
A Theory of Transfer-Based Black-Box Attacks: Explanation and Implications (Supplementary Material)

Anonymous Author(s)

Affiliation

Address

email

1 This article serves as the supplementary material to the central part of our paper. Appendix A includes
2 some further discussions. Complete proofs of the theorems and propositions in Sections 4 and 5 can
3 be found in Appendix B. A multi-class analysis of the manifold attack model is given in Appendix C.

4 A Further Discussions

5 A.1 What makes a good explanatory model?

6 As its title suggests, our paper’s primary effort is to explain the properties of TBAs by the manifold
7 attack model. During the writing of this paper, the following question is discussed repeatedly:

8 *What makes a good explanatory model and how to evaluate an explanatory model?*

9 This subsection provides our answer to this question. First of all, we believe that a good explanatory
10 model should be:

- 11 • **(Criterion 1)** consistent with existing empirical results,
- 12 • **(Criterion 2)** based on reasonable assumptions, and
- 13 • **(Criterion 3)** theoretically tractable.

14 Throughout this paper, we make many efforts to validate our model. Specifically, we try to check
15 whether our model fulfills criteria 1-3. Clearly, our model is theoretically tractable. We theoretically
16 analyze TBAs and provide many explanatory results in Sections 4 and 5.

17 In the rest of this subsection, we briefly discuss criteria 1 and 2. For the first criterion, we discuss the
18 intriguing properties of TBAs (i.e., the empirical results observed by previous works) in Sections
19 1 and 2. Two of the most widely-known properties of TBAs are: 1) TBAs can craft transferable
20 adversarial examples even when the source model is inaccurate [1], and 2) the success rates of
21 TBAs are constantly lower than other methods of black-box adversarial attacks [2–4]. Section 4
22 demonstrates that our model is consistent with the existing empirical results and provides reasonable
23 explanations for these properties.

24 As for criterion 2, our model assumes that the natural data lies on a low-dimensional manifold. This
25 assumption is commonly seen in previous works [5–8]. We also assume that the classifiers (i.e.,
26 the source and target models in TBAs) can be decomposed into the product of a semantic classifier
27 f_b (Definition 4.1) and a concentration multiplier ϕ (Definition 4.2). This assumption is based on
28 the empirical observation that *ML models can capture semantic and geometrical information of the*
29 *natural data* [9, 10]. Here, our concerns are two folds: 1) what are the semantic and geometrical
30 information, and 2) how does an ML model capture such information?

31 **Semantic information** We first focus on semantic information. The following remark explains
 32 what is the semantic information of a dataset by an example.

33 *Remark A.1* (The semantic information of CIFAR-10). Generally speaking, "semantic" refers to
 34 the relationship between natural data and their true label, which should be consistent with human
 35 recognition. For example, the semantic information contained in the CIFAR-10 dataset is the true
 36 labels (e.g., airplane, automobile, and bird) and their corresponding natural images (e.g., images of
 37 airliners, SUVs, and chickens). In this example, an image cannot simultaneously include an airplane
 38 and an automobile since "the classes are completely mutually exclusive" in the CIFAR-10 dataset,
 39 cf. the official website of CIFAR-10. That is, the semantic information provided by CIFAR-10 is
 40 separated. ▲

41 In our paper, we formalize the semantic information of natural data by separated sets $A^1, A^2, \dots, A^k \subset$
 42 \mathcal{M} (for a k -class classification task), see Section 3.2 for the definitions. As is discussed in Remark A.1,
 43 these sets reflect the relationship between true labels and their corresponding natural data, and more
 44 importantly, these sets should be separated. In this paper, we define separated sets in Definition 3.2
 45 and assume that $A^1, A^2, \dots, A^k \subset \mathcal{M}$ are separated. It is worth noting that the definition of "semantic
 46 information" in our paper is motivated by that of the "concept" in classical learning theory [11, 12].
 47 In these works, learning a concept is equivalent to approximating the decision boundary of ML
 48 models to the concept sets (i.e., subsets in the sample space).

49 Our model captures the semantic information in a similar way as [11]. We let $A_f^1, A_f^2, \dots, A_f^k \subset \mathcal{M}$ be
 50 the semantic information learned by f . Note that we do not assume these sets to be regions or to have
 51 any compactness or connectedness restriction. Instead, we only assume that these sets are separated
 52 (as the semantic information of natural data). The "similarity" between A_f^i and A^i ($1 \leq i \leq k$) reflects
 53 how well the ML model f has learned the semantic information of the training data.

54 **Geometrical information** As for the geometrical information, we are motivated by the methods in
 55 OOD detection [13–16]. In these works, the scores of the OOD samples are lower than in-distribution
 56 samples. In our setting, by the low-dimensional manifold assumption, we know that the off-manifold
 57 data are also outside of the data distribution. Thus, by approximating the shape of the manifold, the
 58 concentration multiplier ϕ should assign lower scores to those off-manifold samples, see Definition
 59 4.2 for a formal definition. In summary, our paper assumes that natural data is drawn from a
 60 low-dimensional manifold and the source and target models capture the semantic and geometrical
 61 information in the way we have discussed above. Our assumption is intuitive, reasonable, and milder
 62 than previous works that theoretically analyze TBAs. Our model fulfills criterion 2.

63 Last but not least, the following remark explains why our paper does not present any experiments.
 64 *Remark A.2* (Experiments are unnecessary for validating our model). As mentioned in Section 2.1,
 65 most of the recent studies on TBAs focus on empirically improving the success rates of TBAs [3, 17].
 66 However, to the best of our knowledge, existing theoretical analyses of TBAs [18–20] are either based
 67 on simple models (e.g., linear classifiers) or strong assumptions (e.g., natural data are drawn from the
 68 spherical Gaussian distribution). The theoretical studies of TBAs are falling behind the engineering
 69 practice, which motivates us to propose an explanatory model that analyzes and explains the existing
 70 empirical results. As is discussed in Appendix A.1, we argue that conducting experiments (on either
 71 real-world or synthetic datasets) is unnecessary for evaluating an explanatory model. Therefore, we
 72 do not include experiments in our paper. ▲

73 A.2 Visualization of the Non-Adversarial Region

74 We provide a visualization of Example 4.10 in Figure A.1.

75 B Complete Proofs

76 **Proposition 4.3** (semantic classifier, binary case). *Given 2λ -separated sets $A_f, B_f \subset \mathcal{M}$. Define:*

$$f_b(\mathbf{x}) = f_b(\mathbf{x}; A_f, B_f) := \frac{d_p(\mathbf{x}, B_f) - d_p(\mathbf{x}, A_f)}{d_p(\mathbf{x}, B_f) + d_p(\mathbf{x}, A_f)}. \quad (\text{B.1})$$

77 *Then, f_b is a semantic classifier. In particular, we can obtain from Equation (B.1) that $f_b(\mathbf{x}) > 0$ if \mathbf{x}*
 78 *is closer (w.r.t. d_p) to A_f than B_f and $f_b(\mathbf{x}) < 0$ otherwise.*

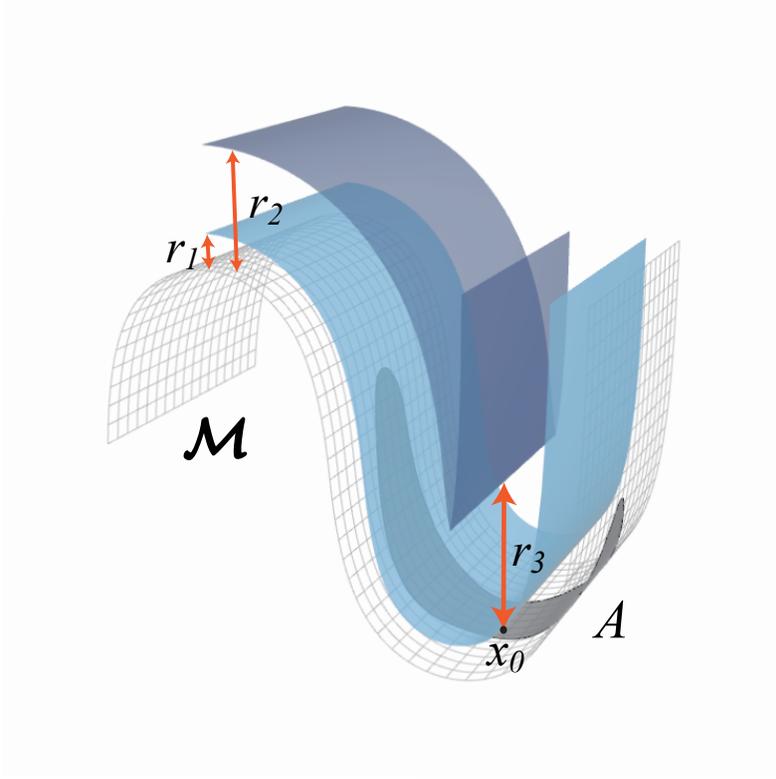


Figure A.1: A visualization of Example 4.10. The data manifold \mathcal{M} is represented by the grid surface. Let the surface in light blue (or dark blue) be the contour surface that $\phi_1 = 0$ (or $\phi_2 = 0$). The distance between \mathbf{x}_0 and the dark blue surface is r_3 , which is greater than δ and r_2 .

