
Uncovering the Connections Between Adversarial Transferability and Knowledge Transferability

Kaizhao Liang^{*1} Jacky Y. Zhang^{*1} Boxin Wang¹ Zhuolin Yang¹ Oluwasanmi Koyejo¹ Bo Li¹

Abstract

Knowledge transferability, or transfer learning, has been widely adopted to allow a pre-trained model in the source domain to be effectively adapted to downstream tasks in the target domain. It is thus important to explore and understand the factors affecting knowledge transferability. In this paper, as the first work, we analyze and demonstrate the connections between knowledge transferability and another important phenomenon—adversarial transferability, *i.e.*, adversarial examples generated against one model can be transferred to attack other models. Our theoretical studies show that adversarial transferability indicates knowledge transferability, and vice versa. Moreover, based on the theoretical insights, we propose two practical adversarial transferability metrics to characterize this process, serving as bidirectional indicators between adversarial and knowledge transferability. We conduct extensive experiments for different scenarios on diverse datasets, showing a positive correlation between adversarial transferability and knowledge transferability. Our findings will shed light on future research about effective knowledge transfer learning and adversarial transferability analyses. All code and data are available [here](#).

1. Introduction

Knowledge transfer is quickly becoming the standard approach for fast learning adaptation across domains. Also known as transfer learning or learning transfer, knowledge transfer has been a critical technology for enabling several real-world applications, including object detection (Zhang et al., 2014), image segmentation (Kendall et al., 2018),

multi-lingual machine translation (Dong et al., 2015), and language understanding evaluation (Wang et al., 2019a), among others. For example, since the release of ImageNet (Russakovsky et al., 2015), pretrained ImageNet models (e.g., on TensorFlow Hub or PyTorch-Hub) have become the default option for the knowledge transfer source due to its broad coverage of visual concepts and compatibility with various visual tasks (Huh et al., 2016). Motivated by its importance, many studies have explored the factors associated with knowledge transferability. Most recently, Salman et al. (2020) showed that more robust pretrained ImageNet models transfer better to downstream tasks, which reveals that *adversarial training* helps to improve knowledge transferability.

In the meantime, *adversarial transferability* has been extensively studied—a phenomenon that an adversarial instance generated against one model has high probability attack another one without additional modification (Papernot et al., 2016; Goodfellow et al., 2014; Joon Oh et al., 2017). Hence, adversarial transferability is widely exploited in black-box attacks (Ilyas et al., 2018; Liu et al., 2016; Naseer et al., 2019). A line of work has been conducted to bound the adversarial transferability based on model (gradient) similarity (Tramèr et al., 2017b). Given that both *adversarial transferability* and *knowledge transferability* are impacted by certain model similarity and adversarial ML properties, in this work, we aim to conduct the *first* study to analyze the connections between them and ask,

What is the fundamental connection between knowledge transferability and adversarial transferability? Can we measure one and indicate the other?

Technical Contributions. In this paper, we take the *first* step towards exploring the fundamental relation between adversarial transferability and knowledge transferability. We make contributions on both theoretical and empirical fronts.

- We formally define the adversarial transferability for the *first* time by considering all potential adversarial perturbation vectors. We then conduct thorough and novel theoretical analysis to characterize the precise connection between adversarial transferability and knowledge transferability based on our definition.

^{*}Equal contribution ¹Department of Computer Science, the University of Illinois at Urbana-Champaign, Urbana, USA. Correspondence to: Oluwasanmi Koyejo <sanmi@illinois.edu>, Bo Li <lbo@illinois.edu>.

- In particular, we prove that high adversarial transferability will indicate high knowledge transferability, which can be represented as the distance in an inner product space defined by the Hessian of the adversarial loss. In the meantime, we prove that high knowledge transferability will indicate high adversarial transferability.
- Based on our theoretical insights, we propose two practical adversarial transferability metrics that quantitatively measure the adversarial transferability in practice. We then provide simulational results to verify how these metrics connect with the knowledge transferability in a bidirectional way.
- Extensive experiments justify our theoretical insights and the proposed adversarial transferability metrics, leading to our discussion on potential applications and future research.

Related Work There is a line of research studying different factors that affect knowledge transferability (Yosinski et al., 2014; Long et al., 2015; Wang et al., 2019b; Xu et al., 2019; Shinya et al., 2019). Further, empirical observations show that the correlation between learning tasks (Achille et al., 2019; Zamir et al., 2018), the similarity of model architectures, and data distribution are all correlated with different knowledge transfer abilities. Interestingly, recent empirical evidence suggests that adversarially-trained models transfer better than non-robust models (Salman et al., 2020; Utrera et al., 2020), suggesting a connection between the adversarial properties and knowledge transferability. On the other hand, several approaches have been proposed to boost the adversarial transferability (Zhou et al., 2018; Demontis et al., 2019; Dong et al., 2019; Xie et al., 2019). Beyond the above empirical studies, there are a few existing analyses of adversarial transferability, which explore different conditions that may enhance adversarial transferability (Athalye et al., 2018; Tramèr et al., 2017b; Ma et al., 2018; Demontis et al., 2019). In this work, we aim to bridge the connection between adversarial and knowledge transferability, both of which reveal interesting properties of ML model similarities from different perspectives.

2. Adversarial Transferability and Knowledge Transferability

This section introduces the preliminaries and the formal definitions of the knowledge and adversarial transferability, and formally defines our problem of interest.

Notation. Sets are denoted in blackboard bold, e.g., \mathbb{R} , and the set of integers $\{1 \dots n\}$ is denoted as $[n]$. Distributions are denoted in calligraphy, e.g., \mathcal{D} , and the support of a distribution \mathcal{D} is denoted as $\text{supp}(\mathcal{D})$. Vectors are denoted as bold lower case letters, e.g., $\mathbf{x} \in \mathbb{R}^n$, and matrices are denoted as bold uppercase letters, e.g., \mathbf{W} . We denote the

entry-wise product operator between vectors or matrices as \odot . The Moore–Penrose inverse of a matrix \mathbf{W} is denoted as \mathbf{W}^\dagger . We use $\|\cdot\|_2$ to denote Euclidean norm induced by Euclidean inner product $\langle \cdot, \cdot \rangle$. The standard inner product of two matrices is defined as $\langle \mathbf{W}, \mathbf{M} \rangle = \text{tr}(\mathbf{W}^\top \mathbf{M})$, where $\text{tr}(\cdot)$ is the trace of a matrix. The Frobenius norm $\|\cdot\|_F$ is induced by the standard matrix inner product. Moreover, in the (semi-)inner product space defined by a positive (semi-)definite matrix \mathbf{S} , the (semi-)inner product of two vectors or matrices is defined by $\langle \mathbf{v}_1, \mathbf{v}_2 \rangle_{\mathbf{S}} = \mathbf{v}_1^\top \mathbf{S} \mathbf{v}_2$ or $\langle \mathbf{W}, \mathbf{M} \rangle_{\mathbf{S}} = \text{tr}(\mathbf{W}^\top \mathbf{S} \mathbf{M})$, respectively. Given a vector \mathbf{v} , we define its normalization as $\hat{\mathbf{v}} = \mathbf{v} / \|\mathbf{v}\|_2$. When using a denominator $\|\cdot\|_*$ other than Euclidean norm, we denote the normalization as $\hat{\mathbf{v}}|_*$.

Given a (vector-valued) function f , we denote $f(\mathbf{x})$ as its evaluated value at \mathbf{x} , and f represents the function itself in the corresponding Hilbert space. Composition of functions is denoted as $g \circ f(\mathbf{x}) = g(f(\mathbf{x}))$. We use $\langle \cdot, \cdot \rangle_{\mathcal{D}}$ to denote the inner product induced by distribution \mathcal{D} and inherited from Euclidean inner product, i.e., $\langle f_1, f_2 \rangle_{\mathcal{D}} = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \langle f_1(\mathbf{x}), f_2(\mathbf{x}) \rangle$. Accordingly, we use $\|\cdot\|_{\mathcal{D}}$ to denote the norm induced by the inner product $\langle \cdot, \cdot \rangle_{\mathcal{D}}$, i.e., $\|f\|_{\mathcal{D}} = \sqrt{\langle f, f \rangle_{\mathcal{D}}}$. When the inherited inner product is defined by \mathbf{S} , we denote $\langle f_1, f_2 \rangle_{\mathcal{D}, \mathbf{S}} = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \langle f_1(\mathbf{x}), f_2(\mathbf{x}) \rangle_{\mathbf{S}}$, and similarly for $\|f\|_{\mathcal{D}, \mathbf{S}}$.

Knowledge Transferability Given a pre-trained *source* model $f_S : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and a *target* domain $\mathbf{x} \in \mathbb{R}^n$ with data distribution $\mathbf{x} \sim \mathcal{D}$ and *target* labels $y(\mathbf{x}) \in \mathbb{R}^d$, *knowledge transferability* is defined as the performance of fine-tuning f_S on \mathcal{D} to predict y . Concretely, knowledge transferability can be represented as a loss $\mathcal{L}(\cdot, y, \mathcal{D})$ after fine-tuning by composing the fixed source model with a trainable function $g : \mathbb{R}^m \rightarrow \mathbb{R}^d$, typically from a small function class $g \in \mathbb{G}$, i.e.,

$$\min_{g \in \mathbb{G}} \mathcal{L}(g \circ f_S, y, \mathcal{D}), \quad (1)$$

where the loss function \mathcal{L} measures the error between $g \circ f_S$ and the ground truth y under the *target* data distribution \mathcal{D} . For example, for neural networks it is usual to stack on and fine-tune a linear layer; here \mathbb{G} is the affine function class. We will focus on the affine setting in this paper.

For our purposes, a more useful measure of transfer is to compare the quality of the fine-tuned model to a model trained directly on the target domain $f_T : \mathbb{R}^n \rightarrow \mathbb{R}^d$. Thus, we study the following surrogate of knowledge transferability, where the ground truth target is replaced by a reference target model f_T :

$$\min_{g \in \mathbb{G}} \mathcal{L}(g \circ f_S, f_T, \mathcal{D}). \quad (2)$$

Adversarial Attacks. For simplicity we consider untargeted attacks that seeks to maximize the deviation of model

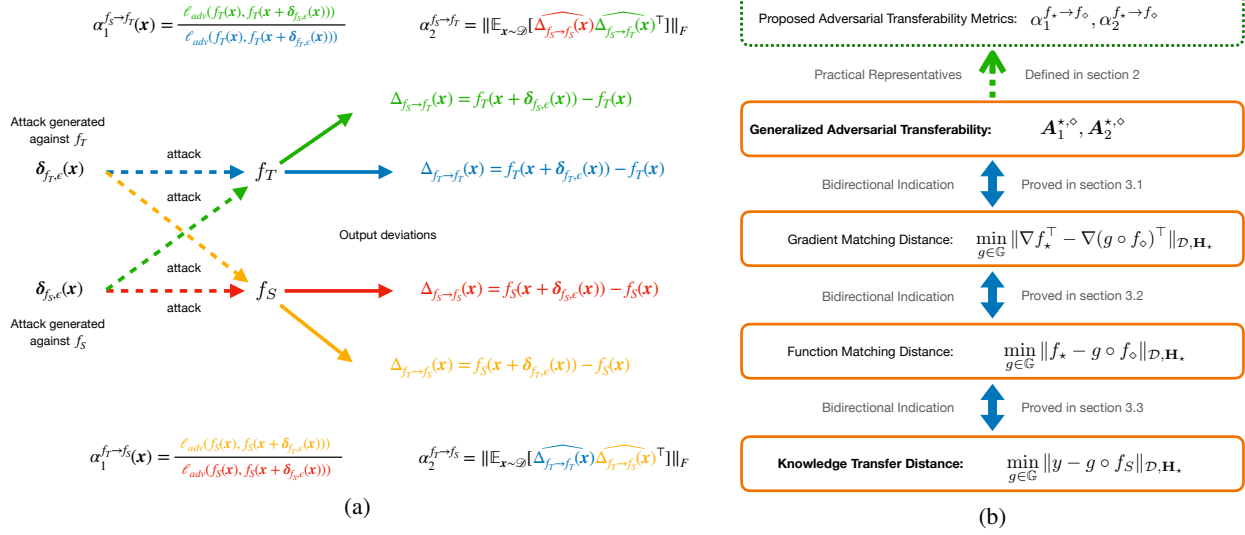


Figure 1. (a) An illustration of the two proposed adversarial transferability metrics α_1, α_2 under different adversarial transferability settings, i.e., $\alpha_1^{f_S \rightarrow f_T}, \alpha_1^{f_T \rightarrow f_S}, \alpha_2^{f_S \rightarrow f_T}$, and $\alpha_2^{f_T \rightarrow f_S}$. (b) An overview of the theoretical analysis framework, and its practical inspirations, where $*, \diamond \in \{T, S\}$ and $* \neq \diamond$. The three blue double-headed arrows are the bidirectional indication relationships proved in our theory section, and the dashed green arrow shows in practice how the two proposed adversarial transferability metrics are measured as representatives of the generalized adversarial transferability based on our theory.

output as measured by a given adversarial loss function $\ell_{adv}(\cdot, \cdot)$. The targeted attack can be viewed as a special case. Without loss of generality, we assume the adversarial loss is *non-negative*. Given a datapoint \mathbf{x} and model f , an adversarial example of magnitude ϵ is denoted by $\delta_{f, \epsilon}(\mathbf{x})$, computed as:

$$\delta_{f, \epsilon}(\mathbf{x}) = \arg \max_{\|\delta\| \leq \epsilon} \ell_{adv}(f(\mathbf{x}), f(\mathbf{x} + \delta)). \quad (3)$$

We note that in theory $\delta_{f, \epsilon}(\mathbf{x})$ may not be unique, and its generalized definition and its discussion are provided in our theoretical analysis (Section 3).

