

MRCert: Towards Post-deployment Patch Robustness Certification for Adversarially Patched Samples via Type-specific Masking

795 A Formal Proof

796 **Property** (Infeasibility of Re-certifying Adversarially
 797 Patched Samples for Existing Masking-based Recovery
 798 Works). $\forall x' \in \mathbb{A}_{\mathbb{P}}(x), c_r(x) \wedge f(x') \neq f(x) \implies \neg c_r(x')$.

799 *Proof.* We know in [Xiang *et al.*, 2022], $c_r(\hat{x}) :=$
 800 $[\forall M_1, M_2 \in \mathbb{M}_{\mathbb{P}}, f(\hat{x} \odot M_1 \odot M_2) = f(\hat{x})]$. Then if $c_r(x) =$
 801 *True*, we know $[\forall M_1, M_2 \in \mathbb{M}_{\mathbb{P}}, f(x \odot M_1 \odot M_2) = f(x)]$.
 802 We also know $\forall P \in \mathbb{P}, \exists M \in \mathbb{M}_{\mathbb{P}}, M \odot P = P$ by definition,
 803 which makes $\forall x' \in \mathbb{A}_{\mathbb{P}}(x), x' \odot M = x \odot M$ (see Fig. 2).
 804 Therefore, $\exists x' \in \mathbb{A}_{\mathbb{P}}(x), f(x' \odot M) = f(x \odot M) = f(x)$.
 805 We have $f(x') \neq f(x)$ from the antecedent; therefore,
 806 $\exists x' \in \mathbb{A}_{\mathbb{P}}(x), f(x' \odot M) \neq f(x')$. Since $c_r(x')$ requires
 807 $[\forall M_1, M_2 \in \mathbb{M}_{\mathbb{P}}, f(x' \odot M_1 \odot M_2) = f(x')]$, $c_r(x') = \text{False}$
 808 holds. \square

809 **Lemma** (Certification of benign samples). *Given an arbitrary sample \hat{x} , if $[\forall M_1, M_2, M_3 \in \mathbb{M}_{\mathbb{P}}, f(\hat{x} \odot M_1 \odot M_2 \odot M_3) = f(\hat{x})] \text{ holds, } \forall \hat{x}' \in \mathbb{A}_{\mathbb{P}}(\hat{x}), g(\hat{x}') = g(\hat{x})$.*

812 *Proof.* We first analyze $g(\hat{x})$. Since $[\forall M_1, M_2, M_3 \in \mathbb{M}_{\mathbb{P}}, f(\hat{x} \odot M_1 \odot M_2 \odot M_3) = f(\hat{x})] \implies [\forall M_1, M_2 \in \mathbb{M}_{\mathbb{P}}, f(\hat{x} \odot M_1 \odot M_2) = f(\hat{x})]$, we know its prediction label $g(\hat{x}) = f(\hat{x})$ output in Case ①. We then analyze $g(\hat{x}')$. For $\forall \hat{x}' \in \{\hat{x}' \mid \hat{x}' = (J - P) \odot \hat{x} + P \odot \hat{x}' \wedge P \in \mathbb{P}\}$ (i.e., $\forall \hat{x}' \in \mathbb{A}_{\mathbb{P}}(\hat{x})$), w.l.o.g., we let $M_1 \odot P = O$. Therefore, we get $\hat{x}' \odot M_1 = ((J - P) \odot \hat{x} + P \odot \hat{x}') \odot M_1 = \hat{x} \odot M_1$. Then, we get $[\exists M_1 \in \mathbb{M}_{\mathbb{P}}, \forall M_2, M_3 \in \mathbb{M}_{\mathbb{P}}, f(\hat{x}' \odot M_1 \odot M_2 \odot M_3) = f(\hat{x})]$ (see Fig. 2). Note that the special cases $M_2 = M_3, M_1 = M_2 = M_3$ are included. Case 1: Suppose the returned label of \hat{x}' output in Case ① as $f(\hat{x}')$ (i.e., $\forall M_1, M_2 \in \mathbb{M}_{\mathbb{P}}, f(\hat{x}' \odot M_1 \odot M_2) = f(\hat{x}')$), then we know $g(\hat{x}') = f(\hat{x}') = f(\hat{x}) = g(\hat{x})$. Case 2: Otherwise, its returned label should be output in Case ② as $f(\hat{x}' \odot M_1)$, since $[\exists M_1 \in \mathbb{M}_{\mathbb{P}}, \forall M_2, M_3 \in \mathbb{M}_{\mathbb{P}}, f(\hat{x}' \odot M_1 \odot M_2 \odot M_3) = f(\hat{x}' \odot M_1) = f(\hat{x})]$, and further $g(\hat{x}') = f(\hat{x}' \odot M_1) = f(\hat{x}) = g(\hat{x})$. \square

