

Unified Entropy Optimization for Open-Set Test-Time Adaptation

Supplementary Material

6. Pseudo Code

For a better understanding of our proposed methods, we summarize UniEnt and UniEnt+ as Algorithm 1 and Algorithm 2, respectively.

Algorithm 1: UniEnt

Input: Source model f_{θ_0} pre-trained on the source domain dataset, testing samples
 $\mathcal{B}_t = \{\mathbf{x}\}, t = 1, \dots, T$.
for $t \leftarrow 1$ **to** T **do**
 for $\mathbf{x} \in \mathcal{B}_t$ **do**
 | Compute csOOD score for each testing sample via Eq. (3);
 end
 Obtain $\pi(x)$ via the EM algorithm;
 Split \mathcal{B}_t into $\mathcal{B}_{t,\text{csID}}$ and $\mathcal{B}_{t,\text{csOOD}}$ via Eq. (5);
 Update model via Eq. (8);
end
Output: The predictions $\arg \max_c f_{\theta_t}(\mathbf{x})$ for all $\mathbf{x} \in \mathcal{B}_t, t = 1, \dots, T$.

Algorithm 2: UniEnt+

Input: Source model f_{θ_0} pre-trained on the source domain dataset, testing samples
 $\mathcal{B}_t = \{\mathbf{x}\}, t = 1, \dots, T$.
for $t \leftarrow 1$ **to** T **do**
 for $\mathbf{x} \in \mathcal{B}_t$ **do**
 | Compute csOOD score for each testing sample via Eq. (3);
 end
 Obtain $\pi(x)$ via the EM algorithm;
 Update model via Eq. (9);
end
Output: The predictions $\arg \max_c f_{\theta_t}(\mathbf{x})$ for all $\mathbf{x} \in \mathcal{B}_t, t = 1, \dots, T$.

7. More Analysis

Scalability of large-scale datasets. To demonstrate that our methods can be used for large-scale datasets, we conduct experiments on ImageNet-C [14]. Specifically, we use ResNet-50 [13] pre-trained with AugMix [17], the weights of which can be obtained from RobustBench [7]. For optimization, we use the SGD optimizer [38] with the learning rate of 0.00025 and the batch size of 64. We apply

Method	ImageNet-C			
	Acc↑	AUROC↑	FPR@TPR95↓	OSCR↑
Source [54]	28.21	49.63	94.74	19.81
BN Adapt [33]	43.57	55.89	93.39	30.42
CoTTA [46]	47.67	55.58	94.51	33.80
TENT [44]	45.82	51.34	96.47	30.33
+ UniEnt	47.53 (+1.71)	56.33 (+4.99)	<u>95.21</u> (-1.26)	34.42 (+4.09)
+ UniEnt+	46.87 (+1.05)	<u>55.80</u> (+4.52)	95.10 (-1.37)	<u>33.73</u> (+3.40)
EATA [35]	<u>51.40</u>	53.10	95.18	34.87
+ UniEnt	49.60 (-1.80)	<u>58.29</u> (+5.19)	<u>93.63</u> (-1.55)	<u>36.28</u> (+1.41)
+ UniEnt+	51.57 (+0.17)	59.45 (+6.35)	93.60 (-1.58)	38.27 (+3.40)
OSTTA [27]	<u>47.91</u>	52.93	96.15	32.77
+ UniEnt	47.92 (+0.01)	56.02 (+3.09)	<u>95.23</u> (-0.92)	34.47 (+1.70)
+ UniEnt+	47.47 (-0.44)	<u>55.67</u> (+2.74)	95.16 (-0.99)	<u>34.03</u> (+1.26)

Table 8. Results of different methods on ImageNet-C. ↑ indicates that larger values are better, and vice versa. All values are percentages. The **bold** values indicate the best results, and the underlined values indicate the second best results. The values in parentheses indicate the improvements of our methods over the baseline methods.

common corruptions and perturbations to ImageNet-O [18] through the official code of [14] to construct ImageNet-O-C as csOOD data. From Table 8, we can see that UniEnt and UniEnt+ consistently improve the performance of the existing baseline methods in the open-set setting.

