# Security and Privacy of Machine Learning, 2025 Critique: Enhancing Certified Robustness via Block Reflector Orthogonal Layers and Logit Annealing Loss

Shih-Yu Lai
*National Taiwan University*
Taipei, Taiwan
akinesia112@gmail.com

## I. Summary of the paper

The paper proposes two contributions for certifiably robust image classification under $\ell_2$ attacks: (1) a *Block Reflector Orthogonal* (BRO) layer that parameterizes orthogonal mappings as $W = I - 2V(V^\top V)^{-1}V^\top$, and extends this to convolutions by operating in the Fourier domain; (2) a *Logit Annealing* (LA) loss, $L_{\mathrm{LA}}(z,y) = -T(1 - p_t)^\beta \log p_t$ with an offset in the softmax, to downweight easy (large-margin) examples and re-allocate capacity to small-margin ones. Combined into *BRONet*, the method reports state-of-the-art certified accuracy on CIFAR-10/100, Tiny-ImageNet, and ImageNet, along with favorable runtime/memory versus orthogonal baselines (e.g., SOC [1], LOT [2], CPL [3], AOL [4], Cholesky [5]). The paper additionally analyzes why indiscriminate margin maximization (e.g., CE+CR) is ill-suited for Lipschitz models using a Rademacher-complexity argument, and presents ablations (diffusion data, backbones, rank choices).

## II. Strengths

- **Simple, exact orthogonal parameterization.** The block-reflector form avoids iterative schemes (e.g., Newton steps in LOT or exponential series in SOC), removing approximation drift and numerical instabilities while retaining exact orthogonality. This is a clear engineering and conceptual win.
- **Fourier-domain convolution that stays orthogonal.** Mapping multi-channel circular convolution to per-frequency matrix multiplications is a clean way to guarantee 1-Lipschitz behavior layerwise while remaining implementable at scale.
- **A principled loss for capacity-limited models.** The LA loss explicitly addresses limited capacity in Lipschitz networks by annealing gradients of high-confidence (large-margin) points; the margin distribution analyses (median increase, reduced variance/skew) are convincing indicators that LA reallocates effort where it matters.
- **Strong empirical results and fairer comparisons.** The paper reports results both with and without large synthetic diffusion datasets and performs backbone swaps to isolate the effect of the BRO layer, which strengthens the empirical case.
- **Clarity of limitations.** The paper explicitly notes less consistent gains at large $\epsilon$ and the need for hyperparameter tuning for LA; acknowledging these helps readers scope applicability.

## III. Weaknesses / Concerns

- **Orthogonality class expressiveness.** A single BRO layer has a constrained spectrum (eigenvalues in $\{\pm 1\}$ with multiplicities governed by $\mathrm{rank}(V)$). While stacking increases expressiveness, the paper does not deeply analyze whether this restriction systematically biases representations (e.g., toward reflections across low-rank subspaces) and how much depth/rank are required to match alternative orthogonal parameterizations.
- **Circular convolution and padding effects.** The orthogonality guarantee hinges on circular convolution in the Fourier domain and zero-padding choices. The paper mentions slight norm drops after cropping padded borders; however, the potential impact on both certification tightness and feature learning (e.g., edge artifacts, frequency leakage) deserves a more thorough analysis and alternatives (e.g., orthogonalization under valid/"same" boundary conditions).
- **Fairness and scalability trade-offs.** While comparisons aim for fairness, some baselines require reduced depth due to memory (especially FFT-based methods). This can blur whether gains are due to BRO's form or capacity differences. A matched-MACs or matched-latency comparison would strengthen the claim.
- **Loss design baselines.** LA is compared primarily to CE and CE+CR. Missing are other margin-shaping or calibration losses (e.g., focal with temperature/label-smoothing, entropy maximization variants, margin-based CE, or per-sample reweighting via uncertainty) adapted to the Lipschitz setting; these could narrow or contextualize LA's advantage.

- **Norm-specificity.** The method focuses on $\ell_2$ certification; results for $\ell_\infty$ are incidental and not central. The orthogonal design and LA may not transfer directly to tighter $\ell_\infty$ certificates, limiting generality across common robustness benchmarks.