829 **Lemma** (Certification of adversarially patched samples).
 830 *Given an arbitrary sample \hat{x} , if $[\exists M_1 \in \mathbb{M}_{\mathbb{P}}, \forall M_2, M_3, M_4 \in \mathbb{M}_{\mathbb{P}}, f(\hat{x} \odot M_1 \odot M_2 \odot M_3 \odot M_4) = f(\hat{x} \odot M_1)] \text{ holds, } \forall \hat{x}' \in \mathbb{A}_{\mathbb{P}}(\hat{x}), g(\hat{x}') = g(\hat{x})$.*

833 *Proof.* We first analyze $g(\hat{x})$. Case 1: If \hat{x} meet the condition of Case ①, we can further get $\forall M_1, M_2, M_3, M_4 \in \mathbb{M}_{\mathbb{P}}, f(\hat{x} \odot M_1 \odot M_2 \odot M_3 \odot M_4) = f(\hat{x})$. Note that the special case $M_1 = M_2 = M_3 = M_4$ is included, which means we can get $g(\hat{x}) = f(\hat{x} \odot M_1)$. Case 2: Otherwise, by the given condition $[\exists M_1 \in \mathbb{M}_{\mathbb{P}}, \forall M_2, M_3, M_4 \in \mathbb{M}_{\mathbb{P}}, f(\hat{x} \odot M_1 \odot M_2 \odot M_3 \odot M_4) = f(\hat{x} \odot M_1)] \implies [\exists M_1 \in \mathbb{M}_{\mathbb{P}}, \forall M_2, M_3 \in \mathbb{M}_{\mathbb{P}}, f(\hat{x} \odot M_1 \odot M_2 \odot M_3) = f(\hat{x} \odot M_1)]$, \hat{x} meet the condition in Case ② and its prediction label $g(\hat{x}) = f(\hat{x} \odot M_1)$, same as Case ①. We then analyze $g(\hat{x}')$. For $\hat{x}' \in \{\hat{x}' \mid \hat{x}' = (J - P) \odot \hat{x} + P \odot \hat{x}' \wedge P \in \mathbb{P}\}$ (i.e.,