Scalability of model architecture. Recently, Vision Transformer (ViT) [10] has demonstrated better performance than Convolutional Neural Network (CNN), we also perform experiments with ViT backbone on ImageNet-C. Specifically, we use DeiT-Base [41] designed in [40], which proposes many techniques in the training phase to improve the robustness of the model to common corruptions. The pre-trained weights are also available from RobustBench. We update the affine parameters of the model’s layer normalization. Table 9 shows that our approaches are compatible with ViT.

Performance under long-term open-set test-time adaptation. Models deployed in real-world scenarios are exposed to test samples for long periods and need to make reliable predictions at any time. Recent work [27, 37] points out that most existing TTA methods perform poorly in long-term settings, even worse than non-updating models. Following [27], we simulate long-term TTA by repeating adaptation for 10 rounds. During adaptation, the domain changes continuously and the model is never reset. The results are summarized in Table 10. We observe that in most cases the performance degradation of our methods is very slight compared to the baseline methods.

Effects of learning rate and batch size. We explore the

Method	ResNet-50				DeiT Base			
	Acc↑	AUROC↑	FPR@TPR95↓	OSCR↑	Acc↑	AUROC↑	FPR@TPR95↓	OSCR↑
Source [54]	28.21	49.63	94.74	19.81	56.59	56.01	91.55	36.13
CoTTA [46]	47.67	55.58	94.51	33.80	60.73	53.51	93.14	37.33
TENT [44]	45.82	51.34	96.47	30.33	62.85	59.51	93.47	43.52
+ UniEnt	47.53 (+1.71)	56.33 (+4.99)	<u>95.21</u> (-1.26)	34.42 (+4.09)	<u>58.81</u> (-4.04)	67.10 (+7.59)	<u>90.90</u> (-2.57)	47.40 (+3.88)
+ UniEnt+	<u>46.87</u> (+1.05)	<u>55.86</u> (+4.52)	95.10 (-1.37)	<u>33.73</u> (+3.40)	<u>58.40</u> (-4.45)	<u>66.69</u> (+7.18)	90.43 (-3.04)	<u>46.74</u> (+3.22)
EATA [35]	<u>51.40</u>	53.10	95.18	34.87	65.38	57.95	92.92	44.29
+ UniEnt	49.60 (-1.80)	<u>58.29</u> (+5.19)	<u>93.63</u> (-1.55)	<u>36.28</u> (+1.41)	59.36 (-6.02)	67.22 (+9.27)	<u>91.63</u> (-1.29)	<u>48.23</u> (+3.94)
+ UniEnt+	51.57 (+0.17)	59.45 (+6.35)	93.60 (-1.58)	38.27 (+3.40)	<u>61.50</u> (-3.88)	<u>66.96</u> (+9.01)	89.99 (-2.93)	48.79 (+4.50)
OSTTA [27]	47.91	52.93	96.15	32.77	60.19	60.69	92.42	43.19
+ UniEnt	47.92 (+0.01)	56.02 (+3.09)	<u>95.23</u> (-0.92)	34.47 (+1.70)	<u>58.73</u> (-1.46)	67.62 (+6.93)	<u>90.51</u> (-1.91)	47.64 (+4.45)
+ UniEnt+	47.47 (-0.44)	<u>55.67</u> (+2.74)	95.16 (-0.99)	<u>34.03</u> (+1.26)	58.72 (-1.47)	<u>67.28</u> (+6.59)	90.02 (-2.40)	<u>47.32</u> (+4.13)

Table 9. Results of different methods on ImageNet-C using diverse architectures.