Adversarial Transferability. The process of adversarial transfer involves applying the adversarial example generated against a model f_1 to another model f_2 . Thus, adversarial transferability from f_1 to f_2 measures how well $\delta_{f_1, \epsilon}$ attacks f_2 . We propose two metrics, namely, α_1 and α_2 that characterize adversarial transferability from complementary perspectives. To provide a visual overview of our definitions for the proposed adversarial transferability metrics, we present an illustration in Figure 1 (a).

Definition 1 (The First Adversarial Transferability). *The first adversarial transferability from f_1 to f_2 at data sample $\mathbf{x} \sim \mathcal{D}$, is defined as*

$$\alpha_1^{f_1 \rightarrow f_2}(\mathbf{x}) = \frac{\ell_{adv}(f_2(\mathbf{x}), f_2(\mathbf{x} + \delta_{f_1, \epsilon}(\mathbf{x})))}{\ell_{adv}(f_2(\mathbf{x}), f_2(\mathbf{x} + \delta_{f_2, \epsilon}(\mathbf{x})))}.$$

Taking the expectation, the first adversarial transferability

is defined as

$$\alpha_1^{f_1 \rightarrow f_2} = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\alpha_1^{f_1 \rightarrow f_2}(\mathbf{x})].$$

Observe that the first adversarial transferability characterizes how well the adversarial attacks $\delta_{f_1, \epsilon}$ generated against f_1 perform on f_2 , compared to f_2 's whitebox adversarial attacks $\delta_{f_2, \epsilon}$. Thus, high α_1 indicates high adversarial transferability. Note that the two attacks use the same magnitude constraint ϵ .

Recall that $\ell_{adv}(f(\mathbf{x}), f(\mathbf{x} + \delta))$ measures the effect of the attack δ on the model output $f(\mathbf{x})$. α_1 characterizes the relative magnitude of this deviation. However, this magnitude information is incomplete, as the direction of the deviation also encodes information about the adversarial transfer process. To this end, we propose the second adversarial metric, inspired by our theoretical analysis, which characterizes adversarial transferability from the directional perspective.

Definition 2 (The Second Adversarial Transferability). *The second adversarial transferability from f_1 to f_2 , under data distribution $\mathbf{x} \sim \mathcal{D}$, is defined as*

$$\alpha_2^{f_1 \rightarrow f_2} = \|\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[\widehat{\Delta_{f_1 \rightarrow f_1}}(\mathbf{x}) \widehat{\Delta_{f_1 \rightarrow f_2}}(\mathbf{x})^\top]\|_F,$$

where

$$\Delta_{f_1 \rightarrow f_1}(\mathbf{x}) = f_1(\mathbf{x} + \delta_{f_1, \epsilon}(\mathbf{x})) - f_1(\mathbf{x})$$

$$\Delta_{f_1 \rightarrow f_2}(\mathbf{x}) = f_2(\mathbf{x} + \delta_{f_1, \epsilon}(\mathbf{x})) - f_2(\mathbf{x})$$

are deviations in model output given the adversarial attack $\delta_{f_1, \epsilon}(\mathbf{x})$ generated against f_1 , and $\widehat{\cdot}$ denotes the corresponding unit-length vector.

To further clarify the second adversarial transferability metric, consider the following alternative form of α_2 .

Proposition 2.1. *The $\alpha_2^{f_1 \rightarrow f_2}$ can be reformulated as*

$$(\alpha_2^{f_1 \rightarrow f_2})^2 = \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2} [\theta_{f_1 \rightarrow f_1}(\mathbf{x}_1, \mathbf{x}_2) \theta_{f_1 \rightarrow f_2}(\mathbf{x}_1, \mathbf{x}_2)],$$

where $\mathbf{x}_1, \mathbf{x}_2 \stackrel{i.i.d.}{\sim} \mathcal{D}$, and

$$\theta_{f_1 \rightarrow f_1}(\mathbf{x}_1, \mathbf{x}_2) = \langle \widehat{\Delta_{f_1 \rightarrow f_1}(\mathbf{x}_1)}, \widehat{\Delta_{f_1 \rightarrow f_1}(\mathbf{x}_2)} \rangle$$

$$\theta_{f_1 \rightarrow f_2}(\mathbf{x}_1, \mathbf{x}_2) = \langle \widehat{\Delta_{f_1 \rightarrow f_2}(\mathbf{x}_1)}, \widehat{\Delta_{f_1 \rightarrow f_2}(\mathbf{x}_2)} \rangle$$

We can see that high α_2 indicates that it is more likely for the two inner products (*i.e.*, $\theta_{f_1 \rightarrow f_1}$ and $\theta_{f_1 \rightarrow f_2}$) to have the same sign. Given that the direction of f_1 's output deviation indicates its attack $\delta_{f_1, \epsilon}$, and the direction of f_2 's output deviation indicates the transferred attack $\delta_{f_1, \epsilon}$, high α_2 implies that the two directions will rotate by a similar angle as the data changes.

α_1 and α_2 represent complementary aspects of the adversarial transferability: α_1 can be understood as how often the adversarial attack transfers, while α_2 encodes directional information of the output deviation caused by adversarial attacks. An example is provided in the appendix section A to illustrate the necessity of both the metrics in characterizing the relation between adversarial transferability and knowledge transferability. To jointly take the two adversarial transferability metrics into consideration, we propose the following metric as the combined value of α_1 and α_2 .

$$(\alpha_1 * \alpha_2)^{f_1 \rightarrow f_2} = \left\| \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\alpha_1^{f_1 \rightarrow f_2}(\mathbf{x}) \widehat{\Delta_{f_1 \rightarrow f_1}(\mathbf{x})} \widehat{\Delta_{f_1 \rightarrow f_2}(\mathbf{x})}^\top] \right\|_F.$$

We defer the justification for the combined adversarial transferability metric in the next section, and move on to state a useful proposition.

Proposition 2.2. *The adversarial transferability metrics $\alpha_1^{f_1 \rightarrow f_2}$, $\alpha_2^{f_1 \rightarrow f_2}$ and $(\alpha_1 * \alpha_2)^{f_1 \rightarrow f_2}$ are in $[0, 1]$.*

So far, we have defined knowledge transferability, and two adversarial transferability metrics. We can now analyze their connections more precisely.

Problem of Interest. Given a *source* model $f_S : \mathbb{R}^n \rightarrow \mathbb{R}^m$, the *target* data distribution $\mathbf{x} \sim \mathcal{D}$, the ground truth target $y : \mathbb{R}^n \rightarrow \mathbb{R}^d$, and a *target* reference model $f_T : \mathbb{R}^n \rightarrow \mathbb{R}^d$, we aim to study how the adversarial transferability between f_S and f_T , characterized by the two proposed adversarial transferability metrics, connects to the knowledge transfer loss $\min_{g \in \mathbb{G}} \mathcal{L}(g \circ f_S, y, \mathcal{D})$ with affine functions $g \in \mathbb{G}$ (equation 1).

3. Theoretical Analysis

In this section, we present the theoretical analysis on how the adversarial transferability and the knowledge transfer process are tied together. To simplify the discussion, as the objects studied in this section are specifically focused on the source domain S and the target domain T , we can use \star or \diamond as a placeholder for either S or T throughout this section.

Theoretical Analysis Overview. In subsection 3.1, we define the two *generalized adversarial transferabilities*, (*i.e.*, $\mathbf{A}_1, \mathbf{A}_2$), and present Theorem 3.1 showing that $\mathbf{A}_1, \mathbf{A}_2$ together determine a gradient matching distance $\min_{g \in \mathbb{G}} \|\nabla f_\star - \nabla g \circ f_\diamond\|$, between the Jacobian matrices of the source and target models in an inner product space defined by the Hessian of the adversarial loss function. In the same subsection, we also show that α_1 and α_2 represent the most influential factors in \mathbf{A}_1 and \mathbf{A}_2 , respectively. Next, we explore the connection to knowledge transferability in subsection 3.2 via Theorem 3.2 which shows the gradient matching distance approximates the function matching distance, *i.e.*, $\min_{g \in \mathbb{G}} \|f_\star - g \circ f_\diamond\|$, with a distribution shift up to a Wasserstein distance. Finally, in subsection 3.3 we complete the analysis by outlining the connection between the function matching distance and the knowledge transfer loss. A visual overview is shown in Figure 1 (b).

Setting. As adversarial perturbations are constrained in a small ϵ -ball, it is reasonable to approximate the deviation of model outputs by its first-order Taylor approximation. Specifically, in this section we consider the Euclidean ϵ -ball. Therefore, the output deviation of a function f at \mathbf{x} given a small perturbation $\|\delta_\epsilon\|_2 \leq \epsilon$ can be approximated by

$$f(\mathbf{x} + \delta_\epsilon) - f(\mathbf{x}) \approx \nabla f(\mathbf{x})^\top \delta_\epsilon,$$

where $\nabla f(\mathbf{x})$ is the Jacobian matrix of f at \mathbf{x} .

We consider a convex and twice-differentiable adversarial loss function $\ell_{adv}^*(\cdot)$ that measures the deviation of model output $f_\star(\mathbf{x} + \delta_\epsilon) - f_\star(\mathbf{x})$, with minimum $\ell_{adv}^*(\mathbf{0}) = 0$, for $\star \in \{S, T\}$. We note that we should treat the adversarial loss on f_S and f_T differently, as they may have different output dimensions. Accordingly, the adversarial attack (equation 3) can be written as

$$\delta_{f_\star, \epsilon}(\mathbf{x}) = \arg \max_{\|\delta\|_2 \leq \epsilon} \ell_{adv}^*(\nabla f_\star(\mathbf{x})^\top \delta). \quad (4)$$

Another justification of the small- ϵ approximation follows the literature; since the ideal attack defined in equation 3 is often intractable to compute, much of the literature uses the proposed formulation (4) in practice, *e.g.*, see (Miyato et al., 2018), with experimental results suggesting similar behaviour as the standard definition.

The Small- ϵ Regime. Recall that the adversarial loss $\ell_{adv}^*(\cdot)$ studied in this section is convex, twice-differentiable,

and achieves its minimum at $\mathbf{0}$, thus in the small ϵ regime:

$$\begin{aligned} \ell_{adv}^*(\nabla f_*(\mathbf{x})^\top \delta_\epsilon) &= (\delta_\epsilon^\top \nabla f_*(\mathbf{x}) \mathbf{H}_* \nabla f_*(\mathbf{x})^\top \delta_\epsilon)^{1/2} \\ &= \|\nabla f_*(\mathbf{x})^\top \delta_\epsilon\|_{\mathbf{H}_*}, \end{aligned}$$

which is the norm of f_* 's output deviation in the inner product space defined by the Hessian \mathbf{H}_* of the squared adversarial loss $(\ell_{adv}^*)^2$.

Accordingly, the adversarial attacks (4) can be written as

$$\delta_{f_*,\epsilon}(\mathbf{x}) = \arg \max_{\|\delta\|_2 \leq \epsilon} \|\nabla f_*(\mathbf{x})^\top \delta\|_{\mathbf{H}_*}, \quad (5)$$

and we can measure the output deviation of f_\diamond 's caused by f_* 's adversarial attack $\delta_{f_*,\epsilon}(\mathbf{x})$, denoted as:

$$\Delta_{f_* \rightarrow f_\diamond, \epsilon}(\mathbf{x}) = \nabla f_\diamond(\mathbf{x})^\top \delta_{f_*,\epsilon}(\mathbf{x}). \quad (6)$$

Note that in the small- ϵ regime, the actual value of ϵ becomes trivial (e.g., α_1), consequently we will omit the ϵ for notational ease:

$$\alpha_1^{f_* \rightarrow f_\diamond}(\mathbf{x}) = \frac{\|\Delta_{f_* \rightarrow f_\diamond}(\mathbf{x})\|_{\mathbf{H}_\diamond}}{\|\nabla f_\diamond(\mathbf{x})\|_{\mathbf{H}_\diamond}}.$$

Similarly, α_2 can be computed using (6) in Definition 2, i.e.,

$$\alpha_2^{f_* \rightarrow f_\diamond} = \|\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[\widehat{\Delta_{f_* \rightarrow f_*}(\mathbf{x})} \widehat{\Delta_{f_* \rightarrow f_\diamond}(\mathbf{x})}^\top]\|_F.$$

With these insights, next we will derive our first theorem.

3.1. Adversarial Transfer Indicates the Gradient Matching Distance, and Vice Versa

We present an interesting finding in this subsection, i.e., the generalized adversarial transferabilities $\mathbf{A}_1, \mathbf{A}_2$ have a direct connection to the gradient matching distance between the source model $f_S : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and target model $f_T : \mathbb{R}^n \rightarrow \mathbb{R}^d$. The gradient matching distance is defined as the smallest distance an affine transformation can achieve between their Jacobians $\nabla f_T : \mathbb{R}^n \rightarrow \mathbb{R}^{n \times d}$ and $\nabla f_S : \mathbb{R}^n \rightarrow \mathbb{R}^{n \times m}$ in the inner product space defined by \mathbf{H}_* and data sample distribution $\mathbf{x} \sim \mathcal{D}$, as shown below.

$$\min_{g \in \mathbb{G}} \|\nabla f_*^\top - \nabla(g \circ f_\diamond)^\top\|_{\mathcal{D}, \mathbf{H}_*}, \quad (7)$$

where $g \in \mathbb{G}$ are affine transformations. Note that $g : \mathbb{R}^m \rightarrow \mathbb{R}^d$ if $(\star, \diamond) = (T, S)$, and $g : \mathbb{R}^d \rightarrow \mathbb{R}^m$ if $(\star, \diamond) = (S, T)$. We defer the analysis of how the gradient matching distance approximates the knowledge transfer loss, and focus on its connection to adversarial transfer.

