

LOOKSHARP: ATTENTION ENTROPY MINIMIZATION FOR TEST-TIME ADAPTATION

Yash Mali¹, Evan Shelhamer^{1,2}

University of British Columbia, Department of Computer Science¹

Vector Institute²

ymali@mail.ubc.ca

ABSTRACT

Test-time adaptation (TTA) updates models during inference to reduce error on distribution shifts. While entropy minimization over the output distribution has proven effective as a TTA loss, we study using the intermediate distributions computed by transformers in the attention mechanism. We propose *LookSharp*, which minimizes the entropy of CLS-to-patch attention in the final layer as a novel TTA objective, encouraging the model to maintain focused attention on shifted data. We demonstrate that attention entropy minimization improves robustness on ImageNet-C (Hendrycks & Dietterich, 2019). We also show that it is complementary to output entropy minimization and maintains performance on clean data.

1 INTRODUCTION AND RELATED WORK

Deep networks achieve impressive performance on in-distribution data but often fail catastrophically when deployed on data from shifted distributions (Hendrycks & Dietterich, 2019). Recent TTA methods have explored entropy minimization over the output distribution, which encourages the model to make confident predictions at test time. While effective, this approach treats the feature extractor as a black box and ignores internal representations that could guide adaptation. Vision Transformers (ViTs) (Dosovitskiy et al., 2021), which have become the dominant architecture for visual recognition due to their scalability, offer attention distributions over image patches that explicitly capture spatial relationships and feature importance (Fuller et al., 2025).

We harness these attention distributions for TTA, minimizing the entropy of the attention distributions in vision transformers as an unsupervised loss to update the model parameters. As this sharpens the distribution to focus more on fewer tokens, we call our method *LookSharp*. Specifically, we minimize the entropy of the distribution defined by the attention scores of the CLS token for the patch tokens from the last layer’s attention heads. Our approach is motivated by two key observations. First, Figure 1 (b) shows that accuracy drops sharply if the attention entropy is too diffuse. Second, Modern ViTs like DINOv3 (Siméoni et al., 2025) learn interpretable and object-centric attention maps through internet-scale self-supervised training.

We demonstrate our method for adaptation to corruptions on ImageNet-C in the batch episodic setting. That is, the model updates and then resets on each batch. We also show that combining attention entropy and output entropy leads to further improvement.

Entropy Minimization for Adaptation. Test-time adaptation often relies on entropy minimization. Tent (Wang et al., 2021) updates normalization layer statistics and parameters to minimize output entropy. MEMO (Zhang et al., 2022) extends this by using test-time augmentation to create a batch from a single sample and updates all parameters episodically using the same loss as Tent. Other works like SAR (Niu et al., 2023) and EATA (ETA) (Niu et al., 2022) use output entropy combined with sharpness-aware minimization, data filtering, and anchoring to the source model using regularization of the parameters.

Attention for Adaptation. There has been less use of attention for updates. **Attent** (Kojima et al., 2023) aligns test-time attention statistics with stored source statistics. Unlike **Attent**, our method is

