1  We thank the reviewers for the valuable feedback. We will reflect minor errors instantly and try suggestions in the
2  future. We want to address concerns and to clarify misconceptions. In this rebuttal, CAP denotes Wong et al.
3  **[R2 R3] Under $\ell_2$-norm, IBP does not work well enough.** In IBP, the authors didn't mention the certifiable training
4  for $\ell_2$-norm. In CROWN-IBP, however, the authors explained that IBP can be applied to $\ell_2$-norm. We run the code for
5  IBP provided by the authors of CROWN-IBP on CIFAR-10 under $\ell_2$-setting with $\ell_2$-norm $36/255$, where we considered a
6  wide range of parameters including those suggested by CROWN-IBP and our range of settings described in Section C.
7  IBP achieved the verification (standard) accuracy of 22.6-23.0% (31.3-33.5%) which was inferior to CAP = 50.29%
8  (60.14%) and to LMT = 37.20% (56.49%). It also implies LMT-bound is much tighter than IBP-bound under $\ell_2$-norm.
9  Thus, we compared our method to LMT and CAP rather than IBP under $\ell_2$-case in the main text. **[R2 R3] On the**
10 **novelty.** BCP was carefully designed in a layer-wise manner (Fig1) to obtain the tighter outer bound. One may simply
11 use IBP to get the additional box constraint in (11), but this box constraint is redundant because it is much looser
12 than the $\ell_2$-constraint in (11). BCP is designed to provide a nonredundant box constraint in (11) to tighten the bound.
13 We observed the tightness in Fig2 and Fig3. As a result, BCP outperforms both LMT and IBP in a large margin
14 (>12-28%p) under $\ell_2$-norm. For further intuition behind the layerwise design of BCP, refer to Section B.1. **[R2 R3]**
15 **BCP outperforms the others with a meaningful margin.** In Tab 1, the evaluation results at a single $\epsilon_{eval}$ seem to be
16 a marginal improvement compared to CAP. However, when considering a wide range of $\epsilon_{eval}$ in Fig4 (and in FigS3),
17 BCP outperforms CAP by 3.7-5.6%p in standard accuracy on CIFAR-10. For $\epsilon_{eval} > 36/255$, BCP outperforms CAP in a
18 large margin. For example, when evaluating at $72/255$, BCP (34.2%) defeats CAP (23.9%) by 10.3%p. Moreover, only
19 BCP can achieve a meaningful verification accuracy on Tiny ImageNet, while others cannot. **[R1 R2] Fig2 illustrates**
20 **how BCP can tighten the outer region by introducing the box constraint.** In Fig2 (a)-(c), we can easily visualize
21 the high-dimensional ellipsoid $h_K(\mathbb{B}_2^{(K-1)}) \subset \mathbb{R}^c$ in 2D plane with $\zeta_y$- and $\zeta_{m'}$-axes (line 211-212) by projection.
22 However, a high-dimensional parallelogram $h_K(\mathbb{B}_\infty^{(K-1)})$ for $c > 2$ is hard to visualize in the 2D plane. Thus, we
23 use the red lines in Fig2 to indicate that the projection of the outer region $h_K(\mathbb{B}_2^{(K-1)} \cap \mathbb{B}_\infty^{(K-1)})$ must lie above the
24 red line and inside the ellipsoid. Based on (9), we used the verification boundary ($\zeta_y - \zeta_{m'} \geq \zeta_y^* - \zeta_{m'}^*$), where
25 the red line is obtained from the solution $\zeta_y^*, \zeta_{m'}^*$ of (9) for $\hat{z}(\mathbb{B}(\boldsymbol{x})) = h_K(\mathbb{B}_2^{(K-1)} \cap \mathbb{B}_\infty^{(K-1)})$, and the blue line is
26 for $\hat{z}(\mathbb{B}(\boldsymbol{x})) = h_K(\mathbb{B}_2^{(K-1)})$. Fig2 (d) explicitly illustrates the ellipsoid and the parallelogram with the verification
27 boundary for a toy binary classification ($c = 2$). Fig2 shows typical examples for which the verification succeeds with
28 BCP but fails without BCP (R2-8-3). On comparing the gap between BCP and other baselines, the visualization of the
