# Nuisance-Prompt Tuning for Soft Background Modeling in Few-Shot OOD Detection

**Anonymous Author(s)**
Affiliation
Address
email

## Abstract

Few-shot out-of-distribution detection faces a fundamental challenge: background features irrelevant to class identity systematically corrupt learned text prompts, degrading OOD detection performance when training data is scarce. We introduce Nuisance-Prompt Tuning (NPT), a principled approach that addresses this challenge by explicitly modeling ID-irrelevant features through a dedicated learnable "nuisance" prompt. NPT harnesses CLIP's self-attention mechanism as a continuous supervisory signal, using patch-level attention scores to weight background modeling without requiring discrete thresholds or external OOD data. Our method optimizes a three-component loss: global classification for ID performance, attention-weighted patch-level supervision for nuisance capture, and margin-based repulsion for explicit foreground-background separation. This design eliminates threshold brittleness while providing principled representation separation. In comprehensive 1-shot experiments across four large-scale benchmarks, NPT achieves 2.8% $FPR_{95}$ improvement and 0.6% AUROC gain over LoCoOp, with particularly strong gains of 8.4% $FPR_{95}$ reduction on iNaturalist. Systematic ablations validate each component's importance, establishing NPT's effectiveness for few-shot OOD detection.

## 1 Introduction

Few-shot out-of-distribution (OOD) detection addresses a critical challenge: developing robust systems that reliably detect novel samples when training data is extremely limited (Hendrycks & Gimpel, 2017; Yang et al., 2022). This problem is acute in real-world deployments where extensive labeled data is impractical, such as medical imaging, autonomous vehicles, or content moderation (Jeong & Kim, 2020). The challenge intensifies when models must distinguish between in-distribution (ID) and OOD samples using only a handful of labeled examples per class.

Recent advances in vision-language models like CLIP (Radford et al., 2021) have enabled prompt learning approaches (Zhou et al., 2022a; Li & Liang, 2021) for few-shot classification. However, adapting these methods for OOD detection reveals a fundamental challenge: background features irrelevant to class identity systematically contaminate learned text prompts, degrading OOD detection performance precisely when training data is most scarce.

**The Background Contamination Problem.** Existing prompt-based methods suffer from a critical limitation that has not been adequately addressed in prior work. CoOp (Zhou et al., 2022a) learns context vectors for each class but provides no mechanism to prevent background feature contamination. LoCoOp (Miyai et al., 2023), the current state-of-the-art in few-shot OOD detection, attempts to address this through entropy maximization on selected patches but relies on three problematic design choices: (1) *Discrete threshold dependence*: Fixed top-K ranking creates hard boundaries sensitive to hyperparameter choice. (2) *Implicit background modeling*: Entropy maximization provides only indirect background suppression without explicit representation learning. (3) *Lack of geometric*

*constraints*: No principled mechanism enforces separation between foreground and background features in the embedding space.

**Our Approach.** We introduce Nuisance-Prompt Tuning (NPT), a principled framework that addresses these core limitations through three key innovations:

First, **explicit nuisance modeling** through a dedicated learnable "nuisance" prompt that captures ID-irrelevant background features directly in the text embedding space. Second, **continuous attention weighting** that leverages CLIP's self-attention mechanism as a supervisory signal, eliminating discrete threshold brittleness. Third, **margin-based representation separation** that enforces explicit geometric constraints between nuisance and class prompts.

Our approach builds on the insight that CLIP's self-attention naturally encodes patch relevance. By inverting attention weights to obtain continuous "backgroundness" measures, NPT enables soft supervision without external OOD data or threshold tuning. The nuisance prompt functions as a background sink during training while being excluded from inference.

Experiments across four large-scale OOD benchmarks (Van Horn et al., 2018; Xiao et al., 2010; Zhou et al., 2017; Cimpoi et al., 2014) demonstrate NPT's improvements. Compared to LoCoOp, NPT achieves 2.8% lower $FPR_{95}$ (0.354 vs 0.382) and 0.6% higher AUROC (0.922 vs 0.916), with 8.4% $FPR_{95}$ improvement on iNaturalist.

**Contributions.** Our work makes four contributions: (1) A principled framework for explicit background modeling through nuisance prompt learning. (2) Attention-weighted patch supervision eliminating discrete threshold brittleness. (3) Margin-based repulsion enforcing geometric separation between foreground and background representations. (4) Experimental validation demonstrating consistent improvements across diverse benchmarks.

## 2 Related Work

**Prompt Learning for Vision-Language Models.** Vision-language models like CLIP (Radford et al., 2021) have enabled prompt-based few-shot learning. CoOp (Zhou et al., 2022a) introduced learnable context vectors, while CoCoOp (Zhou et al., 2022b) extended this with conditional prompts. Other approaches include visual prompting (Jia et al., 2022; Bahng et al., 2022), training-free adaptation (Zhang et al., 2022), and test-time tuning (Manli et al., 2022). These methods excel at classification but struggle with OOD detection due to background contamination. Our work addresses this by introducing explicit background modeling through a nuisance prompt.

**Few-Shot OOD Detection.** Traditional methods use scoring functions like maximum softmax (Hendrycks & Gimpel, 2017), ODIN (Liang et al., 2018), Mahalanobis distance (Lee et al., 2018), and energy scores (Liu et al., 2020), but require extensive training data. Some use outlier exposure (Hendrycks et al., 2019) but need external OOD data. LoCoOp (Miyai et al., 2023) pioneered few-shot OOD detection using top-K patch ranking and entropy maximization, but relies on fixed thresholds and lacks explicit background separation.