79 *Proof of Proposition 4.3.* It is easy to check that $f_b(\mathbf{x}) = 1$ when $\mathbf{x} \in A_f$ and $f_b(\mathbf{x}) = -1$ when $\mathbf{x} \in B_f$.
 80 By definition, we know that f_b is a semantic classifier. \square

81 **Proposition 4.4.** Take $A_f = A$ and $B_f = B$ in Equation (B.1) and denote the corresponding classifier
 82 by f_b^* . Then, for any given $\lambda \geq \delta > 0$, we have $R_{\text{std}}(f_b^*) = R_{\text{adv}}(f_b^*, \delta) = 0$.

83 *Proof of Proposition 4.4.* By Equation (B.1), we have

$$f_b^*(\mathbf{x}) = \frac{d_p(\mathbf{x}, B) - d_p(\mathbf{x}, A)}{d_p(\mathbf{x}, B) + d_p(\mathbf{x}, A)}, \quad \forall \mathbf{x} \in \mathbb{R}^d. \quad (\text{B.2})$$

84 Clearly, we have $f_b^*(\mathbf{x}) = 1$ when $\mathbf{x} \in A$ and $f_b^*(\mathbf{x}) = -1$ when $\mathbf{x} \in B$. Then, the standard risk of f_b^*
 85 w.r.t. $D(\mathbf{x})$ is

$$R_{\text{std}}(f_b^*) = \mathbb{P}_D [f_b^*(\mathbf{x})y < 0 \mid \mathbf{x} \in A] + \mathbb{P}_D [f_b^*(\mathbf{x})y < 0 \mid \mathbf{x} \in B] = 0 \quad (\text{B.3})$$

86 Recall that A and B are 2λ -separated (cf. Definition 3.2). For $\forall \mathbf{x} \in A$ and $\mathbf{x}' \in B(\mathbf{x}, \delta)$, we have
 87 $d_p(\mathbf{x}, B) > \delta$, which implies that $d_p(\mathbf{x}', B) - d_p(\mathbf{x}', A) > 0$, and thus $f_b^*(\mathbf{x}')f_b^*(\mathbf{x}) = f_b^*(\mathbf{x}') > 0$. For
 88 $\forall \mathbf{x} \in B$, a similar deduction shows that $f_b^*(\mathbf{x}')f_b^*(\mathbf{x}) > 0$ holds for $\forall \mathbf{x}' \in B(\mathbf{x}, \delta)$. Together, we have

$$\begin{aligned} R_{\text{adv}}(f_b^*, \delta) := & \mathbb{P} [\exists \mathbf{x}' \in B(\mathbf{x}; \delta) \text{ s.t. } f_b^*(\mathbf{x}')f_b^*(\mathbf{x}) < 0 \mid \mathbf{x} \in A] \\ & + \mathbb{P} [\exists \mathbf{x}' \in B(\mathbf{x}; \delta) \text{ s.t. } f_b^*(\mathbf{x}')f_b^*(\mathbf{x}) < 0 \mid \mathbf{x} \in B] = 0, \end{aligned} \quad (\text{B.4})$$

89 which completes the proof. \square

90 *Remark B.1.* The construction of Equation (B.2) can be found in previous works [6, 21]. In particular,
 91 Li et al. [6] uses the ReLU-approximation of f_b^* to study the robust generalization of deep NNs.

92 **Proposition 4.5** (Concentration multiplier, binary case). For any given $r > 0$ and $G \subset \mathbb{R}^d$, denote

$$\phi(\mathbf{x}) = \phi(\mathbf{x}; r, G) := \frac{r - d_p(\mathbf{x}, G)}{r + d_p(\mathbf{x}, G)}, \quad \forall \mathbf{x} \in \mathbb{R}^d. \quad (\text{B.5})$$

93 Then $\phi(\mathbf{x})$ is a concentration multiplier around G .

94 *Proof of Proposition 4.5.* For $\forall \mathbf{x} \in G$, we have $d_p(\mathbf{x}, G) = 0$. That is, $\phi(\mathbf{x}) = 1$ for $\forall \mathbf{x} \in G$. For
95 $\forall \mathbf{x}_1, \mathbf{x}_2$ s.t. $d_p(\mathbf{x}_1, G) > d_p(\mathbf{x}_2, G)$, it is easy to check that $\phi(\mathbf{x}_1) < \phi(\mathbf{x}_2)$. \square

96 **Proposition 4.6.** Let $f = f_b \cdot \phi$ and A_f, B_f be the semantic information of f_b . We can obtain that

- 97 1. if $R_{\text{adv}}(f; \delta) \neq 0$, then f suffers from off-manifold adversarial examples.
98 2. if $R_{\text{adv}}(f; \delta) \neq 0$ and $d_p(A \cup B, (A_f \cup B_f)^c) > \delta$, then all the adversarial examples of f are
99 off the manifold.¹

100 *Proof of Proposition 4.6.* We first prove the first result. By definition, there are $r > 0$ and $G \subset \mathbb{R}^d$
101 such that

$$f(\mathbf{x}) = \frac{d_p(\mathbf{x}, B_f) - d_p(\mathbf{x}, A_f)}{d_p(\mathbf{x}, B_f) + d_p(\mathbf{x}, A_f)} \cdot \frac{r - d_p(\mathbf{x}, G)}{r + d_p(\mathbf{x}, G)} \quad (\text{B.6})$$

102 Given $R_{\text{adv}}(f; \delta) \neq 0$, then $\exists \mathbf{x} \in A \cup B$ and $\mathbf{x}_0 \in B(\mathbf{x}, \delta)$ such that $f(\mathbf{x})f(\mathbf{x}_0) < 0$. If $x_0 \in \mathcal{M}^c$, there is
103 nothing to prove.

104 Otherwise, we have $\mathbf{x}_0 \in \mathcal{M}$. Without loss of generality (WLOG), we assume that such $\mathbf{x} \in A$
105 and $f(\mathbf{x}_0) < 0$, which implies that either $f_b(\mathbf{x}_0) < 0$ or $\phi(\mathbf{x}_0) < 0$. We first consider the case when
106 $f_b(\mathbf{x}_0) < 0$. Then, we have $d_p(\mathbf{x}, B_f) - d_p(\mathbf{x}, A_f) < 0$ and $r - d_p(\mathbf{x}, S) > 0$. Denote

$$r_0 := \frac{1}{3} \min\{|d_p(\mathbf{x}_0, A_f) - d_p(\mathbf{x}_0, B_f)|, |r - d_p(\mathbf{x}_0, S)|, \delta\}. \quad (\text{B.7})$$

Consider the non-empty set

$$B(\mathbf{x}, \delta) \cap B(\mathbf{x}_0, r_0) \cap \mathcal{M}^c.$$

107 For $\forall \mathbf{x}'_0 \in B(\mathbf{x}_0, r_0)$, there is

$$d_p(\mathbf{x}'_0, A_f) \geq d_p(\mathbf{x}_0, A_f) - d_p(\mathbf{x}_0, \mathbf{x}'_0), \quad (\text{B.8})$$

108 and

$$d_p(\mathbf{x}'_0, B_f) \leq d_p(\mathbf{x}_0, B_f) + d_p(\mathbf{x}_0, \mathbf{x}'_0), \quad (\text{B.9})$$

109 which implies that

$$d_p(\mathbf{x}'_0, A_f) - d_p(\mathbf{x}'_0, B_f) \geq d_p(\mathbf{x}_0, A_f) - d_p(\mathbf{x}_0, B_f) - 2d_p(\mathbf{x}_0, \mathbf{x}'_0) \geq r_0 > 0. \quad (\text{B.10})$$

110 Similarly, we can obtain $r - d_p(\mathbf{x}'_0, S) > 0$. Together, these two inequalities lead us to $f(\mathbf{x}'_0) =$
111 $f(\mathbf{x}'_0)f(\mathbf{x}) < 0$, i.e., \mathbf{x}'_0 is an off-manifold adversarial example of x_0 . Some tedious manipulation
112 yields the same result when $\phi(\mathbf{x}_0) < 0$, which is omitted here.