Method	CIFAR-10-C				CIFAR-100-C				Average			
	Acc↑	AUROC↑	FPR@TPR95↓	OSCR↑	Acc↑	AUROC↑	FPR@TPR95↓	OSCR↑	Acc↑	AUROC↑	FPR@TPR95↓	OSCR↑
Source [54]	81.73	77.89	79.45	68.44	53.25	60.55	94.98	39.87	67.49	69.22	87.22	54.16
BN Adapt [33]	84.20	80.40	76.84	72.13	57.16	72.45	84.29	47.10	70.68	76.43	80.57	59.62
CoTTA [46]	85.77	85.89	72.40	77.26	56.46	77.04	80.96	48.95	71.12	81.47	76.68	63.11
TENT [44]	79.38	65.39	95.94	56.73	54.74	65.00	94.79	42.24	67.06	65.20	95.37	49.49
+ UniEnt	84.31 (+4.93)	92.28 (+26.89)	<u>36.74</u> (-59.20)	80.32 (+23.59)	59.07 (+4.33)	<u>89.28</u> (+24.28)	<u>51.14</u> (+43.65)	56.26 (+14.02)	71.69 (+4.63)	<u>90.78</u> (+25.59)	<u>43.94</u> (-51.43)	68.29 (+18.81)
+ UniEnt+	<u>84.03</u> (+4.65)	93.18 (+27.79)	<u>32.74</u> (-63.20)	80.62 (+23.89)	<u>58.58</u> (+3.84)	91.39 (+26.39)	41.09 (-53.70)	56.36 (+14.12)	<u>71.31</u> (+4.25)	92.29 (+27.09)	<u>36.92</u> (-58.45)	68.49 (+19.01)
EATA [35]	80.92	84.32	71.66	72.63	60.63	88.64	50.18	57.24	70.78	86.48	60.92	64.94
+ UniEnt	<u>84.31</u> (+4.39)	97.15 (+12.83)	13.25 (-58.41)	<u>82.99</u> (+10.36)	<u>59.75</u> (-0.88)	<u>93.42</u> (+4.78)	<u>30.36</u> (-19.82)	57.99 (+0.75)	<u>72.03</u> (+1.26)	95.29 (+8.81)	<u>21.81</u> (-39.12)	70.49 (+55.55)
+ UniEnt+	85.18 (+4.26)	96.97 (+12.65)	14.28 (-57.38)	83.67 (+11.04)	59.71 (0.92)	94.23 (+5.59)	26.87 (-23.31)	58.19 (+0.95)	<u>72.45</u> (+1.67)	<u>95.60</u> (+9.12)	20.58 (-40.35)	70.93 (+6.00)
OSTTA [27]	84.44	72.74	77.02	65.17	60.03	75.37	82.75	51.35	72.24	74.06	79.89	58.26
+ UniEnt	82.46 (-1.98)	<u>96.20</u> (+23.46)	<u>16.37</u> (-60.65)	80.51 (+15.34)	58.69 (-1.34)	<u>94.84</u> (+19.47)	<u>22.95</u> (-59.80)	<u>57.28</u> (+5.93)	70.58 (-1.66)	<u>95.52</u> (+21.47)	<u>19.66</u> (-60.23)	68.90 (+10.64)
+ UniEnt+	84.30 (-0.14)	97.38 (+24.64)	11.56 (-65.46)	82.91 (+17.74)	58.93 (-1.10)	95.42 (+20.05)	20.59 (-62.16)	57.69 (+6.34)	<u>71.62</u> (-0.62)	96.40 (+22.35)	16.08 (-63.81)	70.30 (+12.04)

(a) Results after 1 round of adaptation (*i.e.*, short-term test-time adaptation).