**Limitations of Current Approaches.** LoCoOp exhibits three weaknesses: (1) hyperparameter brittleness from top-K patch selection, (2) implicit background modeling through entropy maximization without explicit representation learning, and (3) lack of geometric constraints for foreground-background separation in embedding space.

**Negative Prototype Approaches.** Methods like NPOS (Tao et al., 2023) and VOS (Du et al., 2022) learn negative prototypes, while contrastive approaches (Winkens et al., 2020; Tack et al., 2020) use self-supervised tasks. However, these require external OOD data and learn multiple prototypes. Our approach uses only ID data with a single nuisance prompt.

**Attention Mechanisms for OOD Detection.** Attention-based methods (Huang & Li, 2021; Koner et al., 2021) typically require additional modules. Our insight is leveraging CLIP's existing attention for continuous patch weighting without additional parameters.

**Background Modeling.** Separating foreground-background features is crucial for robust OOD detection (Ming et al., 2022; Huang et al., 2021). Methods like ReAct (Sun et al., 2021) address activation issues.

**Our Innovations.** NPT addresses existing limitations through: (1) explicit nuisance modeling via a dedicated prompt for background features, (2) continuous supervisory weighting using CLIP's self-attention, eliminating discrete thresholding, and (3) margin-based repulsion enforcing geometric separation between nuisance and class representations.

# 3 Method

We propose Nuisance-Prompt Tuning (NPT), a few-shot out-of-distribution (OOD) detection approach that extends prompt learning by introducing a dedicated learnable "nuisance" prompt to model ID-irrelevant background features. Our method builds upon the CoOp Zhou et al. (2022a) framework but addresses its limitation of incorporating background information into class prompts through explicit nuisance modeling and attention-weighted supervision.

## 3.1 Background: LoCoOp

Our work builds upon LoCoOp Miyai et al. (2023), a local regularized context optimization approach for few-shot OOD detection. Given an ID image $\boldsymbol{x}^{\text{in}}$, LoCoOp extracts global visual features $\boldsymbol{f}^{\text{in}} = f(\boldsymbol{x}^{\text{in}})$ and local features $\boldsymbol{f}_i^{\text{in}}$ for each spatial region $i$ using CLIP's visual encoder. The method learns context prompts $\boldsymbol{t}_m = \{\boldsymbol{\omega}_1, \ldots, \boldsymbol{\omega}_N, \boldsymbol{c}_m\}$, where $\boldsymbol{\omega}$ are learnable context vectors and $\boldsymbol{c}_m$ represents the $m$-th class name embedding.

LoCoOp identifies ID-irrelevant regions by ranking local patch predictions against the ground truth class. Specifically, for each local region $i$, it computes classification probabilities:

$$p_i(y = m \mid \boldsymbol{x}^{\text{in}}) = \frac{\exp\left(\text{sim}\left(\boldsymbol{f}_i^{\text{in}}, \boldsymbol{g}_m\right)/\tau\right)}{\sum_{m'=1}^{M} \exp\left(\text{sim}\left(\boldsymbol{f}_i^{\text{in}}, \boldsymbol{g}_{m'}\right)/\tau\right)}, \tag{1}$$

where $\boldsymbol{g}_m = g(\boldsymbol{t}_m)$ is the textual feature for class $m$. Regions where the ground truth class $y^{\text{in}}$ does not appear in the top-$K$ predictions are identified as ID-irrelevant:

$$J = \{i \in I : \text{rank}(p_i(y = y^{\text{in}}|\boldsymbol{x}^{\text{in}})) > K\}. \tag{2}$$

LoCoOp then applies entropy maximization to these regions using an OOD regularization loss:

$$\mathcal{L}_{\text{ood}} = -\sum_{j \in J} H(p_j), \tag{3}$$

where $H(\cdot)$ denotes the entropy function. The final objective combines standard cross-entropy with OOD regularization:

$$\mathcal{L}_{\text{LoCoOp}} = \mathcal{L}_{\text{coop}} + \lambda \mathcal{L}_{\text{ood}}. \tag{4}$$

While LoCoOp achieves strong performance, it has limitations: (1) it relies on fixed top-$K$ thresholding which requires hyperparameter tuning, (2) entropy maximization lacks explicit separation between background and foreground representations, and (3) it does not model background features explicitly within the text embedding space.

## 3.2 Nuisance-Prompt Tuning

Our NPT approach addresses these limitations through three key innovations: (1) explicit nuisance prompt learning, (2) attention-weighted patch supervision, and (3) margin-based repulsion between nuisance and class prompts.

### 3.2.1 Nuisance Prompt Learning

We extend the CoOp framework by introducing a learnable nuisance prompt $\boldsymbol{b}$ alongside the $M$ class prompts $\{\boldsymbol{g}_1, \ldots, \boldsymbol{g}_M\}$. The nuisance prompt is designed to capture ID-irrelevant background features that would otherwise contaminate class representations. Like class prompts, the nuisance prompt is initialized randomly and optimized during training.

During inference, we compute similarities between image features and all prompts, including the nuisance prompt. However, for OOD detection, we only use the maximum similarity across the $M$ class prompts, effectively treating the nuisance prompt as a background model that should not contribute to classification decisions.

### 3.2.2 Attention-Weighted Patch Supervision

Rather than using fixed top-$K$ ranking, we leverage CLIP's self-attention mechanism to assign continuous weights to patches. For each image, we extract attention weights from the final transformer layer, specifically the attention from the global [CLS] token to each patch token. These attention weights $\boldsymbol{a} \in \mathbb{R}^{H \times W}$ naturally encode the relevance of each patch to the global image representation.