A Full Picture of Adversarial Transferability. A key observation is that the adversarial attack (equation 5) is the singular vector corresponding to the largest singular value of the Jacobian $\nabla f_*(\mathbf{x})$ in the \mathbf{H}_* inner product space.

Thus, information regarding other singular values that are not revealed by the adversarial attack. Therefore, we can consider other singular values, corresponding to smaller signals than the one revealed by adversarial attacks, to complete the analysis. We denote $\sigma_{f_*, \mathbf{H}_*}(\mathbf{x}) \in \mathbb{R}^n$ as the descending (in absolute value) singular values of the Jacobian $\nabla f_*(\mathbf{x})^\top \in \mathbb{R}^{n \times n}$ in the \mathbf{H}_* inner product space. In other words, we denote $\sigma_{f_*, \mathbf{H}_*}(\mathbf{x}) \in \mathbb{R}^n$ as the square root of the descending eigenvalues of $\nabla f_*(\mathbf{x}) \mathbf{H}_* \nabla f_*(\mathbf{x})^\top$, i.e.,

$$\sigma_{f_*, \mathbf{H}_*}(\mathbf{x}) = [\sigma_{f_*}^{(1)}(\mathbf{x}), \dots, \sigma_{f_*}^{(n)}(\mathbf{x})]^\top. \quad (8)$$

Note that the number of non-zero singular values may be less than n , in which case we fill the rest with zeros such that vector is n -dimensional.

Since the adversarial attack $\delta_{f_*,\epsilon}(\mathbf{x})$ corresponds to the largest singular value $\sigma_{f_*}(\mathbf{x})^{(1)}$, we can also generalize the adversarial attack by including all the singular vectors. i.e.,

$$\delta_{f_*}^{(i)}(\mathbf{x}) \quad \text{corresponds to} \quad \sigma_{f_*}^{(i)}(\mathbf{x}), \quad \forall i \in [n]. \quad (9)$$

Loosely speaking, one could think $\delta_{f_*}^{(i)}(\mathbf{x})$ as the adversarial attack of $f_*(\mathbf{x})$ in the subspace orthogonal to all the previous attacks, i.e., $\delta_{f_*}^{(j)}(\mathbf{x})$ for $\forall j \in [i-1]$.

Accordingly, for $\forall i \in [i]$ we denote the output deviation as

$$\Delta_{f_* \rightarrow f_\diamond}^{(i)}(\mathbf{x}) = \nabla f_\diamond(\mathbf{x})^\top \delta_{f_*}^{(i)}(\mathbf{x}). \quad (10)$$

As a consequence, we generalize the first adversarial transferability to be a n -dimensional vector $\mathbf{A}_1^{*,\diamond}(\mathbf{x})$ including the adversarial losses of all of the generalized adversarial attacks, where the i^{th} element in the vector is

$$\mathbf{A}_1^{*,\diamond}(\mathbf{x})^{(i)} = \frac{\|\Delta_{f_* \rightarrow f_\diamond}^{(i)}(\mathbf{x})\|_{\mathbf{H}_\diamond}}{\|\nabla f_\diamond(\mathbf{x})\|_{\mathbf{H}_\diamond}}. \quad (11)$$

Note that the first entry of $\mathbf{A}_1^{*,\diamond}(\mathbf{x})$ is the original adversarial transferability, i.e., $\mathbf{A}_1^{*,\diamond}(\mathbf{x})^{(1)}$ is the same as the $\alpha_1^{f_* \rightarrow f_\diamond}(\mathbf{x})$ in Definition 1.

With the above generalization that captures the full picture of the adversarial transfer process, we able to derive the following theorem.

Theorem 3.1. *Given the target and source models f_*, f_\diamond , where $(\star, \diamond) \in \{(S, T), (T, S)\}$, the gradient matching distance (equation 7) can be written as*

$$\begin{aligned} \min_{g \in \mathbb{G}} \|\nabla f_*^\top - \nabla(g \circ f_\diamond)^\top\|_{\mathcal{D}, \mathbf{H}_*} &= \quad (12) \\ \sqrt{1 - \frac{\mathbb{E}[\mathbf{v}^{*,\diamond}(\mathbf{x}_1)^\top \mathbf{A}_2^{*,\diamond}(\mathbf{x}_1, \mathbf{x}_2) \mathbf{v}^{*,\diamond}(\mathbf{x}_2)]}{\|\nabla f_*^\top\|_{\mathcal{D}, \mathbf{H}_*}^2 \cdot \|\mathbf{J}^\dagger\|_{\mathbf{H}_\diamond}^{-1}}} \|\nabla f_*^\top\|_{\mathcal{D}, \mathbf{H}_*}, \end{aligned}$$

where the expectation is taken over $\mathbf{x}_1, \mathbf{x}_2 \stackrel{i.i.d.}{\sim} \mathcal{D}$, and

$$\begin{aligned} \mathbf{v}^{*,\diamond}(\mathbf{x}) &= \sigma_{f_\diamond, \mathbf{H}_\diamond}^{(1)}(\mathbf{x}) \sigma_{f_*, \mathbf{H}_*}(\mathbf{x}) \odot \mathbf{A}_1^{*,\diamond}(\mathbf{x}) \\ \mathbf{J} &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[\nabla f_\diamond(\mathbf{x})^\top \nabla f_\diamond(\mathbf{x})]. \end{aligned}$$

Moreover, $\mathbf{A}_2^{*,\diamond}(\mathbf{x}_1, \mathbf{x}_2)$ is a matrix, and its element in the i^{th} row and j^{th} column is

$$\mathbf{A}_2^{*,\diamond}(\mathbf{x}_1, \mathbf{x}_2)^{(i,j)} = \langle \widehat{\Delta_{f_\star \rightarrow f_\star}^{(i)}}(\mathbf{x}_1)|_{\mathbf{H}_\star}, \widehat{\Delta_{f_\star \rightarrow f_\star}^{(j)}}(\mathbf{x}_2)|_{\mathbf{H}_\star} \rangle \cdot \langle \widehat{\Delta_{f_\star \rightarrow f_\diamond}^{(i)}}(\mathbf{x}_1)|_{\mathbf{H}_\diamond}, \widehat{\Delta_{f_\star \rightarrow f_\diamond}^{(j)}}(\mathbf{x}_2)|_{\mathbf{H}_\diamond} \rangle_{\widehat{\mathbf{J}^\dagger}|_{\mathbf{H}_\diamond}}.$$

Recall the alternative representation of the second adversarial transferability α_2 , and we can immediately observe that α_2 is determined by \mathbf{A}_2 . Therefore, both α_1 and α_2 appear in this relation. Let us interpret the theorem, and justify the two proposed adversarial transferability metrics.

Interpretation of Theorem 3.1. First, we consider components that are not directly related to the adversarial transfer in the RHS of (12). The $\|\nabla f_\star^\top\|_{\mathcal{D}, \mathbf{H}_\star}$ outside represents the overall magnitude of the loss. In the fraction, the $\|\nabla f_\star^\top\|_{\mathcal{D}, \mathbf{H}_\star}$ in the denominator normalizes the σ_{f_\star} in the numerator. Similarly, though more complicated, the $\|\mathbf{J}^\dagger\|_2^{-1}$ in the denominator corresponds to the $\sigma_{f_\diamond}^{(1)}$ in the numerator. We note that these are properties of f_\star, f_\diamond .

Next, observe that the components directly related to the adversarial transfer process are the *generalized adversarial transferability* \mathbf{A}_1 and \mathbf{A}_2 . Let us neglect the superscript (i) or (j) for now, so we can see that their interpretations are the same as we introduced for α_1 and α_2 in section 2. That is, \mathbf{A}_1 captures the magnitude of the deviation in model outputs caused by adversarial attacks, while \mathbf{A}_2 captures the direction of the deviation. A minor difference between α_2 and \mathbf{A}_2 is that the second inner product in the elements of \mathbf{A}_2 is defined by a positive semi-definite matrix $\widehat{\mathbf{J}^\dagger}$. For practical implementation, we choose to neglect this term, and use the standard Euclidean inner product in α_2 , which can be understood as a stretched version of the $\widehat{\mathbf{J}^\dagger}$ inner product space.

Moreover, as the singular vector σ_{f_\star} has descending entries, we can see that in the vector \mathbf{A}_1 and the matrix \mathbf{A}_2 , the elements with superscript (1) have the most influence in the relations. In other words, the two proposed adversarial transferability metrics, α_1 and α_2 , are the most influential factors in equation 12. We can also see that the combined metric $(\alpha_1 * \alpha_2)$ also stems from here by only considering the components with the first superscript.

To interpret the relation between the gradient matching distance and the adversarial transferabilities, we introduce the following proposition. This shows that, in general, \mathbf{A}_1 and \mathbf{A}_2 with their elements closer to 1 can serve as a bidirectional indicator of a smaller gradient matching distance.

Proposition 3.1. *In Theorem 3.1,*

$$0 \leq \frac{\mathbb{E}[\mathbf{v}^{*,\diamond}(\mathbf{x}_1)^\top \mathbf{A}_2^{*,\diamond}(\mathbf{x}_1, \mathbf{x}_2) \mathbf{v}^{*,\diamond}(\mathbf{x}_2)]}{\|\nabla f_\star^\top\|_{\mathcal{D}, \mathbf{H}_\star}^2 \cdot \|\mathbf{J}^\dagger\|_{\mathbf{H}_\diamond}^{-1}} \leq 1.$$

In conclusion, Theorem 3.1 reveals a bidirectional relation between the adversarial transfer process and the gradient matching distance, where the adversarial transfer process can be encoded by the generalized adversarial transferabilities, *i.e.*, \mathbf{A}_1 and \mathbf{A}_2 . Moreover, α_1 and α_2 play the most influential role in their generalization, *i.e.*, \mathbf{A}_1 and \mathbf{A}_2 .

3.2. The Gradient Matching Distance indicates the Function Matching Distance, and Vice Versa

To bridge the gradient matching distance to the knowledge transfer loss, an immediate step is to connect the gradient distance to the function distance which directly serves as a surrogate knowledge transfer loss as defined in (equation 2). Specifically, in this subsection, we present a connection between the function matching distance, *i.e.*,

$$\min_{g \in \mathbb{G}} \|f_\star - g \circ f_\diamond\|_{\mathcal{D}, \mathbf{H}_\star}, \quad (13)$$

and the gradient matching distance, *i.e.*,

$$\min_{g \in \mathbb{G}} \|\nabla f_\star^\top - \nabla(g \circ f_\diamond)^\top\|_{\mathcal{D}, \mathbf{H}_\star}, \quad (14)$$

where $g \in \mathbb{G}$ are affine transformations.

For intuition, consider a point \mathbf{x}_0 in the input space \mathbb{R}^n , a path $\gamma_{\mathbf{x}} : [0, 1] \rightarrow \mathbb{R}^n$ such that $\gamma_{\mathbf{x}}(0) = \mathbf{x}_0$ and $\gamma_{\mathbf{x}}(1) = \mathbf{x}$. Then, denoting γ as the function of \mathbf{x} , we can write the difference between the two functions as

$$f_\star - g \circ f_\diamond = \int_0^1 (\nabla f_\star(\gamma(t)) - \nabla(g \circ f_\diamond(\gamma(t))))^\top \dot{\gamma}(t) dt + (f_\star(\mathbf{x}_0) - g \circ f_\diamond(\mathbf{x}_0)).$$

Noting that the function difference is a path integral of the gradient difference, we should expect a distribution shift when characterizing their connection, *i.e.*, the integral path affects the original data distribution \mathcal{D} . Accordingly, as the integral path may leave the support of \mathcal{D} , it is necessary to assume the smoothness of the function, as shown below.

Denoting the optimal $g \in \mathbb{G}$ in (13) as \tilde{g} , and one of the optimal $g \in \mathbb{G}$ in (14) as \tilde{g}' , we define

$$h_{\star, \diamond} := f_\star - \tilde{g} \circ f_\diamond \quad \text{and} \quad h'_{\star, \diamond} := f_\star - \tilde{g}' \circ f_\diamond, \quad (15)$$

and we can see that the gradient matching distance and the function matching distance can be written as

$$(13) = \|h_{\star, \diamond}\|_{\mathcal{D}, \mathbf{H}_\star} \quad \text{and} \quad (14) = \|\nabla h'_{\star, \diamond}\|_{\mathcal{D}, \mathbf{H}_\star}.$$

Assumption 1 (β -smoothness). *We assume $h_{\star, \diamond}$ and $h'_{\star, \diamond}$ are both β -smooth, *i.e.*,*

$$\|\nabla h_{\star, \diamond}^\top(\mathbf{x}_1) - \nabla h_{\star, \diamond}^\top(\mathbf{x}_2)\|_{\mathbf{H}_\star} \leq \|\mathbf{x}_1 - \mathbf{x}_2\|_2,$$

and similarly for $h'_{\star, \diamond}$.

With this assumption, we can prove that the gradient matching distance and the function matching distance can bound each other.

Theorem 3.2. *With the notation defined in equation 15, assume the β -smoothness assumption holds. Given a data distribution \mathcal{D} and $\tau > 0$, there exist distributions $\mathcal{D}_1, \mathcal{D}_2$ such that the type-1 Wasserstein distance $W_1(\mathcal{D}, \mathcal{D}_1) \leq \tau$ and $W_1(\mathcal{D}, \mathcal{D}_2) \leq \tau$ satisfying*

$$\begin{aligned} \frac{1}{2B^2} \|h_{\star, \diamond}\|_{\mathcal{D}, \mathbf{H}_\star}^2 &\leq \|\nabla h_{\star, \diamond}^\top\|_{\mathcal{D}_1, \mathbf{H}_\star}^2 + \beta^2(B - \tau)^2 \\ \frac{1}{3n} \|\nabla h_{\star, \diamond}^\top\|_{\mathcal{D}, \mathbf{H}_\star}^2 &\leq \frac{2}{\tau^2} \|h_{\star, \diamond}\|_{\mathcal{D}_2, \mathbf{H}_\star}^2 + \beta^2\tau^2, \end{aligned}$$

where n is the dimension of $\mathbf{x} \sim \mathcal{D}$, and $B = \inf_{\mathbf{x}_0 \in \mathbb{R}^n} \sup_{\mathbf{x} \in \text{supp}(\mathcal{D})} \|\mathbf{x} - \mathbf{x}_0\|_2$ is the radius of $\text{supp}(\mathcal{D})$.