Method	CIFAR-10-C				CIFAR-100-C				Average			
	Acc↑	AUROC↑	FPR@TPR95↓	OSCR↑	Acc↑	AUROC↑	FPR@TPR95↓	OSCR↑	Acc↑	AUROC↑	FPR@TPR95↓	OSCR↑
Source [54]	81.73	77.89	79.45	68.44	53.25	60.55	94.98	39.87	67.49	69.22	87.22	54.16
BN Adapt [33]	84.20	80.40	76.84	72.13	57.16	72.45	84.28	47.09	70.68	76.43	80.57	59.62
CoTTA [46]	35.90	47.27	97.52	19.95	13.34	48.34	91.61	8.19	24.62	47.81	94.57	14.07
TENT [44]	32.61	60.86	93.24	20.86	37.49	53.73	95.07	25.02	35.05	57.30	94.16	22.94
+ UniEnt	<u>84.07</u> (+51.46)	88.53 (+27.67)	51.48 (-41.76)	77.87 (+57.01)	57.93 (+20.44)	<u>90.62</u> (+36.89)	46.18 (-48.89)	55.67 (+30.65)	<u>71.00</u> (+35.95)	89.58 (+32.28)	48.83 (-45.33)	66.77 (+43.83)
+ UniEnt+	84.17 (+51.56)	<u>88.21</u> (+27.35)	<u>52.57</u> (-40.67)	<u>77.75</u> (+56.89)	<u>57.92</u> (+20.43)	90.63 (+36.90)	45.10 (+49.97)	<u>55.59</u> (+30.57)	71.05 (+36.00)	<u>89.42</u> (+32.13)	<u>48.84</u> (+45.32)	<u>66.67</u> (+43.73)
EATA [35]	40.94	64.52	88.41	29.07	48.75	73.26	80.83	41.27	44.85	68.89	84.62	35.17
+ UniEnt	81.22 (+40.28)	<u>91.05</u> (+26.53)	<u>30.59</u> (-57.82)	<u>76.42</u> (+47.35)	<u>57.07</u> (+8.32)	98.59 (+25.33)	5.85 (-74.98)	<u>56.70</u> (+15.43)	<u>69.15</u> (+24.30)	<u>94.82</u> (+25.93)	18.22 (-66.40)	<u>66.36</u> (+31.39)
+ UniEnt+	80.41 (+39.47)	92.49 (+27.97)	30.00 (-58.41)	77.00 (+47.93)	58.02 (+9.27)	<u>98.05</u> (+24.79)	<u>7.92</u> (-72.91)	57.47 (+16.20)	<u>69.22</u> (+24.37)	95.27 (+26.38)	<u>18.96</u> (-65.66)	67.24 (+32.07)
OSTTA [27]	83.83	71.93	76.12	63.90	<u>57.39</u>	75.46	82.47	49.61	70.61	73.70	79.30	56.76
+ UniEnt	80.74 (-3.09)	<u>88.94</u> (+17.01)	<u>35.66</u> (-40.46)	<u>74.52</u> (+10.62)	56.13 (-1.26)	<u>95.20</u> (+19.74)	<u>21.15</u> (-61.32)	<u>54.89</u> (+5.28)	68.44 (-2.18)	<u>92.07</u> (+18.38)	<u>28.41</u> (-50.89)	<u>64.71</u> (+7.95)
+ UniEnt+	82.42 (-1.41)	90.15 (+18.22)	<u>31.18</u> (-44.94)	76.46 (+12.56)	<u>57.45</u> (+0.06)	<u>95.91</u> (+20.45)	<u>17.33</u> (-65.14)	<u>56.32</u> (+6.71)	<u>69.94</u> (-0.67)	93.03 (+19.34)	<u>24.26</u> (-55.04)	66.39 (+9.64)

(b) Results after 10 rounds of adaptation (*i.e.*, long-term test-time adaptation).

Table 10. Results of different methods on CIFAR benchmarks.

impact of learning rate and batch size on our approaches in Table 11. A learning rate that is too large or too small can hurt performance, while a larger batch size results in better performance. Compared to TENT [44] and EATA [35], our methods are more robust to learning rate and batch size. Nonetheless, our methods share the same limitation as the baseline methods: they rely on a large batch size to estimate the distribution accurately. Moreover, we observe that OSTTA [27] is less sensitive to learning rate and batch size.

Effects of OOD score. We use the energy score [32] to measure the model’s detection performance on csOOD data. From Table 12, we can make two observations. First, our methods consistently improve the performance using different OOD scores. Second, compared with MSP [15], using

Max Logit [19] and Energy yields better detection performance.