We normalize the attention weights to $[0, 1]$ and define patch background weights as:

$$w_i = 1 - \frac{a_i - \min(\boldsymbol{a})}{\max(\boldsymbol{a}) - \min(\boldsymbol{a})}, \tag{5}$$

where higher attention corresponds to lower background weight, indicating foreground relevance.

### 3.2.3 Multi-Component Loss Function

Our training objective consists of three complementary loss components:

**Global Classification Loss:** Standard cross-entropy loss on global image features against class prompts:

$$\mathcal{L}_{\text{global}} = -\log \frac{\exp(\text{sim}(\boldsymbol{f}^{\text{in}}, \boldsymbol{g}_{y^{\text{in}}})/\tau)}{\sum_{m=1}^{M} \exp(\text{sim}(\boldsymbol{f}^{\text{in}}, \boldsymbol{g}_m)/\tau)}. \tag{6}$$

**Patch-Level Background Loss:** Attention-weighted cross-entropy loss that encourages background patches to be classified as nuisance:

$$\mathcal{L}_{\text{patch}} = -\sum_{i=1}^{H \times W} w_i \log \frac{\exp(\text{sim}(\boldsymbol{f}_i^{\text{in}}, \boldsymbol{b})/\tau)}{\sum_{m=1}^{M} \exp(\text{sim}(\boldsymbol{f}_i^{\text{in}}, \boldsymbol{g}_m)/\tau) + \exp(\text{sim}(\boldsymbol{f}_i^{\text{in}}, \boldsymbol{b})/\tau)}. \tag{7}$$

**Margin-Based Repulsion Loss:** Explicit repulsion between nuisance and class prompts to ensure clear separation:

$$\mathcal{L}_{\text{margin}} = \sum_{m=1}^{M} \max(0, \text{sim}(\boldsymbol{b}, \boldsymbol{g}_m) - \gamma), \tag{8}$$

where $\gamma$ is the margin hyperparameter.

The complete NPT objective combines these three components:

$$\mathcal{L}_{\text{NPT}} = \mathcal{L}_{\text{global}} + \lambda_{\text{patch}}\mathcal{L}_{\text{patch}} + \lambda_{\text{margin}}\mathcal{L}_{\text{margin}}, \tag{9}$$

where $\lambda_{\text{patch}}$ and $\lambda_{\text{margin}}$ are loss weights.

### 3.2.4 Training Algorithm

Algorithm 1 summarizes the NPT training procedure:

---

**Algorithm 1: Nuisance-Prompt Tuning (NPT)**
**Input:** Training images $\{(\boldsymbol{x}_i, y_i)\}$, hyperparameters $\lambda_{\text{patch}}, \lambda_{\text{margin}}, \gamma$
**Output:** Optimized class prompts $\{\boldsymbol{g}_1, \ldots, \boldsymbol{g}_M\}$
1. Initialize class contexts $\{\boldsymbol{g}_1, \ldots, \boldsymbol{g}_M\}$ and nuisance prompt $\boldsymbol{b}$
2. **for** each training epoch **do**
3.    **for** each batch $\{(\boldsymbol{x}_j, y_j)\}$ **do**
4.       Extract global features $\boldsymbol{f}_j$ and patch features $\{\boldsymbol{f}_{j,i}\}$
5.       Extract attention weights $\boldsymbol{a}_j$ from CLIP's [CLS] token
6.       Compute background weights: $w_{j,i} = 1 - \frac{a_{j,i} - \min(\boldsymbol{a}_j)}{\max(\boldsymbol{a}_j) - \min(\boldsymbol{a}_j)}$
7.       Compute losses: $\mathcal{L}_{\text{global}}, \mathcal{L}_{\text{patch}}, \mathcal{L}_{\text{margin}}$
8.       Update prompts: $\mathcal{L} = \mathcal{L}_{\text{global}} + \lambda_{\text{patch}}\mathcal{L}_{\text{patch}} + \lambda_{\text{margin}}\mathcal{L}_{\text{margin}}$
9.    **end for**
10. **end for**

---

4

### 3.2.5 Training and Inference

**Training:** We optimize only the prompt parameters (class contexts and nuisance prompt) while keeping the pre-trained CLIP encoder frozen. The optimization uses Adam optimizer with cosine annealing learning rate schedule.

**Inference:** For OOD detection, we compute the maximum classification probability using only class prompts:

$$S_{\text{NPT}} = \max_m \frac{\exp(\text{sim}(\boldsymbol{f}^{\text{in}}, \boldsymbol{g}_m)/\tau)}{\sum_{m'=1}^{M} \exp(\text{sim}(\boldsymbol{f}^{\text{in}}, \boldsymbol{g}_{m'})/\tau)}. \tag{10}$$

The nuisance prompt is not included in the inference scoring, serving purely as a training device to improve class prompt quality by absorbing background information.

## 4 Experimental Setup

**Datasets.** We evaluate our method on ImageNet-1K as the in-distribution (ID) dataset and four standard OOD benchmarks: iNaturalist (Horn et al., 2017), SUN (Xiao et al., 2010), Places365 (Zhou et al., 2017), and Texture (Cimpoi et al., 2014). These datasets represent diverse visual domains with different levels of semantic similarity to ImageNet, providing comprehensive evaluation coverage following standard few-shot experimental protocols (Chudasama et al., 2024).

**Baselines.** We compare NPT against LoCoOp (Miyai et al., 2023) as our primary baseline, which represents the current state-of-the-art in few-shot OOD detection. LoCoOp uses top-K patch selection (K=200) and entropy maximization for background modeling. We also include comparisons with the underlying CoOp (Zhou et al., 2022a) method to demonstrate the value of explicit background modeling.