$\forall \hat{x}' \in \mathbb{A}_{\mathbb{P}}(\hat{x})$), Case 1: Suppose $M_1 \odot P = O$. Then we can get $\hat{x}' \odot M_1 = ((J - P) \odot \hat{x} + P \odot \hat{x}') \odot M_1 = \hat{x} \odot M_1$, and further get $[\forall M_2, M_3, M_4 \in \mathbb{M}_{\mathbb{P}}, f(\hat{x}' \odot M_1 \odot M_2 \odot M_3 \odot M_4) = f(\hat{x} \odot M_1)]$, which is the same condition as that on \hat{x} . Therefore, repeating those analysis for $g(\hat{x})$ above can get $g(\hat{x}) = g(\hat{x}')$. Case 2: Otherwise, for M_2, M_3, M_4 , w.l.o.g., we let $M_2 \odot P = O$ ($M_1 \neq M_2$). Then similarly, we get $[\exists M_1, M_2 (\neq M_1) \in \mathbb{M}_{\mathbb{P}}, \forall M_3, M_4 \in \mathbb{M}_{\mathbb{P}}, f(\hat{x}' \odot M_1 \odot M_2 \odot M_3 \odot M_4) = f(\hat{x} \odot M_1)]$. Note that the special cases $M_1 = M_3, M_2 = M_4, M_3 = M_4$ are included. Case 2.1: Suppose the prediction label $g(\hat{x}')$ output in Case ①. Then, since $[\exists M_1, M_2 (\neq M_1) \in \mathbb{M}_{\mathbb{P}}, f(\hat{x}' \odot M_1 \odot M_2) = f(\hat{x} \odot M_1)]$ (special case $M_1 = M_3, M_2 = M_4$), by the condition of Case ①, we know $g(\hat{x}') = f(\hat{x}') = f(\hat{x} \odot M_1) = g(\hat{x})$. Case 2.2: Suppose the prediction label $g(\hat{x}')$ output in Case ②. Then, since $[\exists M_1, M_2 (\neq M_1) \in \mathbb{M}_{\mathbb{P}}, \forall M_3 \in \mathbb{M}_{\mathbb{P}}, f(\hat{x}' \odot M_1 \odot M_2 \odot M_3) = f(\hat{x} \odot M_1)]$ (special case $M_3 = M_4$), by the condition of Case ②, we know $g(\hat{x}') = f(\hat{x}' \odot M_1) = f(\hat{x} \odot M_1) = g(\hat{x})$. Case 2.3: Otherwise, the prediction label of \hat{x}' should output in Case ③ since $[\exists M_1, M_2 (\neq M_1) \in \mathbb{M}_{\mathbb{P}}, \forall M_3, M_4 \in \mathbb{M}_{\mathbb{P}}, f(\hat{x}' \odot M_1 \odot M_2 \odot M_3 \odot M_4) = f(\hat{x} \odot M_1) = f(\hat{x}' \odot M_1 \odot M_2)]$ (special case $M_1 = M_3, M_2 = M_4$), and further $g(\hat{x}') = f(\hat{x}' \odot M_1 \odot M_2) = f(\hat{x} \odot M_1) = g(\hat{x})$. \square

Theorem (Certification of samples). *Given an arbitrary sample \hat{x} , if $c(\hat{x}) = \text{True}$ holds, $\forall \hat{x}' \in \mathbb{A}_{\mathbb{P}}(\hat{x}), g(\hat{x}') = g(\hat{x})$.*

Simply conjoining the antecedent of Lemma 1 and Lemma 2 can prove this theorem.

Theorem (Round-trip certification of samples). *Given a benign sample x , if $[\forall M_1, M_2, M_3, M_4 \in \mathbb{M}_{\mathbb{P}}, f(x \odot M_1 \odot M_2 \odot M_3 \odot M_4) = f(x)]$ (i.e., $c_r^2(x) = \text{True}$), $[\forall x' \in \mathbb{A}_{\mathbb{P}}(x), g(x') = g(x) \wedge c_r(x') = \text{True} \wedge c_r(x) = \text{True}]$.*

Proof. By the condition $[\forall M_1, M_2, M_3, M_4 \in \mathbb{M}_{\mathbb{P}}, f(x \odot M_1 \odot M_2 \odot M_3 \odot M_4) = f(x)]$, we know the returned label $g(x) = f(x)$ in Case ① and $c_r(x) = \text{True}$. Still by this condition, we know $[\forall x' \in \mathbb{A}_{\mathbb{P}}(x), \exists M_1 \in \mathbb{M}_{\mathbb{P}}, \forall M_2, M_3, M_4 \in \mathbb{M}_{\mathbb{P}}, f(x' \odot M_1 \odot M_2 \odot M_3 \odot M_4) = f(x)]$, since $[\exists M_1 \in \mathbb{M}_{\mathbb{P}}, \forall x' \in \mathbb{A}_{\mathbb{P}}(x), x' \odot M_1 = x \odot M_1]$ (see Fig. 2 for illustration). Then by Lemma 2, we know $[\forall x' \in \mathbb{A}_{\mathbb{P}}(x), c_r(x') = \text{True}]$. By Lemma 1, we also know $[\forall x' \in \mathbb{A}_{\mathbb{P}}(x), g(x') = g(x)]$. Finally, we know $[\forall x' \in \mathbb{A}_{\mathbb{P}}(x), g(x') = g(x) \wedge c_r(x') = \text{True} \wedge c_r(x) = \text{True}]$. \square

B The Subtlety in the Design of MRCert

The ordering of the type-specific cases in g . The sequence of the cases in g is crucial. Changing the order will either break the certification guarantee for adversarially patched samples or degrade the clean accuracy. For example, swapping Case ② and Case ③ would require the benign counterpart of an adversarially patched sample to exhibit the same prediction labels under four masks (instead of three) to remain certified; otherwise, a patch attacker could exploit a

894 fourth-order masked mutant that predicts another label to
895 evade the detection of the adversarially patched sample by
896 this case-swapped version of g . In our proposed version of g ,
897 Case 1 is intentionally placed first to prevent benign samples
898 from reaching the recovery stages, thereby improving clean
899 accuracy.