29 logit space can be inappropriate because it is not possible to directly compare the logits from the models obtained by
30 different certifiable training methods. As verification methods, comparing only verification phase using one trained
31 model can be unfair because the performance highly depends on which method the model is trained with. **[R2 R3] We**
32 **computed the radius $\rho^{(k)}$ with the layerwise Lipschitz constants $L^{(i)}$'s.** We computed $L^{(i)}$ for each $i$-th operation
33 (BCP.py line 254,286), and then computed $\rho^{(k)} = \epsilon \Pi_{i=1}^k L^{(i)}$, multiplying the expansion rate (=Lipshitz constant)
34 through each layer (line 114-118,137-139, Section B.2). We are sorry to make R2 confused. Lipschitz constant usually
35 refer to the global Lipschitz constant (gL) rather than a local one (lL) when not specified. We used the gL for the
36 efficient computation (line 127-128). For the layerwise Lipschitz constant, when the layerwise operation $h_k$ is linear, the
37 maximum eigen-value of weight matrix corresponds to both gL and lL (gL = lL). **[R3] On the stabilization of $\rho^{(K-1)}$.**
38 In FigS1 (top), the Lipschitz constant keeps increasing through standard training. To stabilize the Lipschitz constant,
39 we applied a commonly-used scheduling scheme on $\epsilon$ and $\lambda$ during the BCP training (line 195-200). In FigS1 (bottom),
40 $L^{(-1)} = \Pi_{i=1}^{K-1} L^{(i)}$ is about 10 after training with BCP, and it is multiplied by $\epsilon$ to provide $\rho^{(K-1)} = \epsilon L^{(-1)}$. Thus,
41 $\rho^{(K-1)}$ is about $10\epsilon$. In the code, we initialize $\rho^{(0)} = \epsilon$ (BCP.py line 214) and update the radius by $\rho^{(i+1)} = \rho^{(i)} * L^{(i+1)}$
42 (r = r*p) for each layer as implemented in BCP.py line 236-243,270,300. We also tried BN, but it is not effective
43 to improve robustness. Moreover, there is a paper named "Batch Normalization is a Cause of Adversarial Vulnerability".
44 **[R3] On the implementation of BCP.** Our main focus is on $\ell_2$-certifiable training. In the main text, we describe the
45 BCP algorithm and provide the results for $\ell_2$-case, so, in the code, we set `args.linfty=False` as default. However,
46 BCP can be applied to $\ell_p$-norm for any $p \in (0, \infty]$ (line 154-155). In appendix, we presented the results for $\ell_\infty$-norm
47 in Tab S2, and the implementation is available by setting `args.linfty=True` in the code (BCP.py line 216-217).
48 Therefore, the description of the algorithm, including the $\ell_\infty$-case (line 154-159), is consistent with the implementation.
49 Moreover, for reproducibility, we provided how to run the code in one-line command in README.md and explained
50 details of the hyper-parameters in Section C, including the scheduling on $\epsilon$ and $\lambda$. For reference, in the $\ell_\infty$-case, BCP is
51 not IBP but a generalized version of IBP because BCP still uses the $\ell_2$-bound propagation. **[R3] c and $\mathbf{z}^{(K-1)}$ have the**
52 **same shape.** The notation $\mathbf{W}_i$ is for the $i$-th row vector of the matrix $\mathbf{W}$ (line 144). Let $\mathbf{W}^{(K)} \in \mathbb{R}^{m \times n}$, then $\mathbf{W}_i^{(K)}$
53 is an $n$-dimensional row vector, and $\mathbf{z}^{(K-1)}$ is an $n$-dimensional column vector. Therefore, since $\mathbf{c}^T = \mathbf{W}_1^{(K)} - \mathbf{W}_0^{(K)}$,
54 $\mathbf{c}$ has the same dimension to $\mathbf{z}^{(K-1)}$. **[R1] Trade-off between efficiency and accuracy.** BCP is much faster than CAP
55 (Tab1) as well as achieving the better performance in a wide range of $\epsilon_{val}$ (Fig4). BCP can outperform LMT and IBP in
56 a large margin with affordable computational overhead.