We note that the above theorem compromises some tightness in exchange for a cleaner presentation without losing its core message, which is discussed in the proof of the theorem.

Interpretation of Theorem 3.2. The theorem shows that under the smoothness assumption, the gradient matching distance indicates the function matching distance, and vice versa, with a distribution shift bounded in Wasserstein distance. As the distribution shift is in general necessary, we conjecture that using different data distributions for adversarial transfer and knowledge transfer can also be applicable.

3.3. The Function Matching Distance Indicates Knowledge Transferability, and Vice Versa

To complete the story, it remains to connect the function matching distance to knowledge transferability. As the adversarial transfer is symmetric (*i.e.*, either from $f_S \rightarrow f_T$ or $f_T \rightarrow f_S$), we are able to use the placeholders $\star, \diamond \in \{S, T\}$ all the way through. However, as the knowledge transfer is asymmetric (*i.e.*, $f_S \rightarrow y$ to the target ground truth), we need to instantiate the direction of adversarial transfer to further our discussion.

Adversarial Transfer from $f_T \rightarrow f_S$. As we can see from the $\mathbf{A}_1^{\star, \diamond}$ in Theorem 3.1, this direction corresponds to $(\star, \diamond) = (T, S)$. Accordingly, the function matching distance (equation 13) becomes

$$\min_{g \in \mathbb{G}} \|f_T - g \circ f_S\|_{\mathcal{D}, \mathbf{H}_T}. \quad (16)$$

We can see that equation 16 directly translates to the surrogate knowledge transfer loss that uses the ‘‘pseudo ground truth’’ from the target reference model f_T .

In other words, the function matching distance serves as an approximation of the knowledge transfer loss defined as their distance in the inner product space of \mathbf{H}_T , *i.e.*,

$$\min_{g \in \mathbb{G}} \|y - g \circ f_S\|_{\mathcal{D}, \mathbf{H}_T}. \quad (17)$$

The accuracy of the approximation depends on the performance of f_T , as shown in the following theorem.

Theorem 3.3. *The surrogate transfer loss (16) and the true transfer loss (17) are close, with an error of $\|f_T - y\|_{\mathcal{D}, \mathbf{H}_T}$.*

$$-\|f_T - y\|_{\mathcal{D}, \mathbf{H}_T} \leq (17) - (16) \leq \|f_T - y\|_{\mathcal{D}, \mathbf{H}_T}$$

Adversarial Transfer from $f_S \rightarrow f_T$. This direction corresponds to $(\star, \diamond) = (S, T)$. Accordingly, the function matching distance (equation 13) becomes

$$\min_{g \in \mathbb{G}} \|f_S - g \circ f_T\|_{\mathcal{D}, \mathbf{H}_S}. \quad (18)$$

Since the affine transformation g acts on the target reference model, it can not be directly viewed as a surrogate transfer loss. However, interesting interpretations can be found in this direction, depending on the output dimension of $f_S : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and $f_T : \mathbb{R}^n \rightarrow \mathbb{R}^d$.

That is, when the direction of adversarial transfer is from $f_S \rightarrow f_T$, the indicating relation between it and knowledge transferability would possibly be unidirectional, depending on the dimensions. More discussion is included in the appendix section B due to space limitation.

4. Synthetic Experiments

The synthetic experiment aims to bridge the gap between theory and practice by verifying some of the theoretical insights that may be difficult to compute for large-scale experiments. Specifically, the synthetic experiment aims to verify: first, how influential are the two proposed adversarial transferability metrics α_1, α_2 comparing to the other factors in the generalized adversarial attacks (equation 9); Second, how does the gradient matching distance track the knowledge transfer loss. The dataset ($N = 5000$) is generated by a Gaussian mixture of 10 Gaussians. The ground truth target is set to be the sum of 100 radial basis functions. The dimension of \mathbf{x} is 50, and the dimension of the target is 10. Details of the datasets are defer to appendix section F.

Models Both the source model f_S and target model f_T are one-hidden-layer neural networks with sigmoid activation.

Methods First, sample $D = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$ from the distribution, where \mathbf{x} is 50-dimensional, \mathbf{y} is 10-dimensional. Then we train a target model f_T on D . To derive the source models, we first train a target model on D with width $m = 100$. Denoting the weights of a target model as \mathbf{W} , we randomly sample a direction \mathbf{V} where each entry of \mathbf{V} is sampled from $U(-0.5, 0.5)$, and choose a scale $t \in [0, 1]$. Subsequently, we perturb the model weights of the clean source model as $\mathbf{W}' := \mathbf{W} + t\mathbf{V}$, and define the source model f_S to be a one-hidden-layer neural network with weights \mathbf{W}' . Then, we compute each of the quantities we care about, including α_1, α_2 from both $f_S \rightarrow f_T$ and $f_T \rightarrow f_S$, the gradient matching distance (equation 7), and the actual knowledge transfer distance (equation 17). We use the standard ℓ_2 loss as the adversarial loss function.

Results We present two sets of experiment in Figure 2. The indication relations between adversarial transferability and knowledge transferability can be observed. Moreover: 1. the metrics α_1, α_2 are more meaningful if using the regular attacks $\delta_{f_*}^{(1)}$; 2. the gradient matching distance tracks the actual knowledge transferability loss; 3. the directions of $f_T \rightarrow f_S$ and $f_S \rightarrow f_T$ are similar.

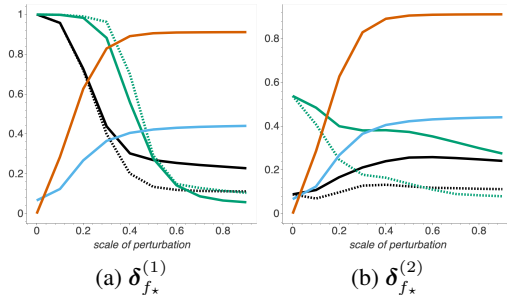


Figure 2. As defined in equation 9, (a) corresponds to the regular adversarial attacks, while (b) the secondary adversarial attack. That is, (b) represents the other information in the adversarial transferring process compared with the first. The x-axis shows the scale of perturbation $t \in [0, 1]$ that controls how much the source model deviates from the target model. There are in total 6 quantities reported. Specifically, $\alpha_1^{f_T \rightarrow f_S}$ is **black solid**; $\alpha_1^{f_S \rightarrow f_T}$ is **black dotted**; $\alpha_2^{f_T \rightarrow f_S}$ is **green solid**; $\alpha_2^{f_S \rightarrow f_T}$ is **green dotted**; the gradient matching loss is **red solid**; and the knowledge transferability distance is **blue solid**.

5. Experimental Evaluation

We present the real-data experiments based on both image and natural language datasets in this section, and discuss the potential applications.

Adversarial Transferability Indicating Knowledge Transferability. In this experiment, we show how to use adversarial transferability to identify the optimal transfer learning candidates from a pool of models trained on the same source dataset. We first train 5 different architectures (AlexNet, Fully connected network, LeNet, ResNet18, ResNet50) on cifar10 (Krizhevsky et al., 2009). Then we perform transfer learning to STL10 (Coates et al., 2011) to obtain the knowledge transferability of each, measured by accuracy. At the same time, we also train one ResNet18 on STL10 as the target model, which has poor accuracy because of the lack of data. To measure the adversarial transferability, we generate adversarial examples with PGD (Madry et al., 2017) on the target model and use the generated adversarial examples to attack each source model. The adversarial transferability is expressed in the form of α_1 and α_2 . Our results in Table 1 indicate that we can use adversarial transferability to forecast knowledge transferability, where the only major computational overheads are training a naive model on the target domain and generating a few adversarial examples.

In the end, We further evaluate the significance of our results by Pearson score. More details about training and generation of adversarial examples can be found in the appendix G.

Model	Knowledge Trans.	α_1	α_2	$\alpha_1 * \alpha_2$
Fully Connected	28.30	0.346	0.189	0.0258
LeNet	45.65	0.324	0.215	0.0254
AlexNet	55.09	0.337	0.205	0.0268
ResNet18	76.60	0.538	0.244	0.0707
ResNet50	77.92	0.614	0.234	0.0899

Table 1. Knowledge transferability (Knowledge Trans.) among different model architectures. Our correlation analysis shows Pearson score of -0.51 between the transfer loss and α_1 . Lower transfer loss corresponds to higher transfer accuracy. More details can be found in fig 4 in the Appendix G

To further validate our idea, we also conduct experiments on the NLP domain. We first finetune 5 different BERT classification models on different data domain (IMDB, Movie Review (MR), Yelp, AG, Fake). We refer the models trained on MR, Yelp, AG and Fake datasets as the source models, and take the model trained on IMDB dataset as the target model. To measure the knowledge transferability, we finetune the source models with new linear layers on the target dataset for one epoch. We report the accuracy of the transferred models on the target test set as the metric to indicate the knowledge transferability. In terms of the adversarial transferability, we generate adversarial examples by the state-of-the-art whitebox attack algorithm T3 (Wang et al., 2020) against the target model and transfer the adversarial examples to source models to evaluate the adversarial transferability. Following our previous experiment, we also calculate α_1 and α_2 . Experimental results are shown in Table 2. We observe that source models with larger adversarial transferability, measured by α_1, α_2 and $\alpha_1 * \alpha_2$, indeed tend to have larger knowledge transferability.

Model	Knowledge Trans.	α_1	α_2	$\alpha_1 * \alpha_2$
MR	89.34	0.743	0.00335	3.00e-3
Yelp	88.81	0.562	0.00135	8.87e-4
AG	87.58	0.295	0.00021	8.56e-5
Fake	84.06	0.028	0.00032	5.58e-6

Table 2. Knowledge transferability (Knowledge Trans.) from the Source Models (MR, Yelp, AG, Fake) to the Target Model (IMDB). Adversarial transferability is measured by using the adversarial examples generated against the Target Model (IMDB) to attack the Source Models and estimate α_1 and α_2 . The correlation analysis shows Pearson Score of 0.27 between the transfer confidence and α_1 . Higher transfer confidence indicates higher knowledge transferability. More details can be found in Figure 6 in Appendix §G.

Knowledge Transferability Indicating Adversarial Transferability. In addition, we are interested in the impact of knowledge transferability on adversarial transferability.

Uncovering the Connections Between Adversarial Transferability and Knowledge Transferability

Similarity	Knowledge Trans.	α_1	α_2	$\alpha_1 * \alpha_2$
0%	45.00	0.310	0.146	0.0169
25%	45.68	0.318	0.305	0.0383
50%	59.09	0.338	0.355	0.0436
75%	71.62	0.337	0.312	0.0402
100%	81.84	0.358	0.357	0.0489

Table 3. Knowledge transferability (Knowledge Trans.) of different source model. Similarity indicates how similar the source distributions are with the target distribution. Our correlation analysis shows Pearson score of -0.06 between the transfer loss and α_1 . Lower transfer loss corresponds to higher knowledge transferability. More details can be found in fig 8 in the Appendix G.

Model	Knowledge Trans.	α_1	α_2	$\alpha_1 * \alpha_2$
MR	89.34	0.584	0.00188	2.32e-3
Yelp	88.81	0.648	0.00120	9.52e-4
AG	87.58	0.293	0.00016	4.35e-6
Fake	84.06	0.150	0.00073	3.55e-5

Table 4. Knowledge transferability (Knowledge Trans.) from the Source Models (MR, Yelp, AG, Fake) to the Target Model (IMDB). Adversarial transferability is measured by using the adversarial examples generated against the Source Models to attack the Target Models and estimate α_1 and α_2 . The correlation analysis shows Pearson Score of 0.27 between the transfer confidence and α_1 . Higher transfer confidence indicates higher knowledge transferability. More details can be found in Figure 9 in Appendix §G.

As predicted by our theory, the more knowledge transferable a source model is to the target domain, the more adversarial transferable it is.

We split cifar10 into 5 different subsets containing different percentages of animals and vehicles. We train a resNet18 on each of them as source models, which are later fine-tuned to obtain the knowledge transferability measured by accuracy. Then we train another resNet18 on a subset of stl10 that only contains vehicles. Different from the last experiment, we generate adversarial examples with PGD on each of the source models and transfer them to the target model. Table 3 shows, the source model that transfers knowledge better generates more transferable adversarial examples. This implies we can use this relation to facilitate blackbox attack against a hidden target model, given some knowledge about the source and target domains. More details of training and generation of adversarial examples can be found in the appendix.

We evaluate the impact of knowledge transferability to adversarial transferability in the NLP domain as well. We mostly follow the setting describe in the previous section, where we have four source models and one target model, and the knowledge transferability from source models to the target model is measured by the accuracy of the transferred models on the target test set. The difference lies on the evaluation of the adversarial transferability, where we generate adversarial examples against the source models and evaluate their attack capability on the target model. As shown in

Table 4, we note that when the source data domain is getting closer to the target data domain, the knowledge transferability grows, and the adversarial transferability also increases. More experimental details can be found in Appendix G.

Ablation Studies Following the settings in table 1, we conduct ablation studies (table 5) on two additional attack methods, MI (Tramèr et al., 2017a), PGD-L2 and two additional ϵ with PGD, 2/225, 4/255, we discover that neither the attack method nor ϵ has significant impact on our conclusion.