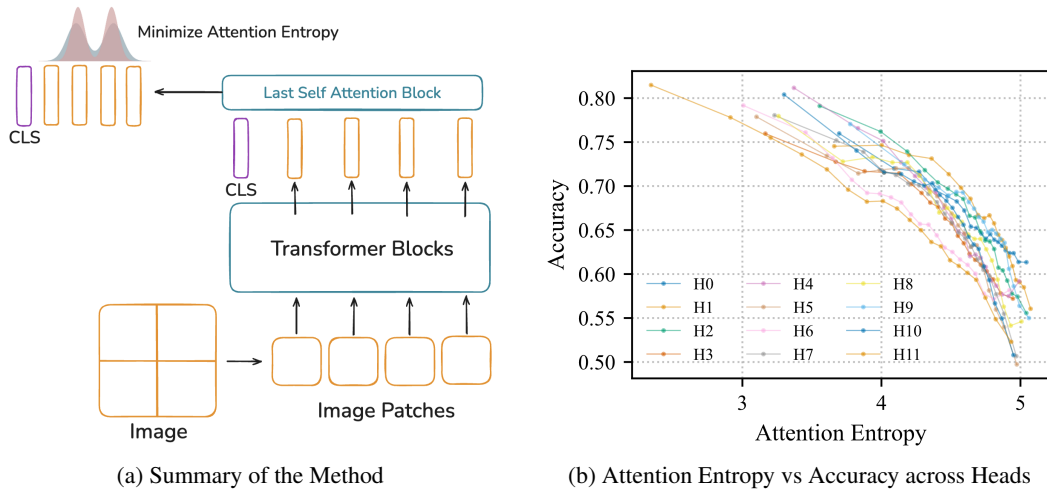


Figure 1: **Left** (a) Our method: attention entropy minimization. We minimize the entropy of the attention distribution from CLS to patch tokens at test time. We combine this with output entropy minimization for best results. **Right** (b) Visualization of attention entropy and accuracy. The entropy of final-layer CLS to patch token attention and the accuracy on shifted data are shown across all heads on a 10% sample of ImageNet-C for the unadapted DINOv3-Base model. Higher entropy (x-axis, right) tends toward lower accuracy (y-axis, bottom).

purely test-time and does not require storing source statistics. Instead, it relies on the confidence of attention during inference alone. We therefore only compare with other fully test-time updates.

2 METHOD: ATTENTION ENTROPY MINIMIZATION

Given a model f_θ trained on a source distribution \mathcal{D}_{source} , we encounter a batch from a shifted distribution \mathcal{D}_{shift} at test time. For each test batch, $\mathcal{B} = \{x_i\}_{i=1}^B$ where $x_i \sim \mathcal{D}_{shift}$, we aim to:

1. Adapt model parameters θ using an unsupervised objective $\mathcal{L}(x_i; \theta)$.
2. Generate prediction $\hat{y}_i = f_{\theta'}(x_i)$ using the adapted model.
3. Reset the model parameters to the original pretrained state (episodic).

Loss: Attention Entropy Minimization. Let $\mathbf{A}(x_i) \in \mathbb{R}^{H \times T \times T}$ denote the post-softmax attention tensor from the final transformer layer for input image x_i , where H is the number of attention heads, T is the sequence length (CLS, register, and patch tokens), and t_{cls} is the index of the CLS token. Let \mathbb{P} denote the set of patch-token indices (excluding the CLS and register tokens), with $P = |\mathbb{P}|$. We extract scores for the CLS token attending to the patch tokens and renormalize them to form a distribution $a^{(h)}(x_i)$:

$$a_j^{(h)}(x_i) = \frac{\mathbf{A}_{h,t_{cls},j}(x_i)}{\sum_{k \in \mathbb{P}} \mathbf{A}_{h,t_{cls},k}(x_i)}, \quad \mathcal{L}_{Attention}(x_i) = -\frac{1}{H} \sum_{h=1}^H \sum_{j \in \mathbb{P}} a_j^{(h)}(x_i) \log a_j^{(h)}(x_i) \quad (1)$$

We **exclude** the CLS token to itself and to register tokens attention scores as we want to focus on the spatial patches of the image as opposed to global information. Minimizing this loss encourages each attention head to place concentrated (low-entropy) focus on a smaller subset of patch tokens, rather than distributing attention more diffusely. Averaging the distributions first, then taking their entropy, was tried but performed worse. This is reasonable as heads tend to specialize (Raghu et al., 2021). We utilize the last layers attention scores as they are the most mature.

