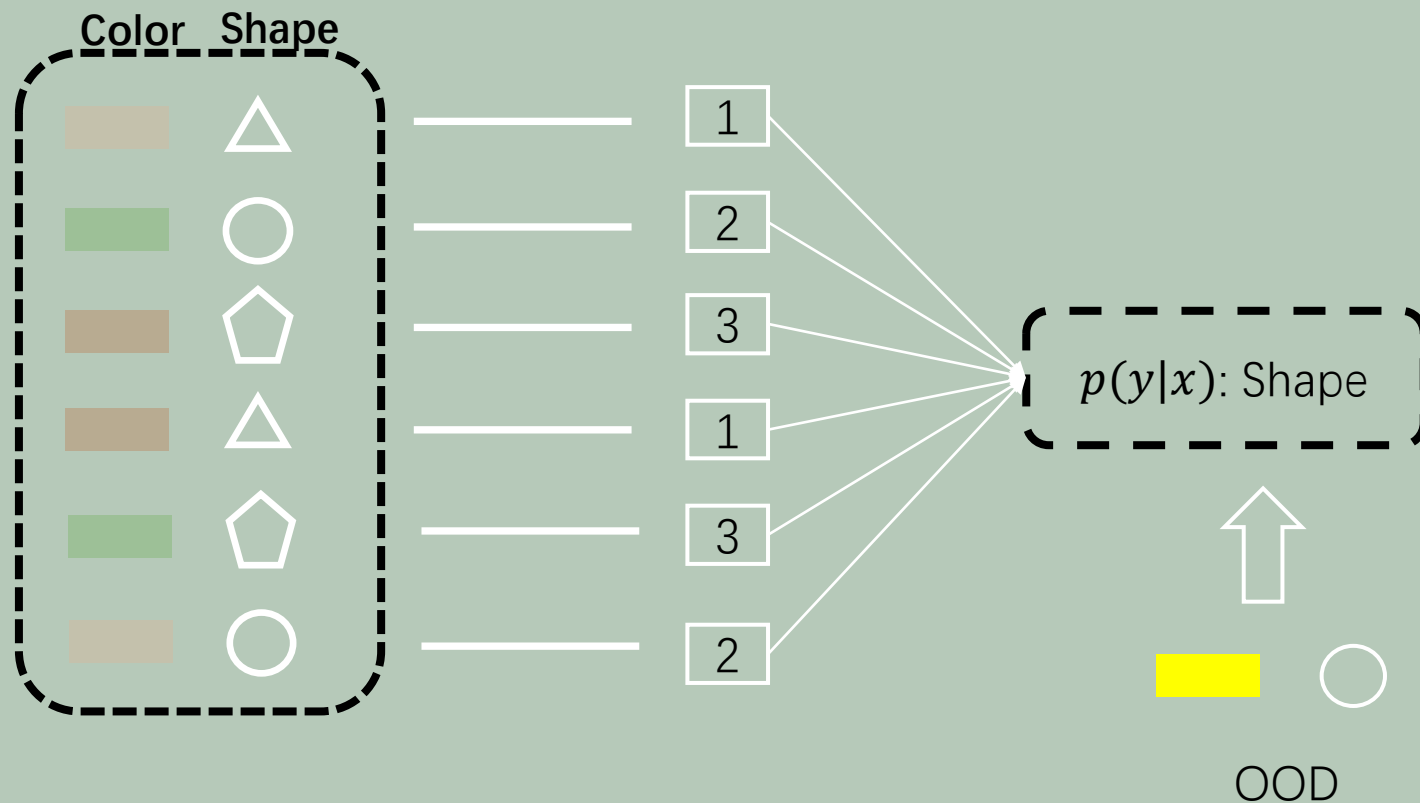
The training objective,  $p(y|x)$ , is biased and does not take into account the ability to handle OOD inputs.



The training objective,  $p(y|x)$ , is biased and does not take into account the ability to handle OOD inputs.



The training objective,  $p(y|x)$ , is biased and does not take into account the ability to handle OOD inputs.