Model	Knowledge Trans.	α_1	α_2	$\alpha_1 * \alpha_2$
Fully Connected	28.30	0.0985	0.0196	0.00027
LeNet	45.65	0.2106	0.0259	0.00158
AlexNet	55.09	0.1196	0.0206	0.00037
ResNet18	76.60	0.2739	0.0413	0.00405
ResNet50	77.92	0.1952	0.0320	0.00172

$\epsilon = 2/255$. Pearson score is -0.45.

Model	Knowledge Trans.	α_1	α_2	$\alpha_1 * \alpha_2$
Fully Connected	28.30	0.0974	0.0225	0.00029
LeNet	45.65	0.2099	0.0309	0.00192
AlexNet	55.09	0.1283	0.0230	0.00048
ResNet18	76.60	0.2853	0.0481	0.00496
ResNet50	77.92	0.2495	0.0414	0.00337

$\epsilon = 4/255$. Pearson score is -0.49.

Model	Knowledge Trans.	α_1	α_2	$\alpha_1 * \alpha_2$
Fully Connected	28.30	0.1678	0.0379	0.0013
LeNet	45.65	0.0997	0.0503	0.0005
AlexNet	55.09	0.1229	0.0506	0.0009
ResNet18	76.60	0.2731	0.0630	0.0052
ResNet50	77.92	0.3695	0.0550	0.0081

Attack with MI. Pearson score is -0.45.

Model	Knowledge Trans.	α_1	α_2	$\alpha_1 * \alpha_2$
Fully Connected	28.30	0.0809	0.0175	0.00018
LeNet	45.65	0.2430	0.0190	0.00149
AlexNet	55.09	0.1101	0.0188	0.00031
ResNet18	76.60	0.3619	0.0303	0.00464
ResNet50	77.92	0.2506	0.0237	0.00179

ℓ_2 attack with $\epsilon = 1$. Pearson score is -0.40.

Table 5. With varying attack methods and ϵ , adversarial transferability is still correlated with knowledge transferability.

6. Conclusion

We theoretically analyze the relation between adversarial transferability and knowledge transferability. We provide empirical experimental justifications in practical settings. Both our theoretical and empirical results show that adversarial transferability can indicate knowledge transferability and vice versa. We expect our work will inspire future work on further exploring other factors that impact knowledge transferability and adversarial transferability.

Acknowledgments This work is partially supported by NSF IIS 1909577, NSF CCF 1934986, NSF CCF 1910100, NSF CNS 20-46726 CAR, Amazon Research Award, and the Intel RSA 2020.

References

- Achille, A., Lam, M., Tewari, R., Ravichandran, A., Maji, S., Fowlkes, C. C., Soatto, S., and Perona, P. Task2vec: Task embedding for meta-learning. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 6430–6439, 2019.
- Athalye, A., Carlini, N., and Wagner, D. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. In *International Conference on Machine Learning*, pp. 274–283, 2018.
- Coates, A., Ng, A., and Lee, H. An analysis of single-layer networks in unsupervised feature learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pp. 215–223, 2011.
- Demontis, A., Melis, M., Pintor, M., Jagielski, M., Biggio, B., Oprea, A., Nita-Rotaru, C., and Roli, F. Why do adversarial attacks transfer? explaining transferability of evasion and poisoning attacks. In *28th {USENIX} Security Symposium ({USENIX} Security 19)*, pp. 321–338, 2019.
- Dong, D., Wu, H., He, W., Yu, D., and Wang, H. Multi-task learning for multiple language translation. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 1723–1732, 2015.
- Dong, Y., Pang, T., Su, H., and Zhu, J. Evading defenses to transferable adversarial examples by translation-invariant attacks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4312–4321, 2019.
- Goodfellow, I. J., Shlens, J., and Szegedy, C. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014.
- Huh, M., Agrawal, P., and Efros, A. A. What makes imagenet good for transfer learning? *arXiv preprint arXiv:1608.08614*, 2016.
- Ilyas, A., Engstrom, L., Athalye, A., and Lin, J. Black-box adversarial attacks with limited queries and information. In *International Conference on Machine Learning*, pp. 2137–2146, 2018.
- Joon Oh, S., Fritz, M., and Schiele, B. Adversarial image perturbation for privacy protection—a game theory perspective. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1482–1491, 2017.
- Kariyappa, S. and Qureshi, M. K. Improving adversarial robustness of ensembles with diversity training. *arXiv preprint arXiv:1901.09981*, 2019.
- Kendall, A., Gal, Y., and Cipolla, R. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 7482–7491, 2018.
- Krizhevsky, A., Hinton, G., et al. Learning multiple layers of features from tiny images. 2009.
- Liu, Y., Chen, X., Liu, C., and Song, D. Delving into transferable adversarial examples and black-box attacks. *arXiv preprint arXiv:1611.02770*, 2016.
- Long, M., Cao, Y., Wang, J., and Jordan, M. Learning transferable features with deep adaptation networks. In *International Conference on Machine Learning*, pp. 97–105, 2015.
- Ma, X., Li, B., Wang, Y., Erfani, S. M., Wijewickrema, S., Schoenebeck, G., Song, D., Houle, M. E., and Bailey, J. Characterizing adversarial subspaces using local intrinsic dimensionality. *arXiv preprint arXiv:1801.02613*, 2018.
- Madry, A., Makelov, A., Schmidt, L., Tsipras, D., and Vladu, A. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.
- Miyato, T., Maeda, S.-i., Koyama, M., and Ishii, S. Virtual adversarial training: a regularization method for supervised and semi-supervised learning. *IEEE transactions on pattern analysis and machine intelligence*, 41(8):1979–1993, 2018.
- Naseer, M. M., Khan, S. H., Khan, M. H., Khan, F. S., and Porikli, F. Cross-domain transferability of adversarial perturbations. In *Advances in Neural Information Processing Systems*, pp. 12885–12895, 2019.
- Papernot, N., McDaniel, P., and Goodfellow, I. Transferability in machine learning: from phenomena to black-box attacks using adversarial samples. *arXiv preprint arXiv:1605.07277*, 2016.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3): 211–252, 2015.
- Salman, H., Ilyas, A., Engstrom, L., Kapoor, A., and Madry, A. Do adversarially robust imagenet models transfer better? *arXiv preprint arXiv:2007.08489*, 2020.
- Shinya, Y., Simo-Serra, E., and Suzuki, T. Understanding the effects of pre-training for object detectors via eigen-spectrum. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pp. 0–0, 2019.

- Tramèr, F., Kurakin, A., Papernot, N., Goodfellow, I., Boneh, D., and McDaniel, P. Ensemble adversarial training: Attacks and defenses. *arXiv preprint arXiv:1705.07204*, 2017a.
- Tramèr, F., Papernot, N., Goodfellow, I., Boneh, D., and McDaniel, P. The space of transferable adversarial examples. *arXiv preprint arXiv:1704.03453*, 2017b.
- Utrera, F., Kravitz, E., Erichson, N. B., Khanna, R., and Mahoney, M. W. Adversarially-trained deep nets transfer better. *arXiv preprint arXiv:2007.05869*, 2020.
- Wang, A., Singh, A., Michael, J., Hill, F., Levy, O., and Bowman, S. R. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In *International Conference on Learning Representations*, 2019a. URL <https://openreview.net/forum?id=rJ4km2R5t7>.
- Wang, B., Pei, H., Pan, B., Chen, Q., Wang, S., and Li, B. T3: Tree-autoencoder constrained adversarial text generation for targeted attack. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 6134–6150, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.495. URL <https://www.aclweb.org/anthology/2020.emnlp-main.495>.
- Wang, Z., Dai, Z., Póczos, B., and Carbonell, J. Characterizing and avoiding negative transfer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 11293–11302, 2019b.
- Xie, C., Zhang, Z., Zhou, Y., Bai, S., Wang, J., Ren, Z., and Yuille, A. L. Improving transferability of adversarial examples with input diversity. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2730–2739, 2019.
- Xu, R., Li, G., Yang, J., and Lin, L. Larger norm more transferable: An adaptive feature norm approach for unsupervised domain adaptation. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1426–1435, 2019.
- Yosinski, J., Clune, J., Bengio, Y., and Lipson, H. How transferable are features in deep neural networks? In *Advances in neural information processing systems*, pp. 3320–3328, 2014.
- Zamir, A. R., Sax, A., Shen, W., Guibas, L. J., Malik, J., and Savarese, S. Taskonomy: Disentangling task transfer learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3712–3722, 2018.
- Zhang, X., Zhao, J., and LeCun, Y. Character-level convolutional networks for text classification. *arXiv preprint arXiv:1509.01626*, 2015.
- Zhang, Z., Luo, P., Loy, C. C., and Tang, X. Facial landmark detection by deep multi-task learning. In *European conference on computer vision*, pp. 94–108. Springer, 2014.
- Zhou, W., Hou, X., Chen, Y., Tang, M., Huang, X., Gan, X., and Yang, Y. Transferable adversarial perturbations. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 452–467, 2018.

Supplementary Material: Uncovering the Connections Between Adversarial Transferability and Knowledge Transferability

Contents Summary

- Section A: An Example Illustrating the Necessity of both α_1, α_2 in Characterizing the Relation Between Adversarial Transferability and Knowledge Transferability.
- Section B: Detailed discussion about the direction of adversarial transfer from $f_S \rightarrow f_T$ in subsection 3.3.
- Section C: Proofs of the propositions in section 2.
 - C.1: Proof of Proposition 2.1
 - C.2: Proof of Proposition 2.2
- Section D: Proofs of the theorems and propositions in section 3.
 - D.1: Proof of Theorem 3.1
 - D.2: Proof of Proposition 3.1
 - D.3: Proof of Theorem 3.2
 - D.4: Proof of Theorem 3.3
 - D.5: Proof of Theorem B.1
- Section E: Auxiliary lemmas.
- Section F: Details and additional results of the synthetic experiments.
- Section G: Details of model training and adversarial examples generations in the experiments section, and ablation study on controlling the adversarial transferability.

A. An Example Illustrating the Necessity of both α_1, α_2 in Characterizing the Relation Between Adversarial Transferability and Knowledge Transferability

α_1 and α_2 (Definition 1&2) represent complementary aspects of the adversarial transferability: α_1 can be understood as how often the adversarial attack transfers, while α_2 encodes directional information of the output deviation caused by adversarial attacks. Recall that $\alpha_1, \alpha_2 \in [0, 1]$ (higher values indicate better adversarial transferability). As we show in our theoretical results reveal that high α_1 alone is not enough, *i.e.*, both the proposed metrics are necessary to characterize adversarial transferability and the relation between adversarial and knowledge transferabilities.

We provide a one-dimensional example showing that large α_1 only is not enough to indicate high knowledge transferability. Suppose the ground truth target function $f_T(x) = x^2$, and the source function $f_S(x) = \text{sgn}(x) \cdot x^2$ where $\text{sgn}(\cdot)$ denotes the sign function. Let the adversarial loss be the deviation in function output, and the data distribution be the uniform distribution on $[-1, 1]$. As we can see, the direction that makes either f_T or f_S deviates the most is always the same, *i.e.*, in this example even with $\alpha_1 = 1$ achieves its maximum and adversarial attacks always transfer, regardless of the choice of $f_1 \rightarrow f_2$ or $f_2 \rightarrow f_1$. However, there does not exist an affine function g (*i.e.*, fine-tuning) making $g \circ f_S$ close to f_T on $[-1, 1]$. Indeed, one can verify that $\alpha_2 = 0$ in this case (either $f_1 \rightarrow f_2$ or $f_2 \rightarrow f_1$), which contributes to the low knowledge transferability. However, if we move the data distribution to $[0, 2]$, we can have $\alpha_1 = \alpha_2 = 1$ (either $f_1 \rightarrow f_2$ or $f_2 \rightarrow f_1$) indicating high adversarial transferability, and indeed it achieves $f_S = f_T$ showing perfect knowledge transferability.

B. Detailed Discussion About the Direction of Adversarial Transfer From $f_S \rightarrow f_T$ in Subsection 3.3

In this section, we present a detailed discussion, in addition to subsection 3.3, about the connection between function matching distance and knowledge transfer distance when the direction of adversarial transfer is from $f_S \rightarrow f_T$.

Recall that, to complete the story, it remains to connect the function matching distance to knowledge transferability. As the adversarial transfer is symmetric (*i.e.*, either from $f_S \rightarrow f_T$ or $f_T \rightarrow f_S$), we are able to use the placeholders $\star, \diamond \in \{S, T\}$ all the way through. However, as the knowledge transfer is asymmetric (*i.e.*, $f_S \rightarrow y$ to the target ground truth), we need to instantiate the direction of adversarial transfer to further our discussion. We have discussed the direction of adversarial transfer from $f_T \rightarrow f_S$ in the main paper, where we show that the function matching distance of this direction, *i.e.*,

$$\min_{g \in \mathbb{G}} \|f_T - g \circ f_S\|_{\mathcal{D}, \mathbf{H}_T}, \quad (16)$$

can both upper and lower bound the knowledge transfer distance, *i.e.*,

$$\min_{g \in \mathbb{G}} \|y - g \circ f_S\|_{\mathcal{D}, \mathbf{H}_T}. \quad (17)$$

The direction of adversarial transfer from $f_S \rightarrow f_T$ corresponds to $(\star, \diamond) = (S, T)$. Accordingly, the function matching distance (equation 13) becomes

$$\min_{g \in \mathbb{G}} \|f_S - g \circ f_T\|_{\mathcal{D}, \mathbf{H}_S}. \quad (18)$$

Since the affine transformation g acts on the target reference model, it can not be directly viewed as a surrogate transfer loss. However, interesting interpretations can be found in this direction, depending on the output dimension of $f_S : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and $f_T : \mathbb{R}^n \rightarrow \mathbb{R}^d$.