900 *The ordering of mask conditions within a type-specific case*
901 *in g .* The internal ordering of masking conditions also
902 matters. Consider Case 2 (indirect testing): Suppose that its
903 quantifiers were reordered from “there exists mask 1 such
904 that for all mask 2, mask3...” to “there exist mask 1 and
905 mask 2 such that for all mask 3 ...”. Then, adversarial sam-
906 ples containing two patches could escape certification, even
907 when their benign counterparts are round-trip certified. In
908 this reordered structure, the adversarial patched sample could
909 simply choose two masks (mask 1 and mask 2) that expose
910 both patches, i.e., they do not cover either patch. Then, even
911 if mask 3 covers one patch, the other patch remains uncov-
912 ered, allowing the attacker to arbitrarily alter the prediction
913 and evade detection. This shows that the original quantifier
914 ordering is essential for preventing such escape cases.

915 *The number of masks within a type-specific cases in g .* Tak-
916 ing Case 2 as an example, the number of masks is also critical.
917 Reducing the number of masks in Case 2 would allow adver-
918 sarial samples with two patches to escape from the detection
919 of indirect testing, for reasons analogous to the ordering is-
920 sue discussed above. Conversely, increasing the number of
921 masks (e.g., by 1) would increase the masking requirements
922 for certification, which makes certifying a benign sample re-
923 quire prediction consistency under four masks, instead of the
924 current three, which directly decreases the total number of
925 samples that can be certified.

926 C Extension of MRCert to Recover and 927 Certify Samples with N-patches

928 We can extend the maximum number of patches from 2 to N
929 (called MRCert-N-patch, a variant of MRCert) following the
930 following idea. First, we apply each set of N masks in the
931 covering mask set \mathbb{M} on the input sample \hat{x} to test whether
932 \hat{x} is harmful. If it is not harmful, MRCert-N-patch returns
933 the label $f(\hat{x})$ (marked as N-Case ①). If \hat{x} is detected as
934 harmful, we then test whether all its first-order mutants are
935 harmful by applying each possible subset with N masks se-
936 lected with replacement from \mathbb{M} on each first-order mutant of
937 \hat{x} . If there exists a first-order mutant, whose all $(N + 1)$ th-
938 order mutants generated from the first-order mutant of \hat{x} are
939 predicted with the same label as this first-order sample, then
940 \hat{x} is deemed as a one-patch harmful sample, this first-order
941 mutant is “clean” and the prediction label of this first-order
942 mutant is returned (marked as N-Case ②). If that is not the
943 case, we then test whether all second-order mutants of \hat{x} are
944 harmful in the same manner, and repeat until the N th-order
945 mutants of \hat{x} are tested. For the certification function c_r with
946 the input sample \hat{x} , it should be extended to the condition that
947 all $(N + 1)$ th-order mutants of \hat{x} are predicted with the same
948 label as \hat{x} (for those input samples whose label returned in
949 Case ①), and the condition that there exists a first-order mu-
950 tant of \hat{x} , whose all $(N + 2)$ th-order mutants are predicted

951 with the same label as this first-order mutant of \hat{x} (for those
952 input samples whose label returned in Case ②), and certify-
953 ing the input samples output in other cases by the condition
954 in the same manner. For the round-trip certification function,
955 it should be the condition that all $2N$ th-order mutants of a
956 benign sample x are predicted with the same label as x . We
957 leave the formal proof and implementation as future work.