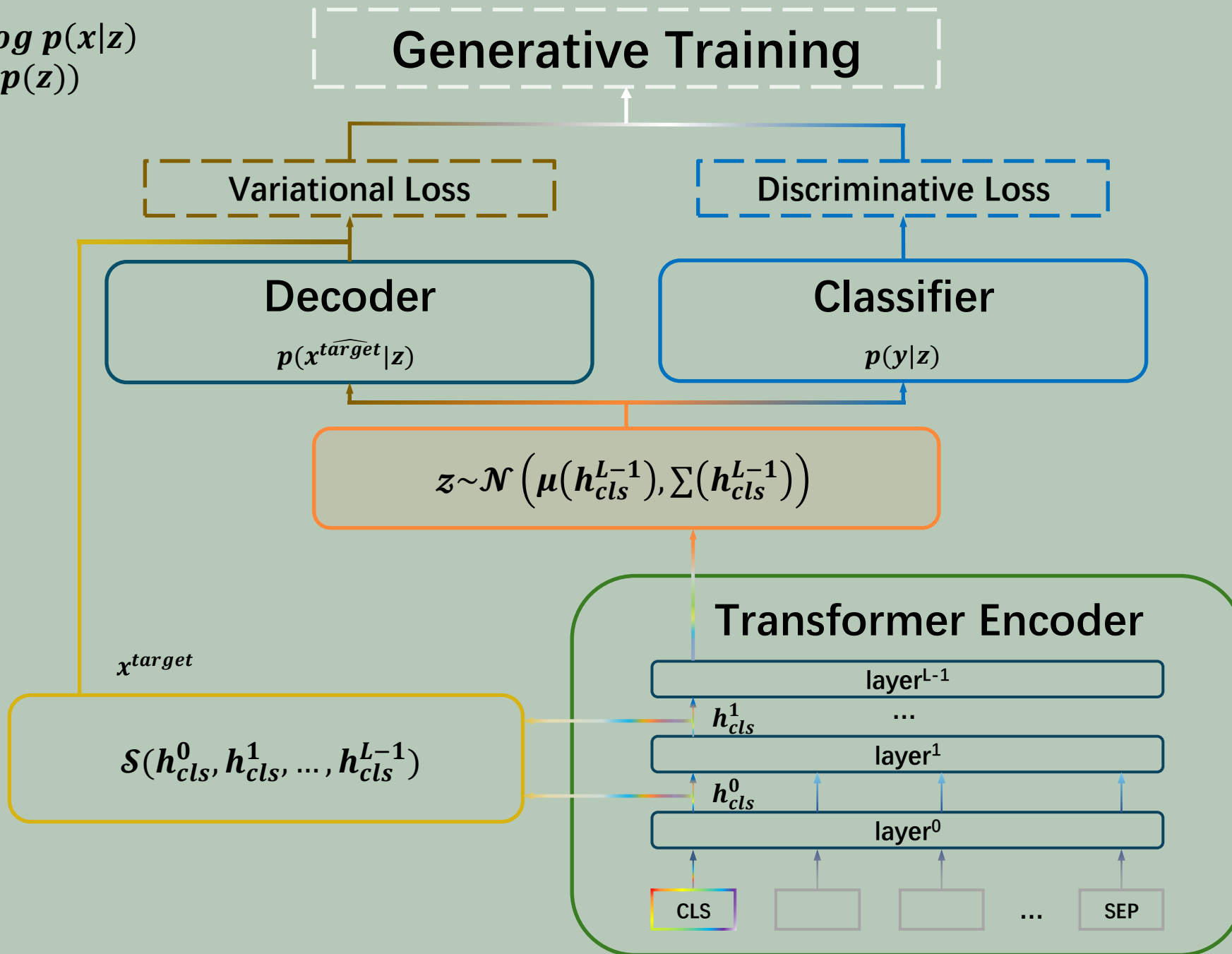


# Methodology

$$p(\mathbf{y}|\mathbf{x}) \rightarrow p(\mathbf{y}, \mathbf{x}) = p(\mathbf{y}|\mathbf{x})p(\mathbf{x})$$

$$\text{ELBO: } p(\mathbf{y}, \mathbf{x}) > \log p(\mathbf{y}|\mathbf{z}) + \log p(\mathbf{x}|\mathbf{z}) - D_{KL}(q(\mathbf{z}|\mathbf{x})|p(\mathbf{z}))$$

$$\log p(y|z) + \log p(x|z) - D_{KL}(q(z|x)|p(z))$$



# Results

# Evaluation

## Metrics:

- Threshold-related: Accuracy, F1-score
- Threshold-free: AUROC, FPR@95, AUPR

## Benchmarks: diverse NLP datasets from natural language understanding tasks

- Sentiment analysis: IMDB, SST-2.
- Topic classification: 20 Newsgroups, Trect-10
- Machine translation: English-German WMT16\ Multi30K
- Natural language inference: RTE, MNLI

## Examples:

Taking IMDB as ID, then 20 Newsgroups, Trect-10... MNLI will be OOD

# Encoder-based Transformer

Methods	SST-2			IMDB		
	AUROC $\uparrow$	FAR@95 $\downarrow$	AUPR $\uparrow$	AUROC $\uparrow$	FAR@95 $\downarrow$	AUPR $\uparrow$
MSP	89.85	66.20	86.40	94.30	41.90	98.80
MSP <sub>Contrast</sub>	85.04	63.42	69.34	94.51	44.69	98.89
<b>MSP<sub>VI</sub></b>	<b>92.85</b>	<b>51.58</b>	<b>89.72</b>	<b>95.95</b>	<b>28.03</b>	<b>99.12</b>
Maha	97.98	11.50	97.30	99.67	0.70	99.95
Maha <sub>Contrast</sub>	<b>99.42</b>	<b>2.98</b>	<b>98.73</b>	<b>99.89</b>	<b>0.05</b>	<b>99.97</b>
<b>Maha<sub>VI</sub></b>	99.33	3.62	98.52	<b>99.90</b>	0.21	<b>99.97</b>
Cosine	95.65	22.65	94.68	99.50	1.53	99.88
Cosine <sub>Contrast</sub>	98.38	8.64	96.36	<b>99.87</b>	1.93	<b>99.96</b>
<b>Cosine<sub>VI</sub></b>	<b>98.87</b>	<b>6.62</b>	<b>98.06</b>	99.57	<b>1.43</b>	99.88
Energy	89.80	67.00	86.53	93.30	56.70	98.63
Energy <sub>Contrast</sub>	84.93	63.16	69.29	94.44	44.46	98.86
<b>Energy<sub>VI</sub></b>	<b>92.79</b>	<b>51.25</b>	<b>89.26</b>	<b>96.05</b>	<b>27.97</b>	<b>99.12</b>