In this subsection in the appendix we provide detailed discussion on the connection between the function matching distance of the direction of adversarial transfer from $f_S \rightarrow f_T$ (equation 18) and the knowledge transfer distance (equation 17). We build this connection by providing the relationships between the two directions of function matching distance, *i.e.*, equation 16 and equation 18. That is being said, since we know equation 17 and equation 16 are tied together, we only need to provide relationships between equation 16 and equation 18 to show the connection between equation 18 and equation 17.

Suppose $g : \mathbb{R}^d \rightarrow \mathbb{R}^m$ is full rank, and loosely speaking we can derive the following intuitions.

- If $d < m$, then g is injective and there exists $g^{-1} : \mathbb{R}^m \rightarrow \mathbb{R}^d$ such that $g^{-1} \circ g$ is the identity function. That is, if g can map f_T to closely track f_S , then reversely g^{-1} can map f_S to f_T , showing equation 18 upper bounds equation 16 in some sense.
- If $d > m$, then g is surjective. By symmetry, equation 16 upper bounds equation 18 in some sense.
- It is when $m = d$ that equation 16 and equation 18 coincide.

Formally, we have the following theorem.

Theorem B.1. Denote $\tilde{g}_{T,S} : \mathbb{R}^m \rightarrow \mathbb{R}^d$ as the optimal solution of equation 16, and $\tilde{g}_{S,T} : \mathbb{R}^d \rightarrow \mathbb{R}^m$ as the optimal solution of equation 18. Suppose the two optimal affine maps $\tilde{g}_{T,S}, \tilde{g}_{S,T}$ are both full-rank. For $\mathbf{v} \in \mathbb{R}^m$, denote the matrix representation of $\tilde{g}_{T,S}$ as $\tilde{g}_{T,S}(\mathbf{v}) = \tilde{\mathbf{W}}_{T,S} \mathbf{v} + \tilde{\mathbf{b}}_{T,S}$. Similarly, for $\mathbf{w} \in \mathbb{R}^d$, denote the matrix representation of $\tilde{g}_{S,T}$ as $\tilde{g}_{S,T}(\mathbf{w}) = \tilde{\mathbf{W}}_{S,T} \mathbf{w} + \tilde{\mathbf{b}}_{S,T}$. We have the following statements.

If $d < m$, then $\tilde{g}_{S,T}$ is injective, and we have:

$$\|f_T - \tilde{g}_{T,S} \circ f_S\|_{\mathcal{D}, \mathbf{H}_T} \leq \sqrt{\|(\tilde{\mathbf{W}}_{S,T}^T \tilde{\mathbf{W}}_{S,T})^{-1}\|_F \cdot \|\mathbf{H}_T\|_F} \cdot \|f_S - \tilde{g}_{S,T} \circ f_T\|_{\mathcal{D}}. \quad (19)$$

If $d > m$, then $\tilde{g}_{T,S}$ is injective, and we have:

$$\|f_S - \tilde{g}_{S,T} \circ f_T\|_{\mathcal{D}, \mathbf{H}_S} \leq \sqrt{\|(\tilde{\mathbf{W}}_{T,S}^T \tilde{\mathbf{W}}_{T,S})^{-1}\|_F \cdot \|\mathbf{H}_S\|_F} \cdot \|f_T - \tilde{g}_{T,S} \circ f_S\|_{\mathcal{D}}. \quad (20)$$

If $d = m$, then both $\tilde{g}_{S,T}$ and $\tilde{g}_{T,S}$ are bijective, and we have both (19) and (20) stand.

That is, when the direction of adversarial transfer is from $f_S \rightarrow f_T$, the indicating relation between the function matching distance if this direction (equation 18) and knowledge transferability would possibly be unidirectional, depending on the dimensions.

C. Proofs in Section 2

In this section, we present proofs for Proposition 2.1 and Proposition 2.2.

C.1. Proof of Proposition 2.1

Proposition C.1 (Proposition 2.1 Restated). *The $\alpha_2^{f_1 \rightarrow f_2}$ can be reformulated as*

$$(\alpha_2^{f_1 \rightarrow f_2})^2 = \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2} [\theta_{f_1 \rightarrow f_1}(\mathbf{x}_1, \mathbf{x}_2) \theta_{f_1 \rightarrow f_2}(\mathbf{x}_1, \mathbf{x}_2)],$$

where $\mathbf{x}_1, \mathbf{x}_2 \stackrel{i.i.d.}{\sim} \mathcal{D}$, and

$$\theta_{f_1 \rightarrow f_1}(\mathbf{x}_1, \mathbf{x}_2) = \langle \widehat{\Delta_{f_1 \rightarrow f_1}}(\mathbf{x}_1), \widehat{\Delta_{f_1 \rightarrow f_1}}(\mathbf{x}_2) \rangle$$

$$\theta_{f_1 \rightarrow f_2}(\mathbf{x}_1, \mathbf{x}_2) = \langle \widehat{\Delta_{f_1 \rightarrow f_2}}(\mathbf{x}_1), \widehat{\Delta_{f_1 \rightarrow f_2}}(\mathbf{x}_2) \rangle$$

Proof. Recall that we want to show

$$\|\mathbb{E}_{\mathbf{x}} [\widehat{\Delta_{f_1 \rightarrow f_1}}(\mathbf{x}) \widehat{\Delta_{f_1 \rightarrow f_2}}(\mathbf{x})^\top]\|_F^2 = (\alpha_2^{f_1 \rightarrow f_2})^2 = \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2} [\theta_{f_1 \rightarrow f_1}(\mathbf{x}_1, \mathbf{x}_2) \theta_{f_1 \rightarrow f_2}(\mathbf{x}_1, \mathbf{x}_2)],$$

and the proof of this proposition is done by applying some trace tricks, as shown below.

$$\begin{aligned} \theta_{f_1 \rightarrow f_1}(\mathbf{x}_1, \mathbf{x}_2) \theta_{f_1 \rightarrow f_2}(\mathbf{x}_1, \mathbf{x}_2) &= \langle \widehat{\Delta_{f_1 \rightarrow f_1}}(\mathbf{x}_1), \widehat{\Delta_{f_1 \rightarrow f_1}}(\mathbf{x}_2) \rangle \cdot \langle \widehat{\Delta_{f_1 \rightarrow f_2}}(\mathbf{x}_1), \widehat{\Delta_{f_1 \rightarrow f_2}}(\mathbf{x}_2) \rangle \\ &= \langle \widehat{\Delta_{f_1 \rightarrow f_1}}(\mathbf{x}_2), \widehat{\Delta_{f_1 \rightarrow f_1}}(\mathbf{x}_1) \rangle \cdot \langle \widehat{\Delta_{f_1 \rightarrow f_2}}(\mathbf{x}_1), \widehat{\Delta_{f_1 \rightarrow f_2}}(\mathbf{x}_2) \rangle \\ &= \widehat{\Delta_{f_1 \rightarrow f_1}}(\mathbf{x}_2)^\top \widehat{\Delta_{f_1 \rightarrow f_1}}(\mathbf{x}_1) \widehat{\Delta_{f_1 \rightarrow f_2}}(\mathbf{x}_1)^\top \widehat{\Delta_{f_1 \rightarrow f_2}}(\mathbf{x}_2) \\ &= \text{tr} \left(\widehat{\Delta_{f_1 \rightarrow f_1}}(\mathbf{x}_2)^\top \widehat{\Delta_{f_1 \rightarrow f_1}}(\mathbf{x}_1) \widehat{\Delta_{f_1 \rightarrow f_2}}(\mathbf{x}_1)^\top \widehat{\Delta_{f_1 \rightarrow f_2}}(\mathbf{x}_2) \right) \\ &= \text{tr} \left(\widehat{\Delta_{f_1 \rightarrow f_1}}(\mathbf{x}_1) \widehat{\Delta_{f_1 \rightarrow f_2}}(\mathbf{x}_1)^\top \widehat{\Delta_{f_1 \rightarrow f_2}}(\mathbf{x}_2) \widehat{\Delta_{f_1 \rightarrow f_1}}(\mathbf{x}_2)^\top \right) \end{aligned} \quad (21)$$

Plugging equation 21 into equation 25, we have

$$\begin{aligned} (\alpha_2^{f_1 \rightarrow f_2})^2 &= \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2} \left[\text{tr} \left(\widehat{\Delta_{f_1 \rightarrow f_1}}(\mathbf{x}_1) \widehat{\Delta_{f_1 \rightarrow f_2}}(\mathbf{x}_1)^\top \widehat{\Delta_{f_1 \rightarrow f_2}}(\mathbf{x}_2) \widehat{\Delta_{f_1 \rightarrow f_1}}(\mathbf{x}_2)^\top \right) \right] \\ &= \text{tr} \left(\mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2} \left[\widehat{\Delta_{f_1 \rightarrow f_1}}(\mathbf{x}_1) \widehat{\Delta_{f_1 \rightarrow f_2}}(\mathbf{x}_1)^\top \widehat{\Delta_{f_1 \rightarrow f_2}}(\mathbf{x}_2) \widehat{\Delta_{f_1 \rightarrow f_1}}(\mathbf{x}_2)^\top \right] \right) \\ &= \text{tr} \left(\mathbb{E}_{\mathbf{x}_1} \left[\widehat{\Delta_{f_1 \rightarrow f_1}}(\mathbf{x}_1) \widehat{\Delta_{f_1 \rightarrow f_2}}(\mathbf{x}_1)^\top \right] \cdot \mathbb{E}_{\mathbf{x}_2} \left[\widehat{\Delta_{f_1 \rightarrow f_2}}(\mathbf{x}_2) \widehat{\Delta_{f_1 \rightarrow f_1}}(\mathbf{x}_2)^\top \right] \right), \end{aligned} \quad (22)$$

where the last equality is because that $\mathbf{x}_1, \mathbf{x}_2$ are *i.i.d.* samples from the same distribution.

Therefore, we can re-write the $\mathbf{x}_1, \mathbf{x}_2$ to be the same $\mathbf{x} \sim \mathcal{D}$ and realize that the two matrices are in fact the same one.

$$\begin{aligned} (22) &= \text{tr} \left(\mathbb{E}_{\mathbf{x}} \left[\widehat{\Delta_{f_1 \rightarrow f_1}}(\mathbf{x}) \widehat{\Delta_{f_1 \rightarrow f_2}}(\mathbf{x})^\top \right] \cdot \mathbb{E}_{\mathbf{x}} \left[\widehat{\Delta_{f_1 \rightarrow f_2}}(\mathbf{x}) \widehat{\Delta_{f_1 \rightarrow f_1}}(\mathbf{x})^\top \right] \right) \\ &= \|\mathbb{E}_{\mathbf{x}} [\widehat{\Delta_{f_1 \rightarrow f_1}}(\mathbf{x}) \widehat{\Delta_{f_1 \rightarrow f_2}}(\mathbf{x})^\top]\|_F^2. \end{aligned}$$

□

C.2. Proof of Proposition 2.2

Proposition C.2 (Proposition 2.2 Restated). *The adversarial transferability metrics $\alpha_1^{f_1 \rightarrow f_2}$, $\alpha_2^{f_1 \rightarrow f_2}$ and $(\alpha_1 * \alpha_2)^{f_1 \rightarrow f_2}$ are in $[0, 1]$.*

Proof. Let us begin with

$$\alpha_1^{f_1 \rightarrow f_2}(\mathbf{x}) = \frac{\ell_{adv}(f_2(\mathbf{x}), f_2(\mathbf{x} + \boldsymbol{\delta}_{f_1, \epsilon}(\mathbf{x})))}{\ell_{adv}(f_2(\mathbf{x}), f_2(\mathbf{x} + \boldsymbol{\delta}_{f_2, \epsilon}(\mathbf{x})))}.$$

Recall that $\ell_{adv}(\cdot) \geq 0$, and the definition of adversarial attack:

$$\boldsymbol{\delta}_{f, \epsilon}(\mathbf{x}) = \arg \max_{\|\boldsymbol{\delta}\| \leq \epsilon} \ell_{adv}(f(\mathbf{x}), f(\mathbf{x} + \boldsymbol{\delta})),$$

and we can see that by definition,

$$0 \leq \ell_{adv}(f_2(\mathbf{x}), f_2(\mathbf{x} + \boldsymbol{\delta}_{f_1, \epsilon}(\mathbf{x}))) \leq \ell_{adv}(f_2(\mathbf{x}), f_2(\mathbf{x} + \boldsymbol{\delta}_{f_2, \epsilon}(\mathbf{x}))).$$

Therefore,

$$0 \leq \frac{\ell_{adv}(f_2(\mathbf{x}), f_2(\mathbf{x} + \boldsymbol{\delta}_{f_1, \epsilon}(\mathbf{x})))}{\ell_{adv}(f_2(\mathbf{x}), f_2(\mathbf{x} + \boldsymbol{\delta}_{f_2, \epsilon}(\mathbf{x})))} \leq 1,$$

where we define $0/0 = 0$ if necessary.

Hence, $\alpha_1^{f_1 \rightarrow f_2} = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[\alpha_1^{f_1 \rightarrow f_2}(\mathbf{x})]$ is also in $[0, 1]$.