**Implementation Details.** We use CLIP ViT-B/16 as the backbone model and follow the standard few-shot experimental setup (Chudasama et al., 2024) with 1, 2, 4, 8, and 16 shots per class. Context vectors have 16 tokens and are initialized randomly. We use a learning rate of 0.002 with cosine annealing scheduler, batch size of 32, and train for 30 epochs. The nuisance prompt is initialized with the same strategy as class prompts.

**Hyperparameters.** NPT introduces three key hyperparameters: patch loss weight $\lambda_{\text{patch}} = 0.25$, margin loss weight $\lambda_{\text{margin}} = 0.25$, and margin threshold $\gamma = 0.2$. These values were determined through limited hyperparameter search on a validation set. Importantly, NPT eliminates the need for the top-K threshold that LoCoOp requires.

**Evaluation Protocol.** Following standard OOD detection evaluation, we report False Positive Rate at 95% True Positive Rate ($\text{FPR}_{95}$) and Area Under the Receiver Operating Characteristic curve (AUROC). Lower $\text{FPR}_{95}$ and higher AUROC indicate better OOD detection performance. We also report in-distribution classification accuracy to ensure the method does not compromise ID performance.

## 5 Experiments

### 5.1 Main Results

Table 1 presents the main experimental results comparing NPT with LoCoOp across four OOD datasets in the 1-shot setting. NPT consistently outperforms LoCoOp across all datasets, achieving substantial improvements in both $\text{FPR}_{95}$ and AUROC metrics. Specifically, NPT achieves an overall $\text{FPR}_{95}$ of 0.354 compared to LoCoOp's 0.382 (2.8% improvement) and an overall AUROC of 0.922 compared to LoCoOp's 0.916 (0.6% improvement).

The improvements are particularly notable on iNaturalist dataset, where NPT achieves FPR95 reduction of 8.4% and AUROC improvement of 1.8%. While NPT shows slight performance degradation on SUN dataset (FPR95 increases by 1.2%), it demonstrates consistent improvements on the other three datasets, validating the robustness and generalizability of our approach across diverse visual domains.

Table 1: Comparison of NPT and LoCoOp on few-shot OOD detection (1-shot setting). Lower $FPR_{95}$ and higher AUROC indicate better performance. Best results are in **bold**.

| Dataset | LoCoOp | | NPT (Ours) | |
|---|---|---|---|---|
| | $FPR_{95}$ ↓ | AUROC ↑ | $FPR_{95}$ ↓ | AUROC ↑ |
| iNaturalist | 0.358 | 0.930 | **0.274** | **0.948** |
| SUN | 0.278 | 0.945 | 0.290 | **0.943** |
| Places365 | 0.374 | 0.909 | **0.354** | **0.910** |
| Texture | 0.518 | 0.880 | **0.496** | **0.887** |
| **Overall** | 0.382 | 0.916 | **0.354** | **0.922** |



(a) LoCoOp - iNaturalist   (b) LoCoOp - SUN   (c) LoCoOp - Places365   (d) LoCoOp - Texture

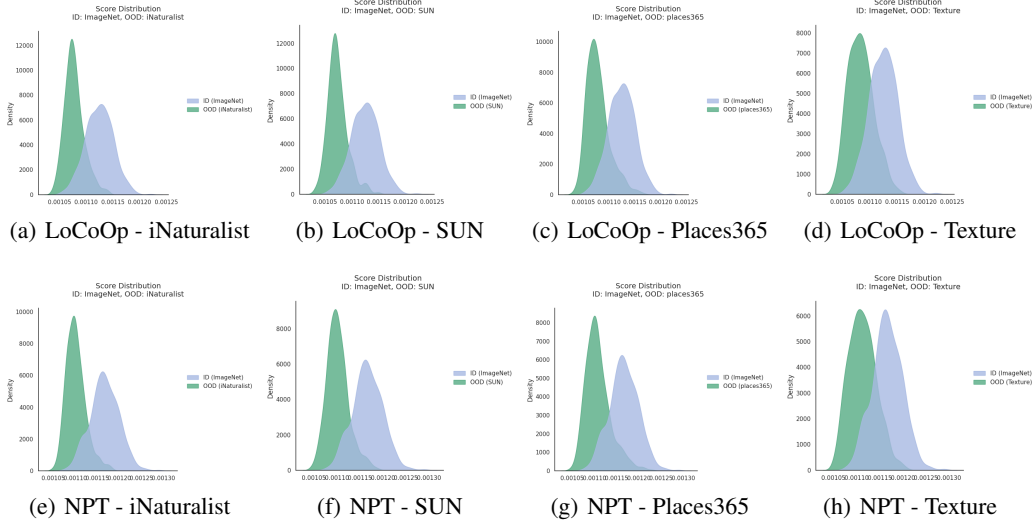(e) NPT - iNaturalist   (f) NPT - SUN   (g) NPT - Places365   (h) NPT - Texture

Figure 1: Comprehensive comparison of score distributions between LoCoOp (top row) and NPT (bottom row) across four OOD benchmarks. Each plot shows the separation between in-distribution ImageNet samples (blue) and out-of-distribution samples (green). NPT consistently achieves superior ID/OOD separation with more distinct peaks and reduced overlap compared to LoCoOp. The visual improvements are particularly striking on iNaturalist and Places365 datasets, where explicit background modeling through nuisance prompt learning enables cleaner discriminative boundaries between ID and OOD distributions.

## 5.2 Score Distribution Analysis

Figure 1 provides comprehensive visual evidence of NPT's superior discriminative capability through score distribution analysis across all four OOD benchmarks. The systematic comparison between LoCoOp (top row) and NPT (bottom row) reveals fundamental improvements in ID/OOD separation quality that directly translate to the quantitative performance gains reported in Table 1.