We found that combining the standard output entropy minimization, as in (Wang et al., 2021), with attention entropy minimization further improves performance. We use the standard output entropy

| Corruption | Source | Tent (Online) | Tent (Episodic) | Output | Attention (<i>Ours</i>) | Combined (<i>Ours</i>) |
|-------------------|--------|------------------|--------------------|------------------|------------------------------|--------------------------------|
| Brightness | 76.06 | 76.48 | 76.54 | 76.92 | 76.79 | 77.08 |
| Contrast | 51.96 | 55.68 | 55.01 | 58.72 | 56.72 | 58.95 |
| Defocus Blur | 42.87 | 44.61 | 45.43 | 46.87 | 48.26 | 47.47 |
| Elastic Transform | 33.94 | 40.57 | 37.51 | 42.35 | 47.64 | 44.15 |
| Fog | 57.20 | 59.77 | 59.16 | 61.55 | 62.91 | 62.31 |
| Frost | 49.72 | 51.71 | 51.23 | 53.92 | 55.38 | 54.58 |
| Gaussian Noise | 33.15 | 36.23 | 36.37 | 40.78 | 32.87 | 40.95 |
| Glass Blur | 21.68 | 31.11 | 26.75 | 37.95 | 38.20 | 38.69 |
| Impulse Noise | 37.05 | 31.08 | 39.66 | 42.14 | 41.01 | 43.19 |
| JPEG Compression | 59.27 | 61.57 | 60.91 | 62.04 | 61.87 | 62.33 |
| Motion Blur | 48.49 | 50.93 | 50.27 | 53.15 | 52.85 | 53.48 |
| Pixelate | 63.27 | 65.34 | 64.88 | 66.57 | 67.42 | 66.91 |
| Shot Noise | 35.33 | 39.33 | 38.70 | 44.69 | 38.99 | 45.06 |
| Snow | 56.94 | 59.37 | 58.82 | 61.70 | 59.15 | 61.96 |
| Zoom Blur | 46.11 | 49.52 | 48.83 | 52.73 | 53.32 | 53.15 |
| Mean | 47.54 | 50.22 (+2.68) | 50.01 (+2.47) | 53.47 (+5.93) | 52.89 (+5.35) | 54.02 (+6.48) |

Table 1: Top-1 Accuracy (%) on ImageNet-C level 5 corruptions. We report the source model and test-time adaptation variants: Tent (online/episodic), output entropy, attention entropy, and their combination (*LookSharp*). All results use batch size 128. Attention, output and combined reset all the parameters.

minimization loss (Wang et al., 2021):

$$\mathcal{L}_{Output}(x_i) = - \sum_{c \in \mathcal{C}} p_{\theta}(c | x_i) \log p_{\theta}(c | x_i) \quad (2)$$

where \mathcal{C} denotes the set of classes and $p_{\theta}(c | x_i)$ is the model’s predicted class probability.

Thus, our loss is:

$$\mathcal{L}_{Combined}(x_i) = \mathcal{L}_{Attention}(x_i) + \mathcal{L}_{Output}(x_i) \quad (3)$$

3 EXPERIMENTS AND RESULTS

We experiment with the standard benchmark for test-time adaptation applied to image classification, using a common architecture and a recent self-supervised backbone. We consider the batch-wise episodic test-time adaptation setting where the parameters are reset after each batch Zhang et al. (2022), and also compare to an online (no resetting) method (Wang et al., 2021).

Dataset: We evaluate on ImageNet-C (Hendrycks & Dietterich, 2019), which augments the standard ImageNet validation set with 15 different corruption types at 5 levels. We only evaluate on level 5, which is the most severe level of shift. We also perform TTA on clean data to ensure our method maintains performance without distribution shift.

Model: We use DINOv3-Base (Siméoni et al., 2025), pretrained on an internet-scale image dataset. We train a linear classification head with this representation on the source data (ImageNet training split) using the standard cross-entropy loss (a.k.a. linear probing). This yields 83.57% top-1 accuracy on the validation set. The images are preprocessed to the standard ImageNet size (224×224) as in Krizhevsky et al. (2012).