Methods	TREC-10			20NG		
	AUROC $\uparrow$	FAR@95 $\downarrow$	AUPR $\uparrow$	AUROC $\uparrow$	FAR@95 $\downarrow$	AUPR $\uparrow$
MSP	97.94	8.43	89.26	<b>93.89</b>	30.49	<b>87.39</b>
MSP <sub>Contrast</sub>	98.43	4.06	<b>91.19</b>	93.19	28.00	83.17
<b>MSP<sub>VI</sub></b>	<b>98.91</b>	<b>2.77</b>	90.39	93.29	<b>25.61</b>	80.09
Maha	98.99	4.87	95.11	98.39	7.77	95.91
Maha <sub>Contrast</sub>	<b>99.57</b>	0.97	<b>98.59</b>	98.78	5.89	97.29
<b>Maha<sub>VI</sub></b>	99.46	<b>0.79</b>	97.67	<b>99.80</b>	<b>0.61</b>	<b>98.93</b>
Cosine	98.89	3.96	94.54	97.73	10.84	88.71
Cosine <sub>Contrast</sub>	99.14	1.42	93.34	98.03	8.86	95.27
<b>Cosine<sub>VI</sub></b>	<b>99.36</b>	<b>1.19</b>	<b>96.09</b>	<b>99.39</b>	<b>2.92</b>	<b>97.19</b>
Energy	97.19	10.07	82.16	95.76	17.93	<b>88.71</b>
Energy <sub>Contrast</sub>	98.45	4.73	<b>91.18</b>	<b>96.04</b>	<b>15.70</b>	88.62
<b>Energy<sub>VI</sub></b>	<b>99.21</b>	<b>2.84</b>	90.84	94.34	17.04	79.67

Average	AUROC $\uparrow$			FAR@95 $\downarrow$			AUPR $\uparrow$					
avg. (MSP / Maha / Cosine / Energy)	94.00	98.78	97.94	94.01	36.76	6.21	9.75	37.93	<b>90.46</b>	97.07	94.45	89.01
avg-Contrast (MSP / Maha / Cosine / Energy)	92.79	99.17	98.86	93.47	35.04	3.93	5.21	32.01	85.65	97.43	96.23	86.99
<b>avg-<sub>VI</sub></b> (MSP / Maha / Cosine / Energy)	<b>95.25</b>	<b>99.62</b>	<b>99.30</b>	<b>95.60</b>	<b>27.00</b>	<b>1.31</b>	<b>3.04</b>	<b>24.78</b>	89.83	<b>98.77</b>	<b>97.81</b>	<b>89.72</b>

Table 1: Main results of our proposed framework. MSP, Maha, Energy, and Cosine are baseline methods trained with the discriminative loss, while each corresponding method with the *VI* subscript denotes the model trained with our *VI* framework. The *Contrast* subscript denotes the method proposed by Zhou et al. (2021b). The best result is marked in bold. At the bottom row, averaged results across four ID datasets are included. All the reported results are presented in percentage values.

# decoder-based Transformer

Methods	SST-2			TREC-10		
	AUROC $\uparrow$	FAR@95 $\downarrow$	AUPR $\uparrow$	AUROC $\uparrow$	FAR@95 $\downarrow$	AUPR $\uparrow$
MSP	78.22	70.40	74.34	<b>99.29</b>	<b>0.64</b>	<b>98.79</b>
<b>MSP<sub>VI</sub></b>	<b>83.05</b>	<b>65.32</b>	<b>69.54</b>	98.69	5.69	94.48
Maha	46.08	84.41	47.34	84.41	71.94	64.13
<b>Maha<sub>VI</sub></b>	<b>95.41</b>	<b>12.47</b>	<b>85.38</b>	<b>99.96</b>	<b>0.12</b>	<b>99.39</b>
Cosine	89.24	26.52	79.47	98.24	11.06	93.05
<b>Cosine<sub>VI</sub></b>	<b>95.95</b>	<b>7.21</b>	<b>85.70</b>	<b>100</b>	<b>0</b>	<b>99.92</b>
Energy	70.13	67.78	64.87	<b>99.86</b>	<b>0.14</b>	<b>99.27</b>
<b>Energy<sub>VI</sub></b>	<b>80.16</b>	<b>66.90</b>	<b>65.76</b>	98.64	7.40	92.44

Table 2: Results of our proposed framework on LLaMA-2-7B, fine-tuned with a classification head. The best result is marked in bold. All the reported results are presented in percentage values.

Thank you very much!