113 As for the second result, since $R_{\text{adv}}(f; \delta) \neq 0$, we can obtain from the first result that off-manifold
114 adversarial examples exist. It remains to show that f has no on-manifold adversarial examples.

115 Since $d_p(A \cup B, (A_f \cup B_f)^c) > \delta$ and by assumption $A_f \cup B_f \subset G$, we have

$$\frac{r - d_p(\mathbf{x}, G)}{r + d_p(\mathbf{x}, G)} = 1 \quad (\text{B.11})$$

116 and

$$\frac{r - d_p(\mathbf{x}', G)}{r + d_p(\mathbf{x}', G)} = 1 \quad (\text{B.12})$$

117 for $\forall \mathbf{x} \in A \cup B$ and $\mathbf{x}' \in B(\mathbf{x}, \delta) \cap \mathcal{M}$. We can easily obtain that $f_b(\mathbf{x}) = f_b(\mathbf{x}')$, which implies that f
118 has no on-manifold adversarial examples. \square

¹**Notice:** In the main part of the paper, we made a typo in this result. Here, we provide the corrected version. The other results in the main paper are based on the corrected version of this result.

119 **Proposition 4.7.** Consider TBAs with perturbation radius $\delta \in (0, \lambda]$, target model $f_t = f_b^* \cdot \phi_{\text{off}}$ and
 120 source model $f_s = f_b \cdot \phi_{\text{off}}$, $f_b \in \mathcal{F}_b$. Denote the semantic information of f_b by A_f and B_f . Then, for
 121 $\forall \mathbf{x} \in A \cup B$, all adversarial examples (if exist) of f_s at \mathbf{x} are transferable if $\mathbf{x} \in A_f \cup B_f$.

122 *Proof of Proposition 4.7.* Consider $\mathbf{x} \in A_f$ WLOG. By Equation (6), denote

$$f_s(\mathbf{x}) = f_b(\mathbf{x}) \cdot \phi_{\text{off}}(\mathbf{x}) = \frac{d_p(\mathbf{x}, B_f) - d_p(\mathbf{x}, A_f)}{d_p(\mathbf{x}, B_f) + d_p(\mathbf{x}, A_f)} \cdot \frac{\alpha\delta - d_p(\mathbf{x}, \mathcal{M})}{\alpha\delta + d_p(\mathbf{x}, \mathcal{M})}. \quad (\text{B.13})$$

123 If f_s is robust against adversarial examples at $\mathbf{x} \in B(\mathbf{x}, \delta) \cap \mathcal{M} \subset A_f$, then there is nothing to prove. If
 124 not, denote the adversarial example of f_s at \mathbf{x} by \mathbf{x}_a , and we have $f_s(\mathbf{x}_a) < 0$. It is not hard to verify
 125 that $d_p(\mathbf{x}_a, B_f) - d_p(\mathbf{x}_a, A_f) > 0$ since $\delta < \lambda$, which implies that $f_b(\mathbf{x}_a) > 0$. To obtain $f_s(\mathbf{x}_a) < 0$,
 126 there must be $\phi_{\text{off}}(\mathbf{x}_a) < 0$. We thus have $\alpha\delta < d_p(\mathbf{x}_a, \mathcal{M})$, which implies that \mathbf{x}_a is off the manifold
 127 and the distance between \mathbf{x}_a and \mathcal{M} is greater than $\alpha\delta$. In particular, we have

$$\phi_{\text{off}}(\mathbf{x})\phi_{\text{off}}(\mathbf{x}_a) < 0, \quad (\text{B.14})$$

128 which is independent of the choice of f_b . Now consider $f_t(\mathbf{x})$ and $f_t(\mathbf{x}_a)$, where

$$f_t(\mathbf{x}) = f_b^*(\mathbf{x}) \cdot \phi_{\text{off}}(\mathbf{x}) = \frac{d_p(\mathbf{x}, B) - d_p(\mathbf{x}, A)}{d_p(\mathbf{x}, B) + d_p(\mathbf{x}, A)} \cdot \frac{\alpha\delta - d_p(\mathbf{x}, \mathcal{M})}{\alpha\delta + d_p(\mathbf{x}, \mathcal{M})}. \quad (\text{B.15})$$

129 No matter $\mathbf{x} \in A$ or $\mathbf{x} \in B$, we have for $\forall \mathbf{x}' \in B(\mathbf{x}, \delta)$, there is $f_b^*(\mathbf{x}) = f_b^*(\mathbf{x}')$ (by the 2λ -separated
 130 property of A and B). By Equation (B.14), we have

$$f_t(\mathbf{x})f_t(\mathbf{x}_a) = f_b^*(\mathbf{x})f_b^*(\mathbf{x}_a) \cdot \phi_{\text{off}}(\mathbf{x})\phi_{\text{off}}(\mathbf{x}_a) < 0, \quad (\text{B.16})$$

131 i.e., \mathbf{x}_a transfers to f_t , which completes the proof. \square

132 **Proposition 4.8.** Consider TBA with perturbation radius $\delta \in (0, \lambda]$, target model $f_t = f_b^* \cdot \phi_{\text{on}}(\cdot, f_b^*)$
 133 and source model $f_s = f_b \cdot \phi_{\text{on}}(\cdot, f_b)$, $f_b \in \mathcal{F}_b$. Denote

$$S_{\text{crit}} := (A \cap A_f) \cup (B \cap B_f), \quad S_{\text{wrg}} := (A \cap B_f) \cup (B \cap A_f). \quad (\text{B.17})$$

134 Then, for $\forall \mathbf{x} \in A \cup B$, we have

- 135 1. if $B(\mathbf{x}, \delta) \cap \mathcal{M} \subset S_{\text{crit}} \cup S_{\text{wrg}}$, then f_t and f_s are both robust against adversarial examples;
- 136 2. if $B(\mathbf{x}, \delta) \cap \mathcal{M} \subset A \cup B$ and $\mathbf{x} \in A_f \cup B_f$, then the adversarial examples of f_s at \mathbf{x} (if exists)
 137 cannot transfer to f_t .

138 *Proof of Proposition 4.8.* The proof of the first result is also straightforward, which is omitted here.
 139 For $\forall \mathbf{x} \in A \cup B$ such that $(B(\mathbf{x}, \delta) \cap \mathcal{M}) \subset S_{\text{crit}}$, it is easy to check that $\phi_{\text{on}}(\mathbf{x}) = 1$ and $\phi_{\text{on}}(\mathbf{x}_a) = 1$ for
 140 $\forall \mathbf{x}_a \in B(\mathbf{x}, \delta)$, which implies that f is robust against adversarial examples. It remains to prove the
 141 second result. By Equation (7), denote

$$f_s(\mathbf{x}) = f_b(\mathbf{x}) \cdot \phi_{\text{on}}(\mathbf{x}; f_b) = \frac{d_2(\mathbf{x}, B_f) - d_2(\mathbf{x}, A_f)}{d_2(\mathbf{x}, B_f) + d_2(\mathbf{x}, A_f)} \cdot \frac{\alpha\delta - d_2(\mathbf{x}, \mathcal{N}_\delta(A_f \cup B_f))}{\alpha\delta + d_2(\mathbf{x}, \mathcal{N}_\delta(A_f \cup B_f))}. \quad (\text{B.18})$$