Next, we use Proposition 2.1 to prove the same property for $\alpha_2^{f_1 \rightarrow f_2}$. Note that

$$(\alpha_2^{f_1 \rightarrow f_2})^2 = \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2} \left[\langle \widehat{\Delta_{f_1 \rightarrow f_1}}(\mathbf{x}_1), \widehat{\Delta_{f_1 \rightarrow f_1}}(\mathbf{x}_2) \rangle \cdot \langle \widehat{\Delta_{f_1 \rightarrow f_2}}(\mathbf{x}_1), \widehat{\Delta_{f_1 \rightarrow f_2}}(\mathbf{x}_2) \rangle \right] \quad (23)$$

is the expectation of the product of two inner products, where each inner product is of two unit-length vector. That is being said, $\langle \widehat{\Delta_{f_1 \rightarrow f_1}}(\mathbf{x}_1), \widehat{\Delta_{f_1 \rightarrow f_1}}(\mathbf{x}_2) \rangle \in [-1, 1]$ and $\langle \widehat{\Delta_{f_1 \rightarrow f_2}}(\mathbf{x}_1), \widehat{\Delta_{f_1 \rightarrow f_2}}(\mathbf{x}_2) \rangle \in [-1, 1]$. Therefore, we know that

$$\mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2} \left[\langle \widehat{\Delta_{f_1 \rightarrow f_1}}(\mathbf{x}_1), \widehat{\Delta_{f_1 \rightarrow f_1}}(\mathbf{x}_2) \rangle \cdot \langle \widehat{\Delta_{f_1 \rightarrow f_2}}(\mathbf{x}_1), \widehat{\Delta_{f_1 \rightarrow f_2}}(\mathbf{x}_2) \rangle \right] \in [-1, 1].$$

In addition, we know from equation 23 that it is non-negative, and hence

$$(\alpha_2^{f_1 \rightarrow f_2})^2 \in [0, 1].$$

As $\alpha_2^{f_1 \rightarrow f_2}$ itself is also non-negative by definition, we can see that $\alpha_2^{f_1 \rightarrow f_2} \in [0, 1]$.

Finally, we move to prove $(\alpha_1 * \alpha_2)^{f_1 \rightarrow f_2} \in [0, 1]$. Recall that

$$(\alpha_1 * \alpha_2)^{f_1 \rightarrow f_2} = \left\| \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\alpha_1^{f_1 \rightarrow f_2}(\mathbf{x}) \widehat{\Delta_{f_1 \rightarrow f_1}}(\mathbf{x}) \widehat{\Delta_{f_1 \rightarrow f_2}}(\mathbf{x})^\top] \right\|_F.$$

If we see $\alpha_1^{f_1 \rightarrow f_2}(\mathbf{x}) \widehat{\Delta_{f_1 \rightarrow f_1}}(\mathbf{x})$ as a whole, we can show exactly the same as the Proposition 2.1 that

$$((\alpha_1 * \alpha_2)^{f_1 \rightarrow f_2})^2 = \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2} [\theta_{f_1 \rightarrow f_1}(\mathbf{x}_1, \mathbf{x}_2) \theta_{f_1 \rightarrow f_2}(\mathbf{x}_1, \mathbf{x}_2)], \quad (24)$$

where

$$\begin{aligned} \theta_{f_1 \rightarrow f_1}(\mathbf{x}_1, \mathbf{x}_2) &= \langle \alpha_1^{f_1 \rightarrow f_2}(\mathbf{x}_1) \widehat{\Delta_{f_1 \rightarrow f_1}}(\mathbf{x}_1), \alpha_1^{f_1 \rightarrow f_2}(\mathbf{x}_2) \widehat{\Delta_{f_1 \rightarrow f_1}}(\mathbf{x}_2) \rangle \\ \theta_{f_1 \rightarrow f_2}(\mathbf{x}_1, \mathbf{x}_2) &= \langle \widehat{\Delta_{f_1 \rightarrow f_2}}(\mathbf{x}_1), \widehat{\Delta_{f_1 \rightarrow f_2}}(\mathbf{x}_2) \rangle. \end{aligned}$$

Similarly, as $\alpha_1^{f_1 \rightarrow f_2}(\mathbf{x}) \in [0, 1]$, we can see that $\theta_{f_1 \rightarrow f_1}(\mathbf{x}_1, \mathbf{x}_2) \theta_{f_1 \rightarrow f_2}(\mathbf{x}_1, \mathbf{x}_2) \in [-1, 1]$, and hence

$$\mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2} [\theta_{f_1 \rightarrow f_1}(\mathbf{x}_1, \mathbf{x}_2) \theta_{f_1 \rightarrow f_2}(\mathbf{x}_1, \mathbf{x}_2)] \in [-1, 1].$$

Noting that equation 24 is non-negative, we conclude that

$$((\alpha_1 * \alpha_2)^{f_1 \rightarrow f_2})^2 \in [0, 1].$$

Since $(\alpha_1 * \alpha_2)^{f_1 \rightarrow f_2}$ itself is non-negative as well, we can see that $(\alpha_1 * \alpha_2)^{f_1 \rightarrow f_2} \in [0, 1]$.

Therefore, the three adversarial transferability metrics are all within $[0, 1]$. \square

D. Proofs in Section 3

In this section, we prove the two theorems and the two propositions presented in section 3, which are our main theories.

D.1. Proof of Theorem 3.1

We introduce two lemmas before proving Theorem 3.1.

Lemma D.1. *The square of the gradient matching distance is*

$$\min_{g \in \mathbb{G}} \|\nabla f_{\star}^{\top} - \nabla(g \circ f_{\diamond})^{\top}\|_{\mathcal{D}, \mathbf{H}_{\star}}^2 = \|\nabla f_{\star}^{\top}\|_{\mathcal{D}, \mathbf{H}_{\star}}^2 - \langle \mathbf{P}^{\top} \mathbf{H}_{\star} \mathbf{P}, \mathbf{J}^{\dagger} \rangle,$$

where $g \in \mathbb{G}$ are affine transformations, and

$$\mathbf{P} = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\nabla f_{\star}(\mathbf{x})^{\top} \nabla f_{\diamond}(\mathbf{x})], \quad \mathbf{J} = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\nabla f_{\diamond}(\mathbf{x})^{\top} \nabla f_{\diamond}(\mathbf{x})].$$

Proof.

$$\begin{aligned} \min_{g \in \mathbb{G}} \|\nabla f_{\star}^{\top} - \nabla(g \circ f_{\diamond})^{\top}\|_{\mathcal{D}, \mathbf{H}_{\star}}^2 &= \min_{\mathbf{W}} \|\nabla f_{\star}^{\top} - \mathbf{W} \nabla f_{\diamond}^{\top}\|_{\mathcal{D}, \mathbf{H}_{\star}}^2 \\ &= \min_{\mathbf{W}} \mathbb{E}_{\mathbf{x} \in \mathcal{D}} \|\nabla f_{\star}(\mathbf{x})^{\top} - \mathbf{W} \nabla f_{\diamond}(\mathbf{x})^{\top}\|_{\mathbf{H}_{\star}}^2, \end{aligned} \quad (25)$$

where \mathbf{W} is a matrix.

We can see that (25) is a convex program, where the optimal solution exists in a closed-form form, as shown in the following. Denote $l(\mathbf{W}) = \|\nabla f_{\star}(\mathbf{x})^{\top} - \mathbf{W} \nabla f_{\diamond}(\mathbf{x})^{\top}\|_{\mathbf{H}_{\star}}^2$, we have

$$\begin{aligned} l(\mathbf{W}) &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\|\nabla f_{\star}(\mathbf{x})^{\top}\|_{\mathbf{H}_{\star}}^2 + \|\mathbf{W} \nabla f_{\diamond}(\mathbf{x})^{\top}\|_{\mathbf{H}_{\star}}^2 - 2\langle \nabla f_{\star}(\mathbf{x})^{\top}, \mathbf{W} \nabla f_{\diamond}(\mathbf{x})^{\top} \rangle_{\mathbf{H}_{\star}}] \\ &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\|\nabla f_{\star}(\mathbf{x})^{\top}\|_{\mathbf{H}_{\star}}^2 + \text{tr}(\nabla f_{\diamond}(\mathbf{x}) \mathbf{W}^{\top} \mathbf{H}_{\star} \mathbf{W} \nabla f_{\diamond}(\mathbf{x})^{\top}) - 2 \text{tr}(\nabla f_{\star}(\mathbf{x}) \mathbf{H}_{\star} \mathbf{W} \nabla f_{\diamond}(\mathbf{x})^{\top})] \\ &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\|\nabla f_{\star}(\mathbf{x})^{\top}\|_{\mathbf{H}_{\star}}^2 + \text{tr}(\mathbf{H}_{\star} \mathbf{W} \nabla f_{\diamond}(\mathbf{x})^{\top} \nabla f_{\diamond}(\mathbf{x}) \mathbf{W}^{\top}) - 2 \text{tr}(\mathbf{H}_{\star} \mathbf{W} \nabla f_{\diamond}(\mathbf{x})^{\top} \nabla f_{\star}(\mathbf{x}))]. \end{aligned}$$

Taking the derivative of $l(\cdot)$ w.r.t. \mathbf{W} , we have

$$\begin{aligned} \frac{\partial l}{\partial \mathbf{W}} &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [2\mathbf{H}_{\star} (\mathbf{W} \nabla f_{\diamond}(\mathbf{x})^{\top} - \nabla f_{\star}(\mathbf{x})^{\top}) \nabla f_{\diamond}(\mathbf{x})] \\ &= 2\mathbf{H}_{\star} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\mathbf{W} \nabla f_{\diamond}(\mathbf{x})^{\top} \nabla f_{\diamond}(\mathbf{x}) - \nabla f_{\star}(\mathbf{x})^{\top} \nabla f_{\diamond}(\mathbf{x})] \\ &= 2\mathbf{H}_{\star} (\mathbf{W} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\nabla f_{\diamond}(\mathbf{x})^{\top} \nabla f_{\diamond}(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\nabla f_{\star}(\mathbf{x})^{\top} \nabla f_{\diamond}(\mathbf{x})]). \end{aligned} \quad (26)$$

Since $l(\cdot)$ is convex, if there exists a $\tilde{\mathbf{W}}$ such that $\frac{\partial l}{\partial \tilde{\mathbf{W}}}|_{\mathbf{W}=\tilde{\mathbf{W}}} = \mathbf{0}$ then we know that $\tilde{\mathbf{W}}$ is an optimal solution. Luckily, we can find such solution easily by using pseudo inverse, *i.e.*,

$$\begin{aligned} \tilde{\mathbf{W}} &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\nabla f_{\star}(\mathbf{x})^{\top} \nabla f_{\diamond}(\mathbf{x})] (\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\nabla f_{\diamond}(\mathbf{x})^{\top} \nabla f_{\diamond}(\mathbf{x})])^{\dagger} \\ &= \mathbf{P} \mathbf{J}^{\dagger}, \end{aligned} \quad (27)$$

where we denote $\mathbf{P} = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\nabla f_{\star}(\mathbf{x})^{\top} \nabla f_{\diamond}(\mathbf{x})]$ and $\mathbf{J} = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\nabla f_{\diamond}(\mathbf{x})^{\top} \nabla f_{\diamond}(\mathbf{x})]$.

We can verify that such $\tilde{\mathbf{W}}$ indeed make the partial derivative (equation 26) zero. In equation 26, we have

$$\tilde{\mathbf{W}} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\nabla f_{\diamond}(\mathbf{x})^{\top} \nabla f_{\diamond}(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\nabla f_{\star}(\mathbf{x})^{\top} \nabla f_{\diamond}(\mathbf{x})] = \mathbf{P} \mathbf{J}^{\dagger} \mathbf{J} - \mathbf{P}. \quad (28)$$

To continue, we can see from Lemma E.2 that $\ker(\mathbf{J}) \subseteq \ker(\mathbf{P})$ which means $\text{rowsp}(\mathbf{P}) \subseteq \text{rowsp}(\mathbf{J})$, where $\ker(\cdot)$ denotes the kernel of a matrix, and $\text{rowsp}(\cdot)$ denotes the row space of a matrix. Therefore, by definition of the pseudo-inverse, we can see that $\mathbf{P} \mathbf{J}^{\dagger} \mathbf{J} = \mathbf{P}$, *i.e.*, (28) = $\mathbf{0}$, and hence $\tilde{\mathbf{W}}$ is indeed the optimal solution.

Plugging (27) into (25), we have the optimal value as

$$\begin{aligned}
 (25) &= l(\tilde{\mathbf{W}}) \\
 &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[\|\nabla f_*(\mathbf{x})^\top\|_{\mathbf{H}_*}^2 + \text{tr} \left(\mathbf{H}_* \tilde{\mathbf{W}} \nabla f_\diamond(\mathbf{x})^\top \nabla f_\diamond(\mathbf{x}) \tilde{\mathbf{W}}^\top \right) - 2 \text{tr} \left(\mathbf{H}_* \tilde{\mathbf{W}} \nabla f_\diamond(\mathbf{x})^\top \nabla f_*(\mathbf{x}) \right) \right] \\
 &= \|\nabla f_*^\top\|_{\mathcal{D}, \mathbf{H}_*}^2 + \text{tr} \left(\mathbf{H}_* \tilde{\mathbf{W}} \mathbf{J} \tilde{\mathbf{W}}^\top - 2 \mathbf{H}_* \tilde{\mathbf{W}} \mathbf{P}^\top \right) \\
 &= \|\nabla f_*^\top\|_{\mathcal{D}, \mathbf{H}_*}^2 + \text{tr} \left(\mathbf{H}_* \mathbf{P} \mathbf{J}^\dagger \mathbf{J} \mathbf{J}^\dagger \mathbf{P}^\top - 2 \mathbf{H}_* \mathbf{P} \mathbf{J}^\dagger \mathbf{P}^\top \right) \\
 &= \|\nabla f_*^\top\|_{\mathcal{D}, \mathbf{H}_*}^2 - \text{tr} \left(\mathbf{H}_* \mathbf{P} \mathbf{J}^\dagger \mathbf{P}^\top \right) \\
 &= \|\nabla f_*^\top\|_{\mathcal{D}, \mathbf{H}_*}^2 - \langle \mathbf{P}^\top \mathbf{H}_* \mathbf{P}, \mathbf{J}^\dagger \rangle.
 \end{aligned}$$

□

Next, we present another lemma to analyze the term $\mathbf{P}^\top \mathbf{H}_* \mathbf{P}$.