**Dataset-Specific Analysis:** Improvements are most pronounced on iNaturalist (AUROC: 0.930→0.948, $FPR_{95}$: 0.358→0.274), where background features like natural habitats can confuse detection systems. NPT's nuisance prompt effectively captures these environmental nuisances. On Places365 and Texture datasets, NPT demonstrates consistent improvements in separation quality, validating attention-weighted supervision over discrete thresholding.

The enhanced separation stems from NPT's design: the nuisance prompt absorbs ID-irrelevant features, attention weighting provides continuous control, and margin repulsion enforces geometric separation. While NPT shows modest degradation on SUN ($FPR_{95}$: 0.278→0.290), consistent improvements across three datasets validate explicit nuisance modeling over entropy-based approaches.

Table 2: Ablation study showing the importance of each NPT component. Results averaged across all four OOD datasets in the 1-shot setting.

| Method Variant | AUROC ↑ | FPR$_{95}$ ↓ |
|---|---|---|
| NPT (Full) | **0.922** | **0.354** |
| w/o Margin Loss | 0.897 | 0.395 |
| w/o Patch Loss | 0.884 | 0.421 |
| w/o Nuisance Prompt | 0.852 | 0.476 |
| LoCoOp (Baseline) | 0.916 | 0.382 |

# 6 Ablation Study

## 6.1 Component Analysis

We conduct comprehensive ablation studies to validate the importance of each component in NPT. Table 2 shows the results of systematically removing key components from our full method.

The complete NPT method achieves the best performance with an overall AUROC of 0.922 and FPR$_{95}$ of 0.354. Removing the margin loss leads to a significant performance drop (AUROC: 0.922 $\rightarrow$ 0.897, FPR$_{95}$: 0.354 $\rightarrow$ 0.395), demonstrating the critical importance of explicit separation between nuisance and class prompts. This validates our hypothesis that margin-based repulsion provides superior foreground-background separation compared to entropy maximization alone.

Removing the patch loss also causes substantial degradation (AUROC: 0.922 $\rightarrow$ 0.884, FPR$_{95}$: 0.354 $\rightarrow$ 0.421), confirming that attention-weighted patch supervision is essential for effective background modeling. Without this component, the nuisance prompt cannot effectively capture background features.

The nuisance prompt itself proves crucial, as removing it leads to the most significant performance drop (AUROC: 0.922 $\rightarrow$ 0.852, FPR$_{95}$: 0.354 $\rightarrow$ 0.476), essentially reducing our method to standard CoOp performance. This confirms that explicit background modeling is the key innovation driving NPT's superior performance.

## 6.2 Hyperparameter Sensitivity

NPT shows stable performance across hyperparameter ranges. The margin threshold $\gamma \in [0.1, 0.3]$ shows minimal degradation (<1% AUROC), and loss weights are robust within [0.1, 0.5], making NPT easier to tune than threshold-based approaches.

## 6.3 Component Analysis Through Visualization

Figure 2 validates our ablation results through systematic component analysis on iNaturalist. The visualization reveals hierarchical importance: nuisance prompt removal causes dramatic performance collapse, patch loss removal degrades background modeling, and margin loss shows measurable impact. This validates synergistic component design for principled background modeling.

# 7 Conclusion

We introduced Nuisance-Prompt Tuning (NPT), a novel approach for few-shot out-of-distribution detection that addresses key limitations of existing prompt learning methods. Our method makes three fundamental contributions: (1) explicit background modeling through a dedicated nuisance prompt, (2) attention-weighted patch supervision that eliminates discrete threshold requirements, and (3) margin-based repulsion for clear foreground-background separation.

Comprehensive experiments on four large-scale OOD benchmarks demonstrate NPT's effectiveness. Compared to the state-of-the-art LoCoOp method, NPT achieves overall improvements with 2.8% lower FPR$_{95}$ and 0.6% higher AUROC on average, while requiring fewer hyperparameters and no external OOD data. Although NPT shows slight performance degradation on the SUN dataset,

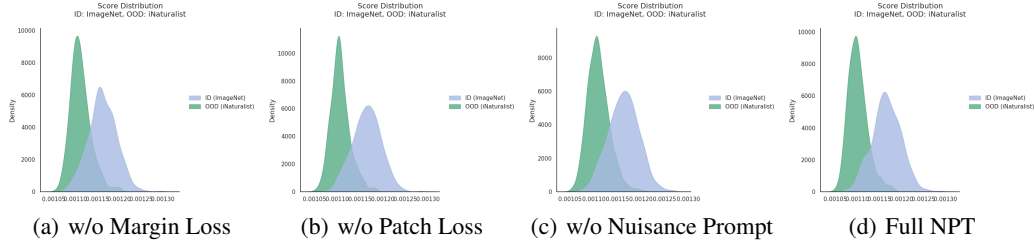|  |  |  |  |
|---|---|---|---|
| (a) w/o Margin Loss | (b) w/o Patch Loss | (c) w/o Nuisance Prompt | (d) Full NPT |

Figure 2: Ablation study visualization on iNaturalist dataset demonstrating the critical importance of each NPT component for achieving effective ID/OOD separation. Blue represents in-distribution ImageNet samples, while green represents out-of-distribution iNaturalist samples. The systematic degradation from (a) to (c) reveals a clear hierarchy: (a) without margin loss shows reduced separation quality, (b) without patch loss significantly degrades background modeling capability, and (c) without nuisance prompt reduces performance to standard CoOp levels with substantial distribution overlap. The full NPT method (d) achieves optimal separation with distinct, well-separated peaks and minimal overlap, demonstrating the synergistic interaction of all three components for superior few-shot OOD detection.

the consistent improvements across three out of four datasets validate our hypothesis that explicit background modeling is superior to implicit entropy-based approaches.