142 For $\forall \mathbf{x} \in A \cup B$ such that $B(\mathbf{x}, \delta) \cap \mathcal{M} \subset A \cup B$ and $\mathbf{x} \in A_f \cup B_f$, denote the unspecific adversarial
 143 example (if exist) of f_s at \mathbf{x} by \mathbf{x}_a . Assume that $\mathbf{x} \in A \cap B_f$ WLOG. By definition, we have

$$f_b(\mathbf{x}) < 0. \quad (\text{B.19})$$

144 By the 2λ -separated assumption of A_f and B_f , we have and

$$f_b(\mathbf{x}_a) < 0. \quad (\text{B.20})$$

145 From $f_s(\mathbf{x})f_s(\mathbf{x}_a) < 0$, we can obtain that $\phi_{\text{on}}(\mathbf{x}; f_b)\phi_{\text{on}}(\mathbf{x}_a; f_b) < 0$. Since $\mathbf{x} \in B_f$, we have

$$d_2(\mathbf{x}, \mathcal{N}_\delta(A_f \cup B_f)) = 0, \quad (\text{B.21})$$

146 i.e., $\phi_{\text{on}}(\mathbf{x}; f_b) = 1$. Combine this with $\phi_{\text{on}}(\mathbf{x}; f_b)\phi_{\text{on}}(\mathbf{x}_a; f_b) < 0$, we have $\phi_{\text{on}}(\mathbf{x}_a; f_b) < 0$, i.e.,

$$\alpha\delta < d_2(\mathbf{x}_a, \mathcal{N}_\delta(A_f \cup B_f)) < d_2(\mathbf{x}_a, \mathbf{x}) \leq \delta. \quad (\text{B.22})$$

147 Since $f_t(\mathbf{x}) = 1$ and \mathbf{x}_a is unspecific, it remains to show that $f_t(\mathbf{x}_a) > 0$. By Equation (7), denote

$$f_t(\mathbf{x}) = f_b^*(\mathbf{x}) \cdot \phi_{\text{on}}(\mathbf{x}; f_b^*) = \frac{d_2(\mathbf{x}, B) - d_2(\mathbf{x}, A)}{d_2(\mathbf{x}, B) + d_2(\mathbf{x}, A)} \cdot \frac{\alpha\delta - d_2(\mathbf{x}, \mathcal{N}_\delta(A \cup B))}{\alpha\delta + d_2(\mathbf{x}, \mathcal{N}_\delta(A \cup B))}. \quad (\text{B.23})$$

148 By $\mathbf{x} \in A$ and the 2λ -separated assumption of A and B , we have $f_b^*(\mathbf{x}_a) > 0$. By $B(\mathbf{x}, \delta) \cap \mathcal{M} \subset A$,
 149 we have $\mathbf{x}_a \in \mathcal{N}_\delta(A \cup B)$, i.e., $\phi_{\text{on}}(\mathbf{x}_a; f_b^*) > 0$. Together, we have $f_t(\mathbf{x}_a) = f_b^*(\mathbf{x}_a) \cdot \phi_{\text{on}}(\mathbf{x}_a; f_b^*) > 0$,
 150 which completes the proof. \square

151 **Proposition 4.12.** Given perturbation radius $\delta \in (0, \lambda]$ and target model $f_t = f_b^* \cdot \phi_{\text{off}}(\cdot; r, \mathcal{M})$. Let Δ
 152 be the constant specified in Lemma 4.11. Then, for $\forall \mathbf{x} \in A \cup B$, the off-manifold adversarial example
 153 of f_t at \mathbf{x} exists if $r < \Delta$.

154 *Proof of Proposition 4.12.* For $\forall \mathbf{x} \in A \cup B$, let $\mathbf{u} \in N_{\mathbf{x}}(\mathcal{M})$ be the normal direction at \mathbf{x} with $\|\mathbf{u}\|_2 = 1$.
 155 Since $r < \Delta$, we can find $r_0 > r$ such that $r_0 < \Delta$ and $r_0 < \delta$. Denote

$$\mathbf{x}_a := \mathbf{x} + r_0 \mathbf{u}. \quad (\text{B.24})$$

156 Clearly, we have $\mathbf{x}_a \in B(\mathbf{x}, \delta)$. Since $\mathcal{N}_\Delta(\mathcal{M})$ is a tubular neighborhood of \mathcal{M} , we have

$$d_2(\mathbf{x}_a, \mathcal{M}) = r_0 > r, \quad (\text{B.25})$$

157 which implies that \mathbf{x}_a is an off-manifold adversarial example of f_t at \mathbf{x} . \square

158 **Corollary 5.3.** Let f_b be a semantic classifier with semantic information A_f and B_f that satisfy a
 159 2λ -separated property. Given $\epsilon > 0$, there is a ReLU network \tilde{f} with $O((1/\lambda\epsilon)^d) \cdot O(d^2 + d \log(1/\epsilon))$
 160 parameters such that $\|f - \tilde{f}\|_\infty \leq \epsilon$.

161 *Proof of Corollary 5.3.* According to Lemma 5.2, our goal is to upper bound the Lipschitz constant l
 162 of f_b . By definition, it suffices to upper bound the supremum of

$$\begin{aligned} s &:= \frac{|f_b(\mathbf{x}_1; A_f, B_f) - f_b(\mathbf{x}_2; A_f, B_f)|}{d_p(\mathbf{x}_1, \mathbf{x}_2)} \\ &= \frac{1}{d_p(\mathbf{x}_1, \mathbf{x}_2)} \cdot \left| \frac{d_p(\mathbf{x}_1, A_f)}{d_p(\mathbf{x}_1, A_f) + d_p(\mathbf{x}_1, B_f)} - \frac{d_p(\mathbf{x}_2, A_f)}{d_p(\mathbf{x}_2, A_f) + d_p(\mathbf{x}_2, B_f)} \right|. \end{aligned} \quad (\text{B.26})$$

163 We only need to consider three cases:

- 164 1. both of $\mathbf{x}_1, \mathbf{x}_2 \in A_f \cup B_f$, or
- 165 2. both of $\mathbf{x}_1, \mathbf{x}_2 \in (A_f \cup B_f)^c$, and
- 166 3. either \mathbf{x}_1 or \mathbf{x}_2 is in $A_f \cup B_f$.

167 When $\mathbf{x}_1, \mathbf{x}_2 \in A_f \cup B_f$, a trivial verification shows that $s \leq \frac{1}{\lambda}$. We now turn to the second case. By
 168 symmetry, let

$$\frac{d_p(\mathbf{x}_1, A_f)}{d_p(\mathbf{x}_1, A_f) + d_p(\mathbf{x}_1, B_f)} - \frac{d_p(\mathbf{x}_2, A_f)}{d_p(\mathbf{x}_2, A_f) + d_p(\mathbf{x}_2, B_f)} > 0. \quad (\text{B.27})$$

169 By simplifying Equation (B.26), we can obtain that

$$\begin{aligned} &\frac{|f_b(\mathbf{x}_1; A_f, B_f) - f_b(\mathbf{x}_2; A_f, B_f)|}{d_p(\mathbf{x}_1, \mathbf{x}_2)} \\ &= \frac{1}{d_p(\mathbf{x}_1, \mathbf{x}_2)} \cdot \left(\frac{d_p(\mathbf{x}_1, A_f)}{d_p(\mathbf{x}_1, A_f) + d_p(\mathbf{x}_1, B_f)} - \frac{d_p(\mathbf{x}_2, A_f)}{d_p(\mathbf{x}_2, A_f) + d_p(\mathbf{x}_2, B_f)} \right) \\ &\leq \frac{1}{2\lambda} \cdot \left(\frac{d_p(\mathbf{x}_1, A_f) - d_p(\mathbf{x}_2, A_f)}{d_p(\mathbf{x}_1, \mathbf{x}_2)} \cdot \frac{d_p(\mathbf{x}_2, B_f)}{d_p(\mathbf{x}_2, A_f) + d_p(\mathbf{x}_2, B_f)} \right. \\ &\quad \left. + \frac{d_p(\mathbf{x}_1, B_f) - d_p(\mathbf{x}_2, B_f)}{d_p(\mathbf{x}_1, \mathbf{x}_2)} \cdot \frac{d_p(\mathbf{x}_2, A_f)}{d_p(\mathbf{x}_2, A_f) + d_p(\mathbf{x}_2, B_f)} \right) \\ &\leq \frac{1}{2\lambda} \cdot (1 \cdot 1 + 1 \cdot 1) = \frac{1}{\lambda}, \end{aligned} \quad (\text{B.28})$$