Method	Learning rate				Δ
	0.005	0.001	0.0005	0.0001	
Source [54]	39.87	39.87	39.87	39.87	0.00
BN Adapt [33]	47.10	47.10	47.10	47.10	0.00
TENT [44]	10.60	42.24	42.38	48.36	37.76
+ UniEnt	<u>53.82</u> (+43.22)	<u>56.20</u> (+13.96)	<u>56.06</u> (+13.68)	<u>54.51</u> (+6.15)	2.38
+ UniEnt+	54.44 (+43.84)	56.36 (+14.12)	56.27 (+13.89)	54.65 (+6.29)	1.92
EATA [35]	40.96	57.00	56.91	<u>53.60</u>	16.04
+ UniEnt	49.36 (+8.40)	<u>57.76</u> (+0.76)	<u>57.10</u> (+0.19)	53.63 (+0.03)	8.40
+ UniEnt+	<u>49.05</u> (+8.09)	58.07 (+1.07)	57.39 (+0.48)	<u>53.40</u> (-0.20)	9.02
OSTTA [27]	49.43	51.35	51.98	52.37	2.94
+ UniEnt	<u>51.41</u> (+1.98)	<u>56.93</u> (+5.58)	<u>57.22</u> (+5.24)	<u>55.58</u> (+3.21)	5.81
+ UniEnt+	53.39 (+3.96)	57.69 (+6.34)	57.68 (+5.70)	56.06 (+3.69)	4.30

(a) OSCR w.r.t. learning rate

Method	Batch size				Δ
	64	32	16	8	
Source [54]	39.87	39.87	39.87	39.87	0.00
BN Adapt [33]	46.38	45.25	42.94	38.61	7.77
TENT [44]	33.27	8.10	2.51	0.95	32.32
+ UniEnt	55.17 (+21.90)	<u>53.05</u> (+44.95)	<u>48.87</u> (+46.36)	31.47 (+30.52)	23.70
+ UniEnt+	<u>55.17</u> (+21.90)	53.13 (+45.03)	49.27 (+46.76)	<u>28.35</u> (+27.40)	26.82
EATA [35]	53.09	47.78	40.57	31.57	21.52
+ UniEnt	57.08 (+3.99)	54.52 (+6.74)	50.71 (+10.14)	43.89 (+12.32)	13.19
+ UniEnt+	<u>56.79</u> (+3.70)	<u>54.29</u> (+6.51)	<u>50.49</u> (+9.92)	<u>43.17</u> (+11.60)	13.62
OSTTA [27]	50.35	48.82	46.07	<u>39.75</u>	10.60
+ UniEnt	<u>54.54</u> (+4.19)	<u>50.49</u> (+1.67)	44.97 (-1.10)	36.72 (-3.03)	17.82
+ UniEnt+	55.76 (+5.41)	52.66 (+3.84)	47.94 (+1.87)	41.45 (+1.70)	14.31

(b) OSCR w.r.t. batch size

Table 11. OSCR of different methods on CIFAR-100-C with diverse learning rates and batch sizes. Δ is the difference between the largest and smallest values. Smaller Δ values represent better robustness.

Method	OOD score			Δ
	MSP [15]	Max Logit [19]	Energy [32]	
Source [54]	39.65	40.24	39.87	0.59
BN Adapt [33]	48.75	48.04	47.10	1.65
CoTTA [46]	49.44	49.73	48.99	0.74
TENT [44]	36.86	41.79	42.24	5.38
+ UniEnt	55.42 (+18.56)	<u>56.20</u> (+14.41)	<u>56.26</u> (+14.02)	0.84
+ UniEnt+	<u>55.24</u> (+18.38)	56.31 (+14.52)	56.36 (+14.12)	1.12
EATA [35]	55.20	57.52	57.55	2.35
+ UniEnt	<u>56.94</u> (+1.74)	<u>57.88</u> (+0.36)	<u>57.87</u> (+0.32)	0.94
+ UniEnt+	57.37 (+2.17)	58.33 (+0.81)	58.33 (+0.78)	0.96
OSTTA [27]	49.14	51.42	51.35	2.28
+ UniEnt	<u>56.52</u> (+7.38)	<u>57.23</u> (+5.81)	<u>57.25</u> (+5.90)	0.73
+ UniEnt+	57.12 (+7.98)	57.69 (+6.27)	57.69 (+6.34)	0.57

Table 12. OSCR of different methods on CIFAR-100-C using diverse OOD scores.