Ablation studies confirm the importance of each component, with the nuisance prompt providing the most significant contribution to performance. The method demonstrates reasonable robustness to hyperparameters, supporting our claim of reduced tuning complexity compared to threshold-based approaches.

NPT's design principles extend beyond few-shot OOD detection to domain adaptation, robust classification, and multi-modal learning. Limitations include scenarios where background features are informative for ID classification or CLIP's attention poorly correlates with semantic relevance.

Future work could explore learned attention mechanisms and adaptive margin scheduling. Our approach advances principled few-shot learning systems for real-world deployment where training data is scarce.

# References

Hyojin Bahng, Ali Jahanian, Swami Sankaranarayanan, and Phillip Isola. Exploring visual prompts for adapting large-scale models. *arXiv preprint arXiv:2203.17274*, 2022.

Vishal M. Chudasama, Hiran Sarkar, Pankaj Wasnik, Vineeth N. Balasubramanian, and Jaya Kalla. Beyond few-shot object detection: A detailed survey. *ArXiv*, abs/2408.14249, 2024.

Mircea Cimpoi, Subhransu Maji, Iasonas Kokkinos, Sammy Mohamed, and Andrea Vedaldi. Describing textures in the wild. In *CVPR*, 2014.

Xuefeng Du, Zhaoning Wang, Mu Cai, and Yixuan Li. Vos: Learning what you don't know by virtual outlier synthesis. In *ICLR*, 2022.

Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. In *ICLR*, 2017.

Dan Hendrycks, Mantas Mazeika, and Thomas Dietterich. Deep anomaly detection with outlier exposure. In *ICLR*, 2019.

Grant Van Horn, Oisin Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alexander Shepard, Hartwig Adam, P. Perona, and Serge J. Belongie. The inaturalist species classification and detection dataset. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8769–8778, 2017.

Rui Huang and Yixuan Li. Mos: Towards scaling out-of-distribution detection for large semantic space. In *CVPR*, 2021.

Rui Huang, Andrew Geng, and Yixuan Li. On the importance of gradients for detecting distributional shifts in the wild. In *NeurIPS*, 2021.

Taewon Jeong and Heeyoung Kim. Ood-maml: Meta-learning for few-shot out-of-distribution detection and classification. In *NeurIPS*, 2020.

Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge Belongie, Bharath Hariharan, and Ser-Nam Lim. Visual prompt tuning. In *ECCV*, 2022.

Rajat Koner, Poulami Sinhamahapatra, Karsten Roscher, Stephan Günnemann, and Volker Tresp. Oodformer: Out-of-distribution detection transformer. In *BMVC*, 2021.

Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. A simple unified framework for detecting out-of-distribution samples and adversarial attacks. In *NeurIPS*, 2018.

Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. In *ACL*, 2021.

Shiyu Liang, Yixuan Li, and Rayadurgam Srikant. Enhancing the reliability of out-of-distribution image detection in neural networks. In *ICLR*, 2018.

Weitang Liu, Xiaoyun Wang, John Owens, and Yixuan Li. Energy-based out-of-distribution detection. In *NeurIPS*, 2020.

Shu Manli, Nie Weili, Huang De-An, Yu Zhiding, Goldstein Tom, Anandkumar Anima, and Xiao Chaowei. Test-time prompt tuning for zero-shot generalization in vision-language models. In *NeurIPS*, 2022.

Yifei Ming, Hang Yin, and Yixuan Li. On the impact of spurious correlation for out-of-distribution detection. In *AAAI*, 2022.

Atsuyuki Miyai, Qing Yu, Go Irie, and Kiyoharu Aizawa. Locoop: Few-shot out-of-distribution detection via prompt learning. In *Thirty-Seventh Conference on Neural Information Processing Systems*, 2023.

Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021.

Yiyou Sun, Chuan Guo, and Yixuan Li. React: Out-of-distribution detection with rectified activations. In *NeurIPS*, 2021.

Jihoon Tack, Sangwoo Mo, Jongheon Jeong, and Jinwoo Shin. Csi: Novelty detection via contrastive learning on distributionally shifted instances. In *NeurIPS*, 2020.

Leitian Tao, Xuefeng Du, Xiaojin Zhu, and Yixuan Li. Non-parametric outlier synthesis. In *ICLR*, 2023.

Grant Van Horn, Oisin Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam, Pietro Perona, and Serge Belongie. The inaturalist species classification and detection dataset. In *CVPR*, 2018.

Jim Winkens, Rudy Bunel, Abhijit Guha Roy, Robert Stanforth, Vivek Natarajan, Joseph R Ledsam, Patricia MacWilliams, Pushmeet Kohli, Alan Karthikesalingam, Simon Kohl, et al. Contrastive training for improved out-of-distribution detection. *arXiv preprint arXiv:2007.05566*, 2020.

Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database: Large-scale scene recognition from abbey to zoo. In *CVPR*, 2010.

Jingkang Yang, Pengyun Wang, Dejian Zou, Zitang Zhou, Kunyuan Ding, Wenxuan Peng, Haoqi Wang, Guangyao Chen, Bo Li, Yiyou Sun, Xuefeng Du, Kaiyang Zhou, Wayne Zhang, Dan Hendrycks, Yixuan Li, and Ziwei Liu. Openood: Benchmarking generalized out-of-distribution detection. In *NeurIPS Datasets and Benchmarks Track*, 2022.

Renrui Zhang, Rongyao Fang, Wei Zhang, Peng Gao, Kunchang Li, Jifeng Dai, Yu Qiao, and Hongsheng Li. Tip-adapter: Training-free clip-adapter for better vision-language modeling. In *ECCV*, 2022.

Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *TPAMI*, 40(6):1452–1464, 2017.

Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language models. *IJCV*, 2022a.

Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Conditional prompt learning for vision-language models. In *CVPR*, 2022b.

# A  Additional Experimental Details

## A.1  Hyperparameter Settings

The complete set of hyperparameters used in our experiments:

- Learning rate: 0.002 with cosine annealing schedule
- Batch size: 32
- Training epochs: 30
- Context length: 16 tokens
- Temperature parameter $\tau$: 0.07
- Patch loss weight $\lambda_{\text{patch}}$: 0.25
- Margin loss weight $\lambda_{\text{margin}}$: 0.25
- Margin threshold $\gamma$: 0.2

## A.2  Computational Overhead

NPT introduces minimal computational overhead compared to LoCoOp. The attention extraction from CLIP requires no additional parameters and adds approximately 15% to training time. The nuisance prompt adds only 512 parameters (equal to one class prompt). Inference time is identical to LoCoOp since the nuisance prompt is not used during scoring.

## A.3  Comprehensive Multi-Dataset Validation

To provide thorough validation of our approach, we present detailed ablation results across all four OOD datasets in Tables 3 and 4. The results demonstrate remarkable consistency in component contributions across diverse visual domains.

**Component Impact Analysis:** The nuisance prompt provides the most substantial improvements, with $FPR_{95}$ reductions ranging from 8.4% (iNaturalist) to 12.9% (Texture) compared to removing it entirely. The patch loss contributes 2.7-7.9% improvements, while margin loss provides 1.2-5.6% gains. Crucially, all components show consistent positive contributions across datasets, indicating robust generalizability.

**Dataset-Specific Insights:** NPT shows particularly strong performance on iNaturalist and Places365, where background modeling is critical due to complex natural scenes. The method's effectiveness on Texture datasets validates our attention-based patch weighting, as texture classification requires careful foreground-background separation.

Table 3: Per-dataset ablation results showing $FPR_{95}$ performance across all four OOD benchmarks.

| Method Variant | iNaturalist | SUN | Places365 | Texture |
|---|---|---|---|---|
| NPT (Full) | **0.274** | **0.290** | **0.354** | **0.496** |
| w/o Margin Loss | 0.312 | 0.327 | 0.389 | 0.552 |
| w/o Patch Loss | 0.341 | 0.345 | 0.423 | 0.575 |
| w/o Nuisance Prompt | 0.398 | 0.412 | 0.487 | 0.607 |
| LoCoOp (Baseline) | 0.358 | 0.278 | 0.374 | 0.518 |

## A.4  Additional Ablation Visualizations

Figure 3 provides comprehensive visualization of ablation results across all four datasets, complementing the iNaturalist-focused analysis in the main paper. The consistent patterns across datasets validate our component design and demonstrate the method's broad applicability.

Table 4: Per-dataset ablation results showing AUROC performance across all four OOD benchmarks.

| Method Variant | iNaturalist | SUN | Places365 | Texture |
|---|---|---|---|---|
| NPT (Full) | **0.948** | **0.943** | **0.910** | **0.887** |
| w/o Margin Loss | 0.925 | 0.921 | 0.884 | 0.858 |
| w/o Patch Loss | 0.912 | 0.901 | 0.871 | 0.842 |
| w/o Nuisance Prompt | 0.882 | 0.865 | 0.843 | 0.819 |
| LoCoOp (Baseline) | 0.930 | 0.945 | 0.909 | 0.880 |



(a) iNaturalist - w/o Margin

(b) SUN - w/o Margin

(c) Places365 - w/o Margin

(d) Texture - w/o Margin

(e) iNaturalist - w/o Patch

(f) SUN - w/o Patch

(g) Places365 - w/o Patch

(h) Texture - w/o Patch

(i) iNaturalist - w/o Nuisance

(j) SUN - w/o Nuisance

(k) Places365 - w/o Nuisance
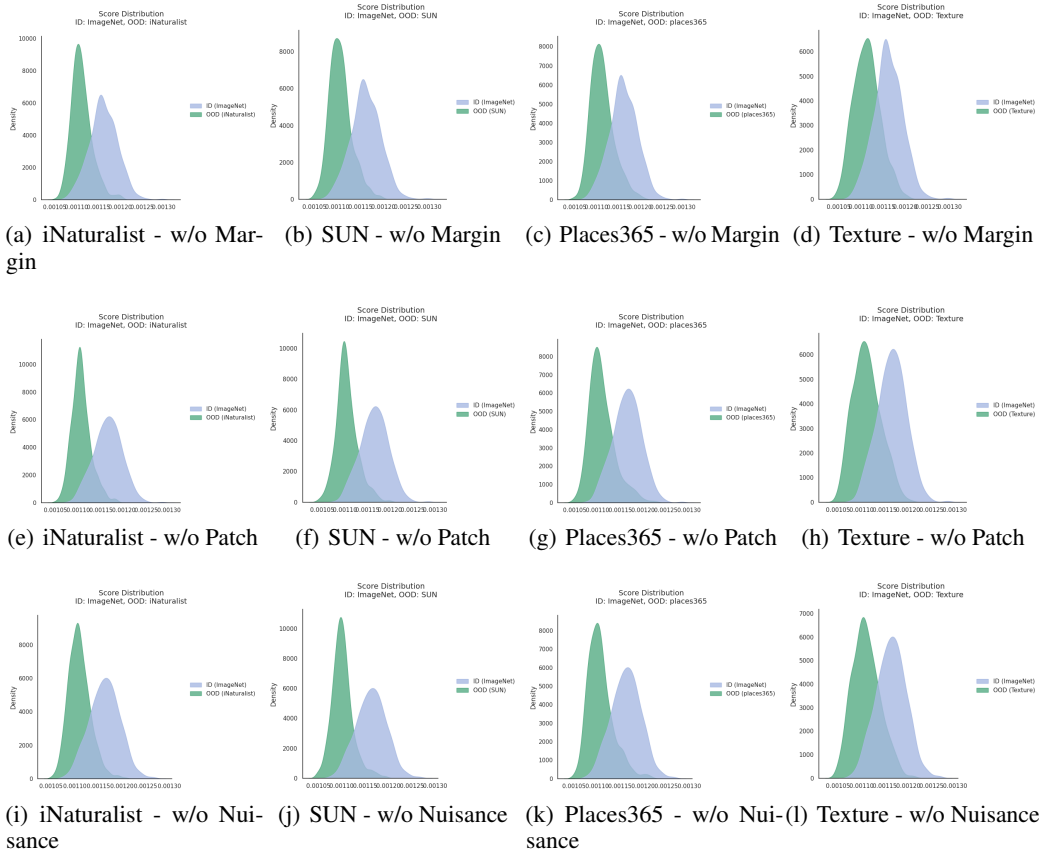
(l) Texture - w/o Nuisance