170 which implies that $s \leq \frac{1}{\lambda}$ in this case. Finally, we consider the third case. We assume WLOG that
 171 $\mathbf{x}_1 \in A_f$ and $\mathbf{x}_2 \in (A_f \cup B_f)^c$. Substitute into Equation (B.26), we have

$$\begin{aligned} s &= \frac{1}{d_p(\mathbf{x}_1, \mathbf{x}_2)} \cdot \left| \frac{d_p(\mathbf{x}_1, A_f)}{d_p(\mathbf{x}_1, A_f) + d_p(\mathbf{x}_1, B_f)} - \frac{d_p(\mathbf{x}_2, A_f)}{d_p(\mathbf{x}_2, A_f) + d_p(\mathbf{x}_2, B_f)} \right| \\ &= \frac{1}{d_p(\mathbf{x}_1, \mathbf{x}_2)} \cdot \frac{d_p(\mathbf{x}_2, A_f)}{d_p(\mathbf{x}_2, A_f) + d_p(\mathbf{x}_2, B_f)} \leq \frac{1}{2\lambda} \end{aligned} \quad (\text{B.29})$$

172 To sum up above, we have $\sup_{\mathbf{x}_1 \neq \mathbf{x}_2} s = \frac{1}{\lambda}$, which implies that f_b is $\frac{1}{\lambda}$ -Lipschitz continuous, as is
 173 required. \square

174 **Corollary 5.4.** Given $\epsilon > 0$, $r > 0$ and $S \subset [0, 1]^d$, there is a ReLU network $\tilde{\phi}$ with $O((1/r\epsilon)^d) \cdot$
 175 $O(d^2 + d \log(1/\epsilon))$ parameters that can approximate $\phi(\cdot; r, S)$ to precision ϵ .

176 *Proof of Corollary 5.4.* We prove this corollary in a similar manner as Corollary 5.3, i.e., we upper
 177 bound the supremum of

$$\begin{aligned} s &:= \frac{|\phi(\mathbf{x}_1; r, S) - \phi(\mathbf{x}_2; r, S)|}{d_p(\mathbf{x}_1, \mathbf{x}_2)} = \frac{1}{d_p(\mathbf{x}_1, \mathbf{x}_2)} \cdot \left| \frac{d_p(\mathbf{x}_1, S)}{r + d_p(\mathbf{x}_1, S)} - \frac{d_p(\mathbf{x}_2, S)}{r + d_p(\mathbf{x}_2, S)} \right| \\ &= \frac{r}{d_p(\mathbf{x}_1, \mathbf{x}_2)} \cdot \left| \frac{1}{r + d_p(\mathbf{x}_1, S)} - \frac{1}{r + d_p(\mathbf{x}_2, S)} \right|. \end{aligned} \quad (\text{B.30})$$

178 We also consider three cases in this proof:

- 179 1. both of $\mathbf{x}_1, \mathbf{x}_2 \in S$, or
- 180 2. both of $\mathbf{x}_1, \mathbf{x}_2 \in S^c$, and
- 181 3. either \mathbf{x}_1 or \mathbf{x}_2 is in S .

182 In case 1, we see at once that $s = 0$. When both of $\mathbf{x}_1, \mathbf{x}_2 \in S^c$, we assume WLOG that $d_p(\mathbf{x}_1, S) >$
 183 $d_p(\mathbf{x}_2, S)$. By Equation (B.30), we have

$$s = \frac{d_p(\mathbf{x}_1, S) - d_p(\mathbf{x}_2, S)}{d_p(\mathbf{x}_1, \mathbf{x}_2)} \cdot \frac{r}{(r + d_p(\mathbf{x}_1, S))(r + d_p(\mathbf{x}_2, S))} \leq \frac{1}{r}. \quad (\text{B.31})$$

184 Analysis similar to Equation (B.31) shows that

$$s = \frac{r}{d_p(\mathbf{x}_1, \mathbf{x}_2)} \cdot \left(\frac{1}{r} - \frac{1}{r + d_p(\mathbf{x}_2, S)} \right) \leq \frac{1}{r}. \quad (\text{B.32})$$

185 To sum up above, we have $\sup_{\mathbf{x}_1 \neq \mathbf{x}_2} s = \frac{1}{\lambda}$, which implies that ϕ is $\frac{1}{r}$ -Lipschitz continuous, as is
 186 required. \square

187 **Proposition 5.6.** Given $\epsilon, \lambda, \delta, r > 0$, for any $f \in \mathcal{F}_{\mathcal{M}}$, there is a ReLU network \tilde{f} with

$$O(\max\{\frac{1}{\lambda\epsilon}, \frac{2}{r\epsilon}\}^d) \cdot O(d^2 + d \log(\frac{1}{\epsilon})) + O(\log^2(\frac{1}{\epsilon})) \quad (\text{B.33})$$

188 parameters that satisfies $\|f - \tilde{f}\|_{\infty} \leq \epsilon$.

189 *Proof of Proposition 5.6.* This proposition can be derived directly from Lemma 5.5, Corollary 5.3,
 190 and Corollary 5.4. \square

191 **Theorem 5.7.** Consider TBAs with perturbation radius $\delta \in (0, \lambda/2]$, target model $f_t = f_b^* \cdot \phi_{\text{off}}$ and
 192 source model $f_s = f_b \cdot \phi_{\text{off}}$, $f_b \in \mathcal{F}_b$. Denote the semantic information of f_b by A_f and B_f . Given
 193 $\epsilon \leq 0.1$, let \tilde{f}_t and \tilde{f}_s be ReLU networks that satisfy

$$\|\tilde{f}_t - f_t\|_{\infty} \leq \epsilon, \|\tilde{f}_s - f_s\|_{\infty} \leq \epsilon \quad (\text{B.34})$$

194 Then, for $\forall \mathbf{x} \in (A \cup B) \cap (A_f \cup B_f)$, the adversarial examples \mathbf{x}_a (if exist) of \tilde{f}_s satisfies

$$\tilde{f}_t(\mathbf{x}) \cdot \tilde{f}_t(\mathbf{x}_a) \leq 2\epsilon(1 + \epsilon)^2 + 2\epsilon^2. \quad (\text{B.35})$$

195 *Proof of Theorem 5.7.* Consider $\mathbf{x} \in A_f$ WLOG. If f_s is robust against adversarial examples at
 196 $\mathbf{x} \in B(\mathbf{x}, \delta) \cap \mathcal{M} \subset A_f$, then there is nothing to prove. If not, denote the adversarial example of f_s at \mathbf{x}
 197 by \mathbf{x}_a . Since $\mathbf{x} \in A_f \subset \mathcal{M}$, there is

$$\tilde{f}_s(\mathbf{x}) \geq \tilde{f}_b(\mathbf{x}) \cdot \tilde{\phi}_{\text{off}}(\mathbf{x}) - \epsilon \geq (1 - \epsilon)^2 - \epsilon > 0 \quad (\text{B.36})$$

198 and we thus have $\tilde{f}_s(\mathbf{x}_a) < 0$. By $\mathbf{x}_a \in B(\mathbf{x}; \delta)$ and the assumption $\delta < \frac{\lambda}{2}$, we have

$$\tilde{f}_b(\mathbf{x}_a) = 1 - \frac{2d_p(\mathbf{x}_a, A_f)}{d_p(\mathbf{x}_a, B_f) + d_p(\mathbf{x}_a, A_f)} \geq 1 - \frac{2\delta}{2\lambda} \geq \frac{1}{2}. \quad (\text{B.37})$$

199 To obtain $\tilde{\chi}(\tilde{f}_b, \tilde{\phi}_{\text{off}})(\mathbf{x}_a) < 0$, there must be $\tilde{f}_b(\mathbf{x}_a) \cdot \tilde{\phi}_{\text{off}}(\mathbf{x}_a) < \epsilon$, which implies that

$$\tilde{\phi}_{\text{off}}(\mathbf{x}_a) < 2\epsilon. \quad (\text{B.38})$$

200 Now consider $\tilde{f}_i(\mathbf{x})$ and $\tilde{f}_i(\mathbf{x}_a)$. By definition, we have $\tilde{f}_b^*(\mathbf{x}) \in [1 - \epsilon, 1 + \epsilon]$, $\tilde{\phi}_{\text{off}}(\mathbf{x}) \in [1 - \epsilon, 1 + \epsilon]$,
201 and thus

$$\tilde{f}_i(\mathbf{x}) = \tilde{\chi}(\tilde{f}_b^*, \tilde{\phi}_{\text{off}})(\mathbf{x}) \leq (1 + \epsilon)^2 + \epsilon. \quad (\text{B.39})$$