Lemma D.2. *In this lemma, we break down the matrix representation of $\mathbf{P}^\top \mathbf{H}_* \mathbf{P}$ into pieces relating to the output deviation caused by the generalized adversarial attacks (defined in equation 10)*

$$\mathbf{P}^\top \mathbf{H}_* \mathbf{P} = \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \stackrel{i.i.d.}{\sim} \mathcal{D}} \sum_{i,j=1}^n \left(\Delta_{f_* \rightarrow f_*}^{(i)}(\mathbf{x}_1)^\top \mathbf{H}_* \Delta_{f_* \rightarrow f_*}^{(j)}(\mathbf{x}_2) \right) \cdot \left(\Delta_{f_* \rightarrow f_\diamond}^{(i)}(\mathbf{x}_1) \Delta_{f_* \rightarrow f_\diamond}^{(j)}(\mathbf{x}_2)^\top \right).$$

Proof. Denote a symmetric decomposition of the positive semi-definitive matrix \mathbf{H}_* as

$$\mathbf{H}_* = \mathbf{T}^\top \mathbf{T},$$

where \mathbf{T} is of the same dimension of \mathbf{H}_* . We note that the choice of decomposition does not matter.

Then, plugging in the definition of \mathbf{P} , we can see that

$$\begin{aligned}
 \mathbf{P}^\top \mathbf{H}_* \mathbf{P} &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[\nabla f_\diamond(\mathbf{x})^\top \nabla f_*(\mathbf{x}) \right] \cdot \mathbf{T}^\top \mathbf{T} \cdot \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[\nabla f_*(\mathbf{x})^\top \nabla f_\diamond(\mathbf{x}) \right] \\
 &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[\nabla f_\diamond(\mathbf{x})^\top \nabla f_*(\mathbf{x}) \mathbf{T}^\top \right] \cdot \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[\mathbf{T} \nabla f_*(\mathbf{x})^\top \nabla f_\diamond(\mathbf{x}) \right].
 \end{aligned} \tag{29}$$

A key observation to connect the above equation to the adversarial attack (equation 5) is that,

$$\begin{aligned}
 \delta_{f_*, \epsilon}(\mathbf{x}) &= \arg \max_{\|\delta\|_2 \leq \epsilon} \|\nabla f_*(\mathbf{x})^\top \delta\|_{\mathbf{H}_*} \\
 &= \arg \max_{\|\delta\|_2 \leq \epsilon} \|\nabla f_*(\mathbf{x})^\top \delta\|_{\mathbf{H}_*}^2 \\
 &= \arg \max_{\|\delta\|_2 \leq \epsilon} \delta^\top \nabla f_*(\mathbf{x}) \mathbf{H}_* \nabla f_*(\mathbf{x})^\top \delta \\
 &= \arg \max_{\|\delta\|_2 \leq \epsilon} \|\mathbf{T} \nabla f_*(\mathbf{x})^\top \delta\|_2^2.
 \end{aligned}$$

That is being said, the adversarial attack is the right singular vector corresponding to the largest singular value (in absolute value) of $\mathbf{T} \nabla f_*(\mathbf{x})^\top$.

Similarly, we can see the singular values $\sigma_{f_*, \mathbf{H}_*}(\mathbf{x}) \in \mathbb{R}^n$, defined as the descending (in absolute value) singular values of the Jacobian $\nabla f_*(\mathbf{x})^\top \in \mathbb{R}^{\times n}$ in the \mathbf{H}_* inner product space (equation 8), are the singular values of $\mathbf{T} \nabla f_*(\mathbf{x})^\top$.

With this perspective, if we write down the singular value decomposition of $\mathbf{T} \nabla f_*(\mathbf{x})^\top$, i.e.,

$$\mathbf{T} \nabla f_*(\mathbf{x})^\top = \mathbf{U}_*(\mathbf{x}) \Sigma_*(\mathbf{x}) \mathbf{V}_*^\top(\mathbf{x}),$$

we can observe that:

1. $\Sigma_*(\mathbf{x})$ is diagonalized singular values $\sigma_{f_*, \mathbf{H}_*}(\mathbf{x})$;
2. The i^{th} column of $\mathbf{V}_*(\mathbf{x})$ is the i^{th} generalized attack $\delta_{f_*}^{(i)}(\mathbf{x})$ (defined in equation 9);

3. The i^{th} column of $\mathbf{U}_*(\mathbf{x})\Sigma(\mathbf{x})$ is $\mathbf{T}\Delta_{f_* \rightarrow f_*}^{(i)}(\mathbf{x})$ where $\Delta_{f_* \rightarrow f_*}^{(i)}(\mathbf{x})$ is the output deviation (defined in equation 10);
4. The i^{th} column of $\nabla f_\diamond(\mathbf{x})^\top \mathbf{V}_*(\mathbf{x})$ is the output deviation $\Delta_{f_* \rightarrow f_\diamond}^{(i)}(\mathbf{x})$ (defined in equation 10).

With the four key observations, we can break down the Jacobian matrices as

$$\begin{aligned} \nabla f_\diamond(\mathbf{x})^\top \nabla f_*(\mathbf{x}) \mathbf{T}^\top &= \left(\Delta_{f_* \rightarrow f_\diamond}^{(1)}(\mathbf{x}) \cdots \Delta_{f_* \rightarrow f_\diamond}^{(n)}(\mathbf{x}) \right) \begin{pmatrix} \Delta_{f_* \rightarrow f_*}^{(1)}(\mathbf{x})^\top \mathbf{T}^\top \\ \vdots \\ \Delta_{f_* \rightarrow f_*}^{(n)}(\mathbf{x})^\top \mathbf{T}^\top \end{pmatrix} \\ &= \sum_{i=1}^n \Delta_{f_* \rightarrow f_\diamond}^{(i)}(\mathbf{x}) \Delta_{f_* \rightarrow f_*}^{(i)}(\mathbf{x})^\top \mathbf{T}^\top. \end{aligned}$$

Therefore, plugging it into the equation 29, we have

$$\begin{aligned} (29) &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[\sum_{i=1}^n \Delta_{f_* \rightarrow f_\diamond}^{(i)}(\mathbf{x}) \Delta_{f_* \rightarrow f_*}^{(i)}(\mathbf{x})^\top \mathbf{T}^\top \right] \cdot \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[\sum_{i=1}^n \mathbf{T} \Delta_{f_* \rightarrow f_*}^{(i)}(\mathbf{x}) \Delta_{f_* \rightarrow f_\diamond}^{(i)}(\mathbf{x})^\top \right] \\ &= \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}} \left[\sum_{i=1}^n \left(\Delta_{f_* \rightarrow f_\diamond}^{(i)}(\mathbf{x}_1) \Delta_{f_* \rightarrow f_*}^{(i)}(\mathbf{x}_1)^\top \mathbf{T}^\top \right) \sum_{j=1}^n \left(\mathbf{T} \Delta_{f_* \rightarrow f_*}^{(j)}(\mathbf{x}_2) \Delta_{f_* \rightarrow f_\diamond}^{(j)}(\mathbf{x}_2)^\top \right) \right] \\ &= \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}} \sum_{i,j=1}^n \left(\Delta_{f_* \rightarrow f_\diamond}^{(i)}(\mathbf{x}_1) \Delta_{f_* \rightarrow f_*}^{(i)}(\mathbf{x}_1)^\top \mathbf{H}_* \Delta_{f_* \rightarrow f_*}^{(j)}(\mathbf{x}_2) \Delta_{f_* \rightarrow f_\diamond}^{(j)}(\mathbf{x}_2)^\top \right) \\ &= \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}} \sum_{i,j=1}^n \left(\Delta_{f_* \rightarrow f_*}^{(i)}(\mathbf{x}_1)^\top \mathbf{H}_* \Delta_{f_* \rightarrow f_*}^{(j)}(\mathbf{x}_2) \right) \cdot \left(\Delta_{f_* \rightarrow f_\diamond}^{(i)}(\mathbf{x}_1) \Delta_{f_* \rightarrow f_\diamond}^{(j)}(\mathbf{x}_2)^\top \right), \end{aligned}$$

where the last equality is due to that $\Delta_{f_* \rightarrow f_*}^{(i)}(\mathbf{x}_1)^\top \mathbf{H}_* \Delta_{f_* \rightarrow f_*}^{(j)}(\mathbf{x}_2)$ is a scalar value. □

Equipped with Lemma D.1 and Lemma D.2, we are able to prove the Theorem 3.1.

Theorem D.1 (Theorem 3.1 Restated). *Given the target and source models f_* , f_\diamond , where $(*, \diamond) \in \{(S, T), (T, S)\}$, the gradient matching distance (equation 7) can be written as*

$$\min_{g \in \mathbb{G}} \|\nabla f_*^\top - \nabla(g \circ f_\diamond)^\top\|_{\mathcal{D}, \mathbf{H}_*} = \left(1 - \frac{\mathbb{E}[\mathbf{v}^{*,\diamond}(\mathbf{x}_1)^\top \mathbf{A}_2^{*,\diamond}(\mathbf{x}_1, \mathbf{x}_2) \mathbf{v}^{*,\diamond}(\mathbf{x}_2)]}{\|\nabla f_*^\top\|_{\mathcal{D}, \mathbf{H}_*}^2 \cdot \|\mathbf{J}^\dagger\|_{\mathbf{H}_\diamond}^{-1}} \right)^{\frac{1}{2}} \|\nabla f_*^\top\|_{\mathcal{D}, \mathbf{H}_*},$$

where the expectation is taken over $\mathbf{x}_1, \mathbf{x}_2 \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}$, and

$$\begin{aligned} \mathbf{v}^{*,\diamond}(\mathbf{x}) &= \sigma_{f_\diamond, \mathbf{H}_\diamond}^{(1)}(\mathbf{x}) \sigma_{f_*, \mathbf{H}_*}(\mathbf{x}) \odot \mathbf{A}_1^{*,\diamond}(\mathbf{x}) \\ \mathbf{J} &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\nabla f_\diamond(\mathbf{x})^\top \nabla f_\diamond(\mathbf{x})]. \end{aligned}$$

Moreover, $\mathbf{A}_2^{*,\diamond}(\mathbf{x}_1, \mathbf{x}_2)$ is a matrix, and its element in the i^{th} row and j^{th} column is

$$\mathbf{A}_2^{*,\diamond}(\mathbf{x}_1, \mathbf{x}_2)^{(i,j)} = \langle \widehat{\Delta_{f_* \rightarrow f_*}^{(i)}}(\mathbf{x}_1) |_{\mathbf{H}_*}, \widehat{\Delta_{f_* \rightarrow f_*}^{(j)}}(\mathbf{x}_2) |_{\mathbf{H}_*} \rangle \cdot \langle \widehat{\Delta_{f_* \rightarrow f_\diamond}^{(i)}}(\mathbf{x}_1) |_{\mathbf{H}_\diamond}, \widehat{\Delta_{f_* \rightarrow f_\diamond}^{(j)}}(\mathbf{x}_2) |_{\mathbf{H}_\diamond} \rangle \widehat{\mathbf{J}^\dagger} |_{\mathbf{H}_\diamond}.$$

Proof. Combining the result from Lemma D.1 and Lemma D.2, and applying the linearity of the inner product, we have

$$\begin{aligned}
 & \min_{g \in \mathbb{G}} \|\nabla f_{\star}^{\top} - \nabla(g \circ f_{\diamond})^{\top}\|_{\mathcal{D}, \mathbf{H}_{\star}}^2 \\
 &= \|\nabla f_{\star}^{\top}\|_{\mathcal{D}, \mathbf{H}_{\star}}^2 - \left\langle \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \sim \text{i.i.d.} \mathcal{D}} \sum_{i,j=1}^n \left(\Delta_{f_{\star} \rightarrow f_{\star}}^{(i)}(\mathbf{x}_1)^{\top} \mathbf{H}_{\star} \Delta_{f_{\star} \rightarrow f_{\star}}^{(j)}(\mathbf{x}_2) \right) \cdot \left(\Delta_{f_{\star} \rightarrow f_{\diamond}}^{(i)}(\mathbf{x}_1) \Delta_{f_{\star} \rightarrow f_{\diamond}}^{(j)}(\mathbf{x}_2)^{\top} \right), \mathbf{J}^{\dagger} \right\rangle \\
 &= \|\nabla f_{\star}^{\top}\|_{\mathcal{D}, \mathbf{H}_{\star}}^2 - \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \sim \text{i.i.d.} \mathcal{D}} \sum_{i,j=1}^n \left(\Delta_{f_{\star} \rightarrow f_{\star}}^{(i)}(\mathbf{x}_1)^{\top} \mathbf{H}_{\star} \Delta_{f_{\star} \rightarrow f_{\star}}^{(j)}(\mathbf{x}_2) \right) \cdot \left\langle \Delta_{f_{\star} \rightarrow f_{\diamond}}^{(i)}(\mathbf{x}_1) \Delta_{f_{\star} \rightarrow f_{\diamond}}^{(j)}(\mathbf{x}_2)^{\top}, \mathbf{J}^{\dagger} \right\rangle \\
 &= \|\nabla f_{\star}^{\top}\|_{\mathcal{D}, \mathbf{H}_{\star}}^2 - \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \sim \text{i.i.d.} \mathcal{D}} \sum_{i,j=1}^n \left(\Delta_{f_{\star} \rightarrow f_{\star}}^{(i)}(\mathbf{x}_1)^{\top} \mathbf{H}_{\star} \Delta_{f_{\star} \rightarrow f_{\star}}^{(j)}(\mathbf{x}_2) \right) \cdot \text{tr} \left(\Delta_{f_{\star} \rightarrow f_{\diamond}}^{(i)}(\mathbf{x}_1) \Delta_{f_{\star} \rightarrow f_{\diamond}}^{(j)}(\mathbf{x}_2)^{\top} \mathbf{J}^{\dagger} \right) \\
 