Evaluation Protocol: For each corruption type, we report per-corruption accuracy and mean corruption accuracy at level 5. We use a batch size of 128. The data is loaded in a randomized order for each shift, and as a result, each batch contains a mix of classes. We optimize by Adam (Kingma & Ba, 2015) with learning rate 5×10^{-5} for all methods except Tent. For Tent, we use 10^{-3} in the episodic setting and 10^{-5} in the online setting. These values are selected by a learning-rate sweep

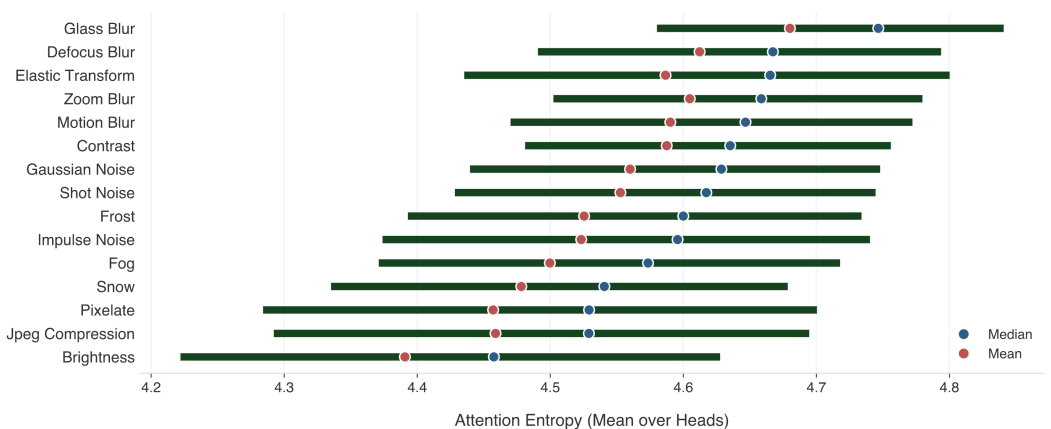


Figure 2: Median and interquartile range (IQR) of per-image mean attention entropy across a 10% sample of ImageNet-C at level 5. For each corruption, attention entropy is first averaged over heads for each image. Blurs and blur-like corruptions tend to have higher attention entropies.

on the level 5 test set with mean accuracy as the metric. We perform 1 gradient update per batch and update all parameters.

Baselines: We evaluate without any test-time updates, to measure the robustness of the source model. We also compare with Tent (Wang et al., 2021), where only the normalization layer parameters are updated, in the episodic and in the online case.

Results. Table 1 shows that our method improves mean accuracy compared to the non-adapted source model on ImageNet-C. The output-head entropy loss alone performs better than attention entropy alone, but combining both losses yields even better results. On clean data, the attention-only loss *slightly* hurts performance (83.57% \rightarrow 82.95%). Using the combined loss *slightly* improves accuracy (83.57% \rightarrow 83.80%).

Overall, our combined objective achieves the best mean corruption accuracy, improving mean accuracy from 47.54% (Source) to 54.02% (+6.48 %). Attention-based entropy minimization works best for blur and blur-like corruptions (elastic transform). We can see from Figure 2 that this is because blurring images makes the attention maps more diffuse, and this is what $\mathcal{L}_{Attention}$ is directly addressing. A visualization of the attention loss is shown in Appendix A.

In our experiments, we found that Tent (Online) is highly sensitive to the learning rate, consistent with Zhao et al. (2023). Larger learning rates improve performance on some corruptions but cause the model to collapse on others, resulting in mean accuracy below the source model. The learning rate we select is the one that achieves maximal mean accuracy on the level 5 test set.

4 CONCLUSION AND FUTURE WORK

We introduce *LookSharp*, a simple test-time adaptation method that minimizes the entropy of CLS-to-patch attention, and show consistent gains on ImageNet-C, especially for blur-like corruptions. Combining attention and output entropy yields the best overall accuracy, suggesting the two signals are complementary.

Limitations. The method incurs computational overhead due to the forward-backward-forward passes needed and requires self-attention in the model architecture. Attention-based adaptation likely also depends on the quality of the learned attention maps, which vary across architectures and pretraining regimes (Darcet et al., 2024).