202 Similar to Equation (B.37), there is

$$\tilde{f}_b^*(\mathbf{x}_a) = 1 - \frac{2d_p(\mathbf{x}_a, A)}{d_p(\mathbf{x}_a, B) + d_p(\mathbf{x}_a, A)} \leq 1 \quad (\text{B.40})$$

203 Combining Equations (B.38) to (B.40) together, we have

$$\tilde{f}_i(\mathbf{x}) \cdot \tilde{f}_i(\mathbf{x}_a) \leq 2\epsilon(1 + \epsilon)^2 + 2\epsilon^2. \quad (\text{B.41})$$

204 as is required. \square

205 **Theorem 5.8.** Consider TBAs with perturbation radius $\delta \in (0, \lambda/2]$, target model $f_t = f_b^* \cdot \phi_{\text{off}}$ and
206 source model $f_s = f_b \cdot \phi_{\text{off}}$, $f_b \in \mathcal{F}_b$. Denote the semantic information of f_b by A_f and B_f . Given
207 $\epsilon \leq 0.1$, let \tilde{f}_t and \tilde{f}_s be ReLU networks that satisfy Equation (13). Then, for $\forall \mathbf{x} \in A \cup B$, we have

- 208 1. if $B(\mathbf{x}, \delta) \cap \mathcal{M} \subset S_{\text{crit}} \cup S_{\text{wrg}}$, then \tilde{f}_t and \tilde{f}_s are both robust against adversarial examples;
- 209 2. if $B(\mathbf{x}, \delta) \cap \mathcal{M} \subset A \cup B$ and $\mathbf{x} \in A_f \cup B_f$, then the adversarial examples of \tilde{f}_s at \mathbf{x} (if exists)
210 cannot transfer to \tilde{f}_t .

211 *Proof of Theorem 5.8.* The proof of the first result is also straightforward, which is omitted here. It
212 remains to prove the second result. By Equation (7), denote

$$f_s(\mathbf{x}) = f_b(\mathbf{x}) \cdot \phi_{\text{on}}(\mathbf{x}; f_b) = \frac{d_2(\mathbf{x}, B_f) - d_2(\mathbf{x}, A_f)}{d_2(\mathbf{x}, B_f) + d_2(\mathbf{x}, A_f)} \cdot \frac{\alpha\delta - d_2(\mathbf{x}, \mathcal{N}_\delta(A_f \cup B_f))}{\alpha\delta + d_2(\mathbf{x}, \mathcal{N}_\delta(A_f \cup B_f))}. \quad (\text{B.42})$$

213 For $\forall \mathbf{x} \in A \cup B$ such that $B(\mathbf{x}, \delta) \cap \mathcal{M} \subset A \cup B$ and $\mathbf{x} \in A_f \cup B_f$, denote the unspecific adversarial
214 example (if exist) of f_s at \mathbf{x} by \mathbf{x}_a . Assume that $\mathbf{x} \in A \cap A_f$ WLOG. By definition, we have

$$\tilde{f}_b(\mathbf{x}) \in [1 - \epsilon, 1 + \epsilon]. \quad (\text{B.43})$$

215 By the 2λ -separated assumption of A_f and B_f , we have and

$$\tilde{f}_b(\mathbf{x}_a) \in [1 - \epsilon, 1 + \epsilon]. \quad (\text{B.44})$$

216 Since $\mathbf{x} \in A_f \cup B_f$, we have $\mathbf{x} \in \mathcal{N}_\delta(A_f \cup B_f)$ and

$$\tilde{\phi}_{\text{on}}(\mathbf{x}; f_b) \in [1 - \epsilon, 1 + \epsilon], \quad (\text{B.45})$$

217 which implies that

$$\tilde{f}_s(\mathbf{x}) = \tilde{\chi}(\tilde{f}_b, \tilde{\phi}_{\text{on}}(\cdot; f_b))(\mathbf{x}) \geq (1 - \epsilon)^2 - \epsilon > 0. \quad (\text{B.46})$$

218 From $\tilde{f}_s(\mathbf{x})\tilde{f}_s(\mathbf{x}_a) < 0$, we can obtain that $\tilde{f}_s(\mathbf{x}_a) < 0$, which implies that

$$\tilde{\phi}_{\text{on}}(\mathbf{x}_a; f_b) \cdot \tilde{f}_b(\mathbf{x}) < \epsilon, \quad (\text{B.47})$$

219 which implies that

$$\tilde{\phi}_{\text{on}}(\mathbf{x}_a; f_b) < \frac{\epsilon}{1 - \epsilon} < 2\epsilon. \quad (\text{B.48})$$

220 By definition, we have

$$\tilde{f}_i(\mathbf{x}) = \tilde{\chi}(\tilde{f}_b^*, \tilde{\phi}_{\text{on}}(\cdot; f_b^*))(\mathbf{x}) \geq (1 - \epsilon)^2 - \epsilon > 0. \quad (\text{B.49})$$

221 By $\mathbf{x} \in A$ and the 2λ -separated assumption of A and B , we have

$$\tilde{f}_b^*(\mathbf{x}_a) = 1 - \frac{2d_p(\mathbf{x}_a, A)}{d_p(\mathbf{x}_a, B) + d_p(\mathbf{x}_a, A)} \geq 1 - \frac{2\delta}{2\lambda} \geq \frac{1}{2}. \quad (\text{B.50})$$

222 By $B(\mathbf{x}, \delta) \cap \mathcal{M} \subset A$, we have $\mathbf{x}_a \in \mathcal{N}_\delta(A \cup B)$, i.e.,

$$\phi_{\text{on}}(\mathbf{x}_a; f_b^*) \in [1 - \epsilon, 1 + \epsilon] \quad (\text{B.51})$$

223 Together, we have

$$\tilde{f}_i(\mathbf{x}_a) = \tilde{\chi}(\tilde{f}_b^*, \tilde{\phi}_{\text{on}}(\cdot; f_b^*))(\mathbf{x}_a) \geq \frac{1 - \epsilon}{2} > 0, \quad (\text{B.52})$$

224 which completes the proof. \square

225 **Proposition 5.9.** For any classifier f^* with $R_{\text{std}}(f^*) = 0$ and perturbation radius $\delta \in (r_\delta(f^*), \lambda)$, there
 226 is $f \in \mathcal{F}_M$ such that

- 227 1. $R_{\text{std}}(f) = R_{\text{std}}(f^*)$, and
 228 2. for $\forall \mathbf{x} \in A \cup B$, if \mathbf{x}_a is an adversarial example of f^* at \mathbf{x} , then exists $\mathbf{x}'_a \in B(\mathbf{x}_a, r_\delta(f^*)/4)$
 229 such that \mathbf{x}'_a is an adversarial example of f .

230 *Proof of Proposition 5.9.* Define the following set

$$S_a := \{\mathbf{x} \in [0, 1]^d \cap (A \cup B)^c : \exists \mathbf{x}' \in A \cup B \text{ s.t. } \mathbf{x}' \in B(\mathbf{x}; \delta), f^*(\mathbf{x})f^*(\mathbf{x}') < 0\}, \quad (\text{B.53})$$

231 and let

$$G = \left(\bigcup_{\mathbf{x} \in S_a} B(\mathbf{x}, r_\delta(f^*)/2) \right)^c. \quad (\text{B.54})$$

232 By definition, $A \cup B \in G$. Consider

$$f(\mathbf{x}) = f_b^*(\mathbf{x}) \cdot \phi(\mathbf{x}; r_\delta(f^*)/4, G). \quad (\text{B.55})$$

233 Since $A \cup B \in G$, we have $R_{\text{std}}(f) = R_{\text{std}}(f_b^*) = 0 = R_{\text{std}}(f^*)$. For $\forall \mathbf{x} \in A \cup B$, if \mathbf{x}_a is an adversarial
 234 example of f^* at \mathbf{x} , then

$$\frac{r_\delta(f^*)}{4} \leq d_p(\mathbf{x}_a, G) \leq \frac{r_\delta(f^*)}{2}, \quad (\text{B.56})$$

235 which implies that exists $\mathbf{x}'_a \in B(\mathbf{x}_a, r_\delta(f^*)/4)$ such that \mathbf{x}'_a is an adversarial example of f . \square

236 C Analyses in Multi-Class Classification Problems

237 This section some of the results in Sections 4 and 5 to k -class classification problems. We first extend
 238 Propositions 4.3 and 4.5 to multi-class classification.