958 D Experimental Setup

959 D.1 Environment and Dataset

960 The evaluation is conducted on an Ubuntu 20.04 machine
961 equipped with four Nvidia 3090 GPUs. Following [Xiang *et*
962 *al.*, 2022], we adopt 1000-class ImageNet [Deng *et al.*, 2009],
963 10-class CIFAR10 [Krizhevsky *et al.*, 2009], and 10-class Im-
964 ageNette as our datasets, encompassing both large-scale, di-
965 verse datasets and efficient benchmarks widely used for rapid
966 experimentation.

967 D.2 Models and Baselines

968 We adopt the Vision Transformer (ViT) as our base model,
969 which achieves the state-of-the-art in many patch robustness
970 certification defenders [Xiang *et al.*, 2024; Li *et al.*, 2022;
971 Salman *et al.*, 2022; Xiang *et al.*, 2022; Huang *et al.*, 2023].

972 We adopt the state-of-the-art masking-based defender
973 PatchCURE (**PC**) [Xiang *et al.*, 2024] and state-of-the-
974 art smoothing-based defender VOT-CrossCert (**VOT**) from
975 CrossCert [Zhou *et al.*, 2024], which shares the common
976 methodology for smoothing-based recovery with ViP [Li *et*
977 *al.*, 2022] and S-ViT [Salman *et al.*, 2022]. We arm VOT with
978 the methodology for runtime verification proposed by [Yat-
979 sura *et al.*, 2023] as our baseline to compare with **MRCert** in
980 our infrastructure.

981 Specifically, MRCert adopts the end-to-end ViT-SRF
982 model (setting 14x1-k6¹) proposed and used by PatchCURE
983 [Xiang *et al.*, 2024] from [https://github.com/inspire-group/
984 PatchCURE](https://github.com/inspire-group/PatchCURE) with the pre-train weights and setting from MAE
985 [He *et al.*, 2022], so that MRCert can use the same strat-
986 egy in the experiment as PatchCURE to calculate the cov-
987 ering mask set, generate mutants with their predictions in-
988 side the ViT-SRF. We use the training scripts with the set-
989 tings and hyperparameters from the official repository of
990 PatchCURE for fine-tuning ViT-SRF, which uses the train-
991 ing samples with their first and second-order mutants only.
992 We then apply the same finetuned ViT-SRF to both MRCert
993 and PatchCURE to ensure fairness. We also adopt this offi-
994 cial repository <https://github.com/inspire-group/PatchCURE>
995 to implement PatchCURE in our infrastructure. We adopt the
996 same pre-trained weights and settings from MAE [He *et al.*,
997 2022] for VOT. We follow the fine-tuning settings and hy-
998 perparameters from the state-of-the-art S-ViT [Salman *et al.*,

¹Other variants ViT-SRF14x2, ViT-SRF2x2, and BagNet33 can-
not handle the situation against two patches, since all their receptive
fields would be inherently masked by VIT-SRF [Xiang *et al.*, 2024]
when generating corresponding mutants. k is the parameter to con-
trol the position of splitting SRF and LRF in VIT-SRF in the range
[0, 12], which can tune the trade-off between computational effi-
ciency and robustness. We adopt the middle one since it is not the
focus of this paper.

2022] for VOT (note that smoothing-based recovery adopts
 a different notion of mutants, where a mutant is a slice of a
 sample; therefore, a different fine-tuning on the base model
 is needed [Li *et al.*, 2022]). We also adopt the column abla-
 tion with the ablation size of 19 pixels from [Salman *et al.*,
 2022]. We have extracted the corresponding results from the
 original papers of ViP [Li *et al.*, 2022] and S-ViT [Salman
et al., 2022] to compare with VOT in Tab. 1 (their results on
 round-trip certification are not reported in their papers).