While this work focuses on succinct experiments to show the effectiveness of attention entropy as an unsupervised TTA loss, future work can explore trying to extract more performance by exploring a dynamic weighting of attention and output entropy based on input characteristics or multi-layer attention losses that span the model from shallow to deep.

ACKNOWLEDGEMENTS

We thank Vivian White for helpful discussions. We also thank the Digital Research Alliance of Canada and the Vector Institute for computational resources.

REFERENCES

- Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. Vision transformers need registers. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2024.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale, 2021. URL <https://arxiv.org/abs/2010.11929>.
- Anthony Fuller, Yousef Yassin, Junfeng Wen, Daniel G. Kyrollos, Tarek Ibrahim, James R. Green, and Evan Shelhamer. Lookwhere? efficient visual recognition by learning where to look and what to see from self-supervision, 2025. URL <https://arxiv.org/abs/2505.18051>.
- Dan Hendrycks and Thomas Dietterich. Benchmarking neural network robustness to common corruptions and perturbations. *Proceedings of the International Conference on Learning Representations*, 2019.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations (ICLR)*, 2015.
- Takuya Kojima, Yusuke Iwasawa, and Yutaka Matsuo. Robustifying vision transformer without retraining from scratch using attention-based test-time adaptation. *New Generation Computing*, 41:5–24, 2023. doi: 10.1007/s00354-022-00197-9.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*, 2012.
- Shuaicheng Niu, Jiayang Wu, Yifan Zhang, Yaofo Chen, Shijian Zheng, Peilin Zhao, and Mingkui Tan. Efficient test-time model adaptation without forgetting. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 16888–16905. PMLR, 17–23 Jul 2022. URL <https://proceedings.mlr.press/v162/niu22a.html>.
- Shuaicheng Niu, Jiayang Wu, Yifan Zhang, Zhiquan Wen, Yaofo Chen, Peilin Zhao, and Mingkui Tan. Towards stable test-time adaptation in dynamic wild world. In *International Conference on Learning Representations*, 2023.
- Maithra Raghu, Thomas Unterthiner, Simon Kornblith, Chiyuan Zhang, and Alexey Dosovitskiy. Do vision transformers see like convolutional neural networks? In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, volume 34, pp. 12116–12128. Curran Associates, Inc., 2021. URL https://proceedings.neurips.cc/paper_files/paper/2021/file/652cf38361a209088302ba2b8b7f51e0-Paper.pdf.
- Oriane Siméoni, Huy V. Vo, Maximilian Seitzer, Federico Baldassarre, Maxime Oquab, Cijo Jose, Vasil Khalidov, Marc Szafraniec, Seungeun Yi, Michaël Ramamonjisoa, Francisco Massa, Daniel Haziza, Luca Wehrstedt, Jianyuan Wang, Timothée Darcet, Théo Moutakanni, Leonel Sentana, Claire Roberts, Andrea Vedaldi, Jamie Tolan, John Brandt, Camille Couprie, Julien Mairal, Hervé Jégou, Patrick Labatut, and Piotr Bojanowski. DINOv3, 2025. URL <https://arxiv.org/abs/2508.10104>.
- Dequan Wang, Evan Shelhamer, Shaoteng Liu, Bruno Olshausen, and Trevor Darrell. Tent: Fully test-time adaptation by entropy minimization. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=uXl3bZLkr3c>.

Marvin Zhang, Sergey Levine, and Chelsea Finn. Memo: Test time robustness via adaptation and augmentation. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 38629–38642. Curran Associates, Inc., 2022. URL

https://proceedings.neurips.cc/paper_files/paper/2022/file/fc28053a08f59fccb48b11f2e31e81c7-Paper-Conference.pdf.

Hao Zhao, Yuejiang Liu, Alexandre Alahi, and Tao Lin. On pitfalls of test-time adaptation. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 42058–42080. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/zhao23d.html>.

A APPENDIX

The figure below shows the CLS-to-patch attention distribution of the model taken from the final layer and averaged across the heads.



Figure 3: This shows the attention map before and after adaptation using our attention entropy loss ($\mathcal{L}_{Attention}$).