239 **Proposition C.1** (Semantic classifier, multi-class case). Given 2λ -separated sets $A_f^1, A_f^2, \dots, A_f^k \subset \mathcal{M}$.

240 Consider $f_b(\mathbf{x}) = (f_b^{(1)}(\mathbf{x}), f_b^{(2)}(\mathbf{x}), \dots, f_b^{(k)}(\mathbf{x}))^T$ and define:

$$f_b^{(i)}(\mathbf{x}) := \frac{\left(\sum_{j \neq i} d_p(\mathbf{x}, A_f^j) \right) - d_p(\mathbf{x}, A_f^i)}{\left(\sum_{j \neq i} d_p(\mathbf{x}, A_f^j) \right) + d_p(\mathbf{x}, A_f^i)} \quad (\text{C.57})$$

241 for $\forall 1 \leq i \leq k$. Then, f_b is a semantic classifier.

242 *Proof of Proposition C.1.* By Equation (C.57), we have

$$f_b^{(i)}(\mathbf{x}) = \frac{\left(\sum_{j=1}^k d_p(\mathbf{x}, A_f^j) \right) - d_p(\mathbf{x}, A_f^i)}{\sum_{j=1}^k d_p(\mathbf{x}, A_f^j)} \quad (\text{C.58})$$

243 for $\forall 1 \leq i \leq k$. Then, there is

$$y(f_b, \mathbf{x}) = \arg \max_{1 \leq i \leq k} f_b^{(i)}(\mathbf{x}) = \arg \max_{1 \leq i \leq k} \left(-d_p(\mathbf{x}, A_f^i) \right) = \arg \min_{1 \leq i \leq k} d_p(\mathbf{x}, A_f^i). \quad (\text{C.59})$$

244 Given that $A_f^1, A_f^2, \dots, A_f^k$ are 2λ -separated, we have

$$0 = d_p(\mathbf{x}, A_f^i) < d_p(\mathbf{x}, A_f^j) \quad (\text{C.60})$$

245 for $\forall j \neq i$ if $\mathbf{x} \in A_f^i$, which completes the proof. \square

246 Next, we specify a family of concentration multipliers for multi-class TBAs.

247 **Proposition C.2** (Concentration multiplier, multi-class case). For any given $r > 0$ and $G \subset \mathbb{R}^d$,
 248 denote

$$\phi(\mathbf{x}) = \phi(\mathbf{x}; r, G) := \frac{r - d_p(\mathbf{x}, G)}{r + d_p(\mathbf{x}, G)}, \quad \forall \mathbf{x} \in \mathbb{R}^d. \quad (\text{C.61})$$

249 Then $\phi(\mathbf{x})$ is a concentration multiplier around G .

250 Note that Equation (C.61) is identical to Equation (5). The proof of Proposition C.2 is therefore
 251 omitted. The following proposition extends Proposition 4.4 to the multi-class case.

252 **Proposition C.3.** Take $A_f^i = A^i$ in Equation (C.57) for $\forall 1 \leq i \leq k$. Denote the corresponding
 253 classifier by f_b^* . Then, for any given $\lambda \geq \delta > 0$, we have $R_{\text{std}}(f_b^*) = R_{\text{adv}}(f_b^*, \delta) = 0$.

254 *Proof of Proposition C.3.* By Equation (C.57), we have

$$f_b^*(\mathbf{x}) = \frac{\left(\sum_{j \neq i} d_p(\mathbf{x}, A^j)\right) - d_p(\mathbf{x}, A^i)}{\left(\sum_{j \neq i} d_p(\mathbf{x}, A^j)\right) + d_p(\mathbf{x}, A^i)} \quad (\text{C.62})$$

255 Apparently, we have $y(f_b^*, \mathbf{x}) = i$ when $\mathbf{x} \in A^i$ for $\forall 1 \leq i \leq k$. The standard risk of f_b^* w.r.t. D is

$$R_{\text{std}}(f_b^*) = \sum_{i=1}^k \left(\mathbb{P}_D \left[y(f_b^*, \mathbf{x}) \neq y(\mathbf{x}) \mid \mathbf{x} \in A^i \right] \right) = 0. \quad (\text{C.63})$$

256 For $\forall i \neq j$, recall that A and B are 2λ -separated (cf. Definition 3.2). For $\forall \mathbf{x} \in A^i$ and $\mathbf{x}' \in B(\mathbf{x}, \delta)$, we
 257 have $d_p(\mathbf{x}', A^j) > \delta$, which implies that $d_p(\mathbf{x}', A_j) > d_p(\mathbf{x}', A^i)$ and

$$\left(\sum_{l=1}^k d_p(\mathbf{x}', A^l) \right) - d_p(\mathbf{x}', A^i) > \left(\sum_{l=1}^k d_p(\mathbf{x}', A^l) \right) - d_p(\mathbf{x}', A^j). \quad (\text{C.64})$$

258 Since j is arbitrarily chosen, and according to Equation (C.58), we have

$$f_b^{(i)}(\mathbf{x}') > f_b^{(j)}(\mathbf{x}') \quad (\text{C.65})$$

259 holds for $\forall j \neq i$, i.e., $y(f_b^*, \mathbf{x}) = y(f_b^*, \mathbf{x}')$ for $\forall \mathbf{x}' \in B(\mathbf{x}, \delta)$. Then, the adversarial risk of f_b^* is

$$R_{\text{adv}}(f_b^*, \delta) = \sum_{i=1}^k \left(\mathbb{P}_D \left[\exists \mathbf{x}_a \in B(\mathbf{x}; \delta) \text{ s.t. } y(f_b^*, \mathbf{x}) \neq y(f_b^*, \mathbf{x}_a) \mid \mathbf{x} \in A^i \right] \right) = 0, \quad (\text{C.66})$$

260 which completes the proof. \square

261 Next, we go straight for the two explanatory results. We first note that the non-existence of off-
 262 manifold adversarial examples is due to the ‘‘sharp curvature’’ of the data manifold. The analyses in
 263 Example 4.10 are regardless of whether the task is binary or multi-class. Here, we extend Proposition
 264 4.12 to multi-class cases. Consider TBAs with perturbation radius $\delta \in (0, \lambda]$ For any unspecified
 265 $\alpha \in (0, 1)$, let

$$\phi_{\text{off}}(\mathbf{x}) := \phi(\mathbf{x}; \alpha\delta, \mathcal{M}) = \frac{\alpha\delta - d_p(\mathbf{x}, \mathcal{M})}{\alpha\delta + d_p(\mathbf{x}, \mathcal{M})}. \quad (\text{C.67})$$

266 Recall that Proposition 4.12 is restricted to $p = 2$. The following proposition provides a sufficient
 267 condition for the existence of off-manifold adversarial examples in multi-class classification tasks.

268 **Proposition C.4.** Given perturbation radius $\delta \in (0, \lambda]$ and target model $f_t = f_b^* \cdot \phi_{\text{off}}$. Let Δ be the
 269 constant specified in Lemma 4.11. Then, for $\forall \mathbf{x} \in \cup_{i=1}^k A^i$, the off-manifold adversarial example of f_t
 270 at \mathbf{x} exists if $\alpha\delta < \Delta$.