## IV. POTENTIAL IMPROVEMENTS OR EXTENSIONS

- **Broader orthogonal families.** Compose BRO with additional unitary/orthogonal factors (e.g., Householder stacks or Givens flows) to relax eigenvalue structure while keeping exactness; study depth/rank–expressivity curves with controlled budgets.
- **Adaptive LA.** Make $\beta$ or $\xi$ *per-sample adaptive* based on running margin statistics or difficulty estimates, or schedule them across depth (early vs. late layers) to reduce tuning burden and better fit capacity.
- **Task diversity and architectures.** Test BRO/LA beyond image classification (e.g., detection/segmentation) and with modern backbones (ConvNeXt, ViT/MLP-Mixer equivalents made Lipschitz) to assess generality and data-regime behavior.
- **Compute-normalized comparisons.** Provide head-to-head curves at matched training time, GPU-hours, or MACs to isolate algorithmic gains from capacity/runtime trade-offs.

## V. QUESTIONS FOR THE AUTHORS

1) **About the certificate:** When using LLN vs. the classical $\epsilon = \max(0, M_f(x))/(\sqrt{2}L)$, how often do bounds disagree, and by how much? Any systematic cases where one is tighter?
2) **Spectrum constraints:** Given a BRO layer yields a fixed number of $-1$ eigenvalues (equal to $\mathrm{rank}(V)$), do you observe representation collapse or directional bias early in training? Would alternating ranks (e.g., $m/8$, then $m/2$) across depth help?
3) **FFT/circularity:** Have you evaluated boundary-condition variants (e.g., symmetric padding with DCT-based orthogonalization) and their effect on both clean accuracy and certified radii, especially for small images (CIFAR) where wrap-around is stronger?
4) **LA vs. other reweightings:** How does LA compare to focal loss with tuned temperature/label smoothing, confidence penalty, or uncertainty-based reweighting in otherwise identical Lipschitz settings?
5) **Hyperparameters:** Can $\beta$ be scheduled or learned (e.g., via meta-gradients) to avoid hand-tuning across datasets? Any failure modes when $\beta$ is too large (e.g., underfitting hard classes)?
6) **Data augmentation dependency:** To what extent do the SOTA results depend on diffusion data volume/quality? If synthetic data is mismatched (domain shift), does LA still improve margin balance, or does it over-anneal?
7) **Transfer to $\ell_\infty$:** Which components of BRO/LA hinder $\ell_\infty$ certification most (architecture vs. bound vs. loss)? Any promising modifications you tried?
8) **Reproducibility at scale:** For ImageNet runs, how sensitive are results to FFT implementations (precision, library), and do you see variability in orthogonality due to kernel conditioning?
9) **Clarification (from the reading notes).** What is a Lipschitz neural network? What is the lower bound? Why use orthogonal layers?

## REFERENCES

[1] S. Singla, S. Singla, and S. Feizi, "Improved deterministic l2 robustness on cifar-10 and cifar-100," 2022. [Online]. Available: https://arxiv.org/abs/2108.04062

[2] X. Xu, L. Li, and B. Li, "LOT: Layer-wise orthogonal training on improving l2 certified robustness," in *Advances in Neural Information Processing Systems*, A. H. Oh, A. Agarwal, D. Belgrave, and K. Cho, Eds., 2022. [Online]. Available: https://openreview.net/forum?id=ZBlaix34YX

[3] L. Meunier, B. J. Delattre, A. Araujo, and A. Allauzen, "A dynamical system perspective for Lipschitz neural networks," in *Proceedings of the 39th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, K. Chaudhuri, S. Jegelka, L. Song, C. Szepesvari, G. Niu, and S. Sabato, Eds., vol. 162. PMLR, 17–23 Jul 2022, pp. 15 484–15 500. [Online]. Available: https://proceedings.mlr.press/v162/meunier22a.html

[4] L. Prach and C. H. Lampert, "Almost orthogonal layers for efficient lipschitz-constrained training," in *European Conference on Computer Vision*, 2022.

[5] K. Hu, K. Leino, Z. Wang, and M. Fredrikson, "A recipe for improved certifiable robustness," 2024. [Online]. Available: https://arxiv.org/abs/2310.02513