1008 D.3 Metrics

1009 Let x be a benign sample with the true label y_0 in a clean
 1010 test dataset \mathbb{D} , and $R = \langle g(x), c(x), c_r^2(x) \rangle$ be a certi-
 1011 fied recovery defender. **Clean accuracy** is the fraction
 1012 of \mathbb{D} that are correctly predicted, defined as $acc_{clean} =$
 1013 $\frac{|\{x \in \mathbb{D} | g(x) = y_0\}|}{|\mathbb{D}|}$, which evaluates the standard performance
 1014 of certified recovery as a classifier. **Certified accuracy** is
 1015 the fraction of \mathbb{D} that are correctly predicted and certified
 1016 robust, whose adversarially patched samples should be pre-
 1017 dicted with the same label as the benign sample, defined
 1018 as $acc_{cert} = \frac{|\{x \in \mathbb{D} | g(x) = y_0 \wedge c_r(x) = True\}|}{|\mathbb{D}|}$. **Round-trip cer-**
 1019 **tified accuracy** is the fraction of \mathbb{D} that are correctly pre-
 1020 dicted and round-trip certified, whose adversarially patched
 1021 samples should be predicted with the same label as the be-
 1022 nign sample and gain a provable robust verdict, defined
 1023 as $acc_{cert^2} = \frac{|\{x \in \mathbb{D} | g(x) = y_0 \wedge c_r^2(x) = True\}|}{|\mathbb{D}|}$. The clean accu-
 1024 racy of the base classification model f is $\frac{|\{x \in \mathbb{D} | f(x) = y_0\}|}{|\mathbb{D}|}$.
 1025 Clean accuracy and certified accuracy are commonly adopted
 1026 by peers [Levine and Feizi, 2020a; Salman *et al.*, 2022;
 1027 Li *et al.*, 2022; Xiang *et al.*, 2024; Xiang *et al.*, 2022;
 1028 Xiang *et al.*, 2021], and round-trip certified accuracy is pro-
 1029 posed by [Yatsura *et al.*, 2023].

1030 In RQ1, the test dataset \mathbb{D} is a set of patched sam-
 1031 ples from attacking a subset of the test dataset (valida-
 1032 tion dataset for ImageNet), denoted as \mathbb{D}_{pat} . Clean ac-
 1033 curacy and certified accuracy are measured in this set of
 1034 patched samples to evaluate certified recovery defenders
 1035 on their ability to recover predictions and recover predic-
 1036 tions with a provably robust verdict, respectively. We mea-
 1037 sure each defender’s ability to certify adversarially patched
 1038 samples by **adversarial certified accuracy** $acc_{cert}^{adv} =$
 1039 $|\{\hat{x} \in \mathbb{D}_{pat} | g(\hat{x}) = y_0 \wedge c_r(\hat{x}) = True \wedge f(\hat{x}) \neq y_0\}| / |\{\hat{x} \in \mathbb{D}_{pat} | f(\hat{x}) \neq y_0\}|$. In
 1040 RQ2, we use the original test dataset (validation dataset for
 1041 ImageNet) as \mathbb{D} .

1042 D.4 Experimental Procedure

1043 We adopt the patch sizes of 16 pixels (measured as a square
 1044 patch region with a side length of 16 pixels) and 32 pixels for
 1045 all three datasets. All samples are rescaled to 224x224 [Zhou
 1046 *et al.*, 2024; Li *et al.*, 2022; Xiang *et al.*, 2022; Xiang *et al.*,
 1047 2024].

1048 In RQ1, we perform an actual adversarial patch attack
 1049 IFGSM adopted from [Levine and Feizi, 2020a] on PC and
 1050 MRCert, which uses their shared base model as the attack
 1051 model for gradient-based attacks without the knowledge of
 1052 the non-differentiable label recovery function for fairness.

Note that the recent proposed attacks are more toward prac-
 1053 tical, such as limiting the access times or can only get the
 1054 returned prediction label, while our IFGSM attack has direct
 1055 and full access to the base model, which is a powerful exam-
 1056 against attackers. We set 80 random starts, 150 iterations per
 1057 random start, and a step size of 0.05 following [Levine and
 1058 Feizi, 2020a]. We randomly select 500 test samples from
 1059 each test dataset for attacks. We follow the practice in [Levine
 1060 and Feizi, 2020a] to return the worst patched sample for each
 1061 benign sample and place it into \mathbb{D}_{pat} . In RQ2, we directly
 1062 adopt the entire test dataset (validation dataset for ImageNet)
 1063 as \mathbb{D} .

In both RQ1 (as a follow-up to the attack) and RQ2, for
 1065 each defender on each sample in each test dataset, we first
 1066 generate and evaluate corresponding mutants using the corre-
 1067 sponding base model (PC and MRCert share the same ones),
 1068 then apply each defender to the prediction results collected
 1069 for these mutants, and measure the metric values.
 1070