271 *Proof of Proposition C.4.* For $\forall \mathbf{x} \in \cup_{i=1}^k A^i$, let $\mathbf{u} \in N_{\mathbf{x}}(\mathcal{M})$ be the normal direction at \mathbf{x} with $\|\mathbf{u}\|_2 = 1$.
 272 Since $\alpha\delta < \Delta$, we can find $r_0 > \alpha\delta$ such that $r_0 < \Delta$ and $r_0 < \delta$. Denote

$$\mathbf{x}_a := \mathbf{x} + r_0 \mathbf{u}. \quad (\text{C.68})$$

273 Clearly, we have $\mathbf{x}_a \in B(\mathbf{x}, \delta)$. Since $\mathcal{N}_{\Delta}(\mathcal{M})$ is a tubular neighborhood of \mathcal{M} , we have

$$d_2(\mathbf{x}_a, \mathcal{M}) = r_0 > \alpha\delta. \quad (\text{C.69})$$

274 For $\forall i \in \{1, 2, \dots, k\}$, by definition, we have

$$f_t^{(i)}(\mathbf{x}_a) = f_b^{*(i)}(\mathbf{x}_a) \cdot \phi_{\text{off}}(\mathbf{x}_a) = \frac{\left(\sum_{l \neq i} d_2(\mathbf{x}_a, A^l)\right) - d_2(\mathbf{x}_a, A^i)}{\left(\sum_{l \neq i} d_2(\mathbf{x}_a, A^l)\right) + d_2(\mathbf{x}_a, A^i)} \cdot \frac{\alpha\delta - d_2(\mathbf{x}_a, \mathcal{M})}{\alpha\delta + d_2(\mathbf{x}_a, \mathcal{M})}. \quad (\text{C.70})$$

275 Since $A_f^1, A_f^2, \dots, A_f^k$ are 2λ -separated, we can easily obtain that

$$y(f_b^*, \mathbf{x}_a) = y(f_b^*, \mathbf{x}). \quad (\text{C.71})$$

276 Combining Equations (C.69) to (C.71) together, we can obtain that \mathbf{x}_a is an off-manifold adversarial
 277 example of f_t at \mathbf{x} , since ϕ_{off} is negative and thus turn the arg max of f_b^* to the arg min. \square

278 Let \mathcal{F}_b and Φ be the function class defined in Proposition C.1 and Proposition C.2, respectively. The
 279 following proposition extends Proposition 4.7 to multi-class classification tasks based on the results
 280 in Proposition C.4.

281 **Proposition C.5.** Denote the target model by $f_t = f_b^* \cdot \phi_{\text{off}}$ and the source model $f_s = f_b \cdot \phi_{\text{off}}$, $f_b \in \mathcal{F}_b$.
 282 With some abuse of notation, let the semantic information of f_b be $A_f^1, A_f^2, \dots, A_f^k$. Let Δ be the
 283 constant specified in Lemma 4.11. Assume that $\alpha\delta < \Delta$. Then, for $\forall \mathbf{x} \in \cup_{i=1}^k A_f^i$, exists adversarial
 284 example of f_s at \mathbf{x} that is transferable if $\mathbf{x} \in \cup_{i=1}^k A_f^i$.

285 In the main part of our paper, Proposition 4.7 proves that adversarial examples are transferable even
 286 if the source model is accurate, which is consistent with the empirical results in Papernot et al. [1].
 287 Proposition C.5 also explains this phenomenon, even though it is weaker than Proposition 4.7.

288 *Proof of Proposition C.5.* For $\forall i \in \{1, 2, \dots, k\}$, we first consider those $\mathbf{x} \in A_f^i$. By definition, we
 289 have

$$f_s^{(l)}(\mathbf{x}) = f_b^{(l)}(\mathbf{x}) \cdot \phi_{\text{off}}(\mathbf{x}) = \frac{(\sum_{l \neq i} d_p(\mathbf{x}, A_f^l)) - d_p(\mathbf{x}, A_f^i)}{(\sum_{l \neq i} d_p(\mathbf{x}, A_f^l)) + d_p(\mathbf{x}, A_f^i)} \cdot \frac{\alpha\delta - d_p(\mathbf{x}, \mathcal{M})}{\alpha\delta + d_p(\mathbf{x}, \mathcal{M})}. \quad (\text{C.72})$$

290 According to Proposition C.4, we know that off-manifold adversarial examples exist. In fact, for
 291 $\forall \mathbf{x} \in \cup_{i=1}^k A_f^i$, let \mathbf{x}_a be as defined in Equation (C.68). Similar to the proof of Proposition C.4, we have
 292 $\phi_{\text{off}}(\mathbf{x}_a) < 0$ and

$$y(f_b^*, \mathbf{x}_a) = y(f_b^*, \mathbf{x}), \quad y(f_b, \mathbf{x}_a) = y(f_b, \mathbf{x}). \quad (\text{C.73})$$

293 which implies that \mathbf{x}_a is a transferable adversarial example. \square

294 References

- 295 [1] Nicolas Papernot, Patrick D. McDaniel, Ian J. Goodfellow, Somesh Jha, Z. Berkay Celik, and
 296 Ananthram Swami. Practical black-box attacks against machine learning. In *AsiaCCS*, 2017.
- 297 [2] Yinpeng Dong, Shuyu Cheng, Tianyu Pang, Hang Su, and Jun Zhu. Query-efficient black-box
 298 adversarial attacks guided by a transfer-based prior. *IEEE Trans. Pattern Anal. Mach. Intell.*,
 299 2022.
- 300 [3] Shuman Fang, Jie Li, Xianming Lin, and Rongrong Ji. Learning to learn transferable attack. In
 301 *AAAI*, 2022.
- 302 [4] Zhichao Huang and Tong Zhang. Black-box adversarial attack with transferable model-based
 303 embedding. In *ICLR*, 2020.
- 304 [5] Minshuo Chen, Haoming Jiang, Wenjing Liao, and Tuo Zhao. Efficient approximation of deep
 305 relu networks for functions on low dimensional manifolds. In *NeurIPS*, pages 8172–8182, 2019.
- 306 [6] Binghui Li, Jikai Jin, Han Zhong, John E. Hopcroft, and Liwei Wang. Why robust generalization
 307 in deep learning is difficult: Perspective of expressive power. *CoRR*, abs/2205.13863, 2022.
- 308 [7] Wei-An Lin, Chun Pong Lau, Alexander Levine, Rama Chellappa, and Soheil Feizi. Dual
 309 manifold adversarial robustness: Defense against lp and non-lp adversarial attacks. In *NeurIPS*,
 310 2020.
- 311 [8] Cassidy Laidlaw, Sahil Singla, and Soheil Feizi. Perceptual adversarial robustness: Defense
 312 against unseen threat models. In *ICLR*, 2021.
- 313 [9] Andrew Ilyas, Shibani Santurkar, Dimitris Tsipras, Logan Engstrom, Brandon Tran, and
 314 Aleksander Madry. Adversarial examples are not bugs, they are features. In *NeurIPS*, 2019.
- 315 [10] Zeyuan Allen-Zhu and Yuanzhi Li. Feature purification: How adversarial training performs
 316 robust deep learning. In *FOCS*, 2021.
- 317 [11] Sanjeev R. Kulkarni, Sanjoy K. Mitter, and John N. Tsitsiklis. Active learning using arbitrary
 318 binary valued queries. *Mach. Learn.*, 1993.
- 319 [12] Sanjoy Dasgupta. Coarse sample complexity bounds for active learning. In *NIPS*, 2005.

- 320 [13] Ev Zisselman and Aviv Tamar. Deep residual flow for out of distribution detection. In *CVPR*,
321 2020.
- 322 [14] Shiyu Liang, Yixuan Li, and R. Srikant. Enhancing the reliability of out-of-distribution image
323 detection in neural networks. In *ICLR*, 2018.
- 324 [15] Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution
325 examples in neural networks. In *ICLR*, 2017.
- 326 [16] Kimin Lee, Kibok Lee, Honglak Lee, and Jinwoo Shin. A simple unified framework for
327 detecting out-of-distribution samples and adversarial attacks. In *NeurIPS*, 2018.
- 328 [17] Zhipeng Wei, Jingjing Chen, Micah Goldblum, Zuxuan Wu, Tom Goldstein, and Yu-Gang Jiang.
329 Towards transferable adversarial attacks on vision transformers. In *AAAI*, 2022.
- 330 [18] Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adver-
331 sarial examples. In *ICLR*, 2015.
- 332 [19] Justin Gilmer, Luke Metz, Fartash Faghri, Samuel S. Schoenholz, Maithra Raghu, Martin
333 Wattenberg, and Ian J. Goodfellow. Adversarial spheres. In *ICLR*, 2018.
- 334 [20] Andrew Ilyas, Logan Engstrom, and Aleksander Madry. Prior convictions: Black-box adversar-
335 ial attacks with bandits and priors. In *ICLR*, 2019.
- 336 [21] Yao-Yuan Yang, Cyrus Rashtchian, Hongyang Zhang, Ruslan Salakhutdinov, and Kamalika
337 Chaudhuri. A closer look at accuracy vs. robustness. In *NeurIPS*, 2020.