

1 We thank the reviewers for their insightful feedback. We are encouraged that they find our motivation and idea to
2 be appealing (R1), theoretically insightful (R2), new (R2) and clearly presented (R2, R4). We are equally glad they
3 found our results strong (R1), sound (R2), promising and thorough (R3), appropriately positioned w.r.t. prior work
4 (R1, R3), and outperforming the state-of-the-art (R2, R4). Moreover, we are pleased that R2 and R4 appreciated our
5 mathematical insights on why the energy score is better than the softmax score for out-of-distribution detection. We
6 address reviewer comments below and will incorporate all feedback.

7 **[R1, R2] Can you evaluate on additional in-distribution dataset?** We have run experiments using SVHN as in-
8 distribution dataset. Using energy score consistently outperforms softmax score by **8.56%** (FPR at 95% TPR), which
9 suggests its general applicability beyond CIFAR datasets. We will definitely include these results in the final version.

10 **[R2, R3] How to distinguish methodological differences w.r.t. prior work?** As discussed in L240–L248, there are
11 several major distinctions. We have moved the expanded discussions to Section 3 for clarity, as R3 suggested. To recap:

- 12 • The idea of using the energy score for OOD detection is novel and theoretically motivated. JEM [10] used
13 likelihood score $\log p(\mathbf{x})$ from **generative** modeling or softmax score; we instead derive the energy score from a
14 **pure discriminative** model. We showed both mathematical insights (L103–L117) and empirical success that energy
15 is better than softmax score—which was not presented in OE [14] or JEM [10].
- 16 • JEM’s optimization is intractable as it requires estimating the normalized densities w.r.t. all inputs \mathbf{x} . Ours is easy to
17 optimize using standard SGD, and doesn’t involve proper normalization.
- 18 • JEM only uses in-distribution data. Our training objective (6) leverages auxiliary outlier data and explicitly enlarges
19 the energy gap, which is also different from OE.

20 **[R1] Is it realistic to assume we have OOD data to optimize the objective? Does one source of OOD data also
21 help another?** Thank you for the question! In fact, we are not the first to use auxiliary data for OOD detection. Previous
22 work (Outlier Exposure [14]) showed one can leverage a large, diverse auxiliary outlier dataset for better estimation on
23 the decision boundary for OOD detection. Since the auxiliary data (i.e., 80 million TinyImages) is sufficiently diverse,
24 this fine-tuning process helps generalize to different test-time OOD data, as evidenced in Table 1.

25 **[R1] Is there a train/test split of the OOD data for fine-tuning?** Following common practice, we used 80M
26 TinyImages as auxiliary OOD training data; we tested on six OOD datasets (L143) that are different from training.

27 **[R2] Why are the $\log p(x)$ results in JEM [10] much worse?** Very insightful question! While the energy score can
28 be interpreted from the likelihood view, the optimization processes of JEM vs. ours are very different. We believe
29 the difficulty in optimization led to the worse results we see in JEM. Specifically, JEM’s objective estimates the
30 joint distribution $\log p_\theta(\mathbf{x}, y) = \log p_\theta(\mathbf{x}) + \log p_\theta(y|\mathbf{x})$. The first term maximizes the likelihood from a generative
31 modeling perspective, which requires estimating normalized densities. As pointed out by Grathwohl et al. [10] in the
32 limitations (Sec. 6), “*Energy based models can be very challenging to work with. Since normalized likelihoods cannot
33 be computed, it can be hard to verify that learning is taking place at all . . . Furthermore, the gradient estimators we
34 use to train JEM are quite unstable and are prone to diverging if the sampling and optimization parameters are not
35 tuned correctly*”. In contrast, our training objective is purely discriminative and can be more easily optimized using
36 standard SGD.

37 **[R2] Why having two margin loss parameters?** One problem with the single margin is that the difference of the
38 energy scores between in- and out-of-distribution is relatively fixed; however, the energy scores can change over the
39 training process. This results in lower accuracy of in-distribution classification (-0.57% for CIFAR-10 and -1.11% for
40 CIFAR-100) because the training is less stable due to the drift of energy scores from batch to batch. Instead, using two
41 margin hyper-parameters results in overall good performance on both in-distribution classification and OOD detection.

42 **[R3] Does the method work on more complex (CIFAR/SVHN) → simpler dataset (MNIST)?** Yes. On the CIFAR-
43 10-pretrained WideResNet, using the energy score yields **14.50%** FPR when evaluated against MNIST as OOD data,
44 outperforming 51.73% using the softmax score.

45 **[R3, R4] Other architectures/toy example/more comparison between energy and JEM:** Yes, consistent improve-
46 ment was also observed on other architectures such as AllConv and DenseNet. We will include these results and the
47 suggested toy example, more comparisons with JEM in the final version.

48 **[R3] Were these 10 runs done over random hyper-parameter configurations or 10 runs with the best selected
49 hyperparameters’ values?** The latter. We will make sure to report variances in the final version.

50 **[R4] What if you use self-normalized output, replacing logits with log score?** We found the suggestion of self-
51 normalization very interesting, and investigated further. Mathematically, the log score is equivalent to logit shifted by the
52 free energy, which is *sample specific* and therefore might change the overall score distribution. This shifting in actuality
53 leads to the free energy score (on self-normalized output) collapsing to 0, i.e., $\text{LogSumExp}(\log \text{softmax}_i) = 0$, and
54 did not work well. We believe the normalizing factor needs to be *sample independent* to ensure $E(\mathbf{x}, f)$ is proportional
55 to the density. We are excited to explore this as part of future work.

56 **[R2, R3, R4] Typo/Writing clarity:** All fixed. We thank the reviewers for the careful read and helpful suggestions!