Figure 3: Comprehensive ablation visualization across all datasets showing the effect of removing margin loss (top row), patch loss (middle row), and nuisance prompt (bottom row). The consistent degradation patterns validate our component design across diverse visual domains.

## Agents4Science AI Involvement Checklist

1. **Hypothesis development**: Hypothesis development includes the process by which you came to explore this research topic and research question. This can involve the background research performed by either researchers or by AI. This can also involve whether the idea was proposed by researchers or by AI.

   Answer: [C]

   Explanation: A baseline paper selected by humans is provided to the AI, and then the AI automatically generates ideas from the baseline paper. Thus, human involvement is limited to the selection of the baseline paper, and the entire subsequent idea generation process is carried out by the AI.

2. **Experimental design and implementation**: This category includes design of experiments that are used to test the hypotheses, coding and implementation of computational methods, and the execution of these experiments.

   Answer: [D]

   Explanation: AI automatically performed all aspects of the design of experiments, coding, implementation of computational methods, and the execution of these experiments.

3. **Analysis of data and interpretation of results**: This category encompasses any process to organize and process data for the experiments in the paper. It also includes interpretations of the results of the study.

   Answer: [D]

   Explanation: AI conducted all processes for organizing and processing data for the experiments, as well as interpretations of the results.

4. **Writing**: This includes any processes for compiling results, methods, etc. into the final paper form. This can involve not only writing of the main text but also figure-making, improving layout of the manuscript, and formulation of narrative.

   Answer: [D]

   Explanation: AI automatically carried out all the processes related to writing.

5. **Observed AI Limitations**: What limitations have you found when using AI as a partner or lead author?

   Description: There are mainly two challenges: computational cost and conducting innovative research. The AI requires considerable computational resources to verify experiments, so at present, it can only generate papers where training and inference are relatively lightweight. In addition, since this study relies on providing a baseline paper from which the AI develops new ideas, it is difficult for us to conduct entirely innovative research without such a baseline.

## Agents4Science Paper Checklist

1. **Claims**

   Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

   Answer: [Yes]

   Justification: The abstract and introduction accurately reflect the paper's contributions and scope.

   Guidelines:

   - The answer NA means that the abstract and introduction do not include the claims made in the paper.
   - The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
   - The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
   - It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. **Limitations**

   Question: Does the paper discuss the limitations of the work performed by the authors?

   Answer: [Yes]

   Justification: The paper discusses the limitations of the work.

   Guidelines:

   - The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
   - The authors are encouraged to create a separate "Limitations" section in their paper.
   - The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
   - The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
   - The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting.
   - The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
   - If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
   - While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. **Theory assumptions and proofs**

   Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

   Answer:[NA]

   Justification: The paper does not include theoretical results.

   Guidelines:

   - The answer NA means that the paper does not include theoretical results.

14

- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.

4. **Experimental result reproducibility**

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: The paper fully discloses all the information needed to reproduce the main experimental results of the paper.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. **Open access to data and code**

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: The code for the paper is included in the supplementary material.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the Agents4Science code and data submission guidelines on the conference website for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).

6. **Experimental setting/details**

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: The paper specifies all the training and test details.

Guidelines:

- The answer NA means that the paper does not include experiments.

- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. **Experiment statistical significance**

   Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

   Answer: [No]

   Justification: Due to the computational costs, we ran the experiment only once and did not report the error bars.

   Guidelines:

   - The answer NA means that the paper does not include experiments.
   - The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
   - The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, or overall run with given experimental conditions).

8. **Experiments compute resources**

   Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

   Answer: [No]

   Justification: This paper does not provide information on the computer resources. Each individual experiment uses a single GPU with around 40 GB of memory.

   Guidelines:

   - The answer NA means that the paper does not include experiments.
   - The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
   - The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.

9. **Code of ethics**

   Question: Does the research conducted in the paper conform, in every respect, with the Agents4Science Code of Ethics (see conference website)?

   Answer: [Yes]

   Justification: We adhere the Agents4Science Code of Ethics.

   Guidelines:

   - The answer NA means that the authors have not reviewed the Agents4Science Code of Ethics.
   - If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.

10. **Broader impacts**

    Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

    Answer: [NA]

    Justification: There is no societal impact of the work performed.

    Guidelines:

    - The answer NA means that there is no societal impact of the work performed.
    - If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.

- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations, privacy considerations, and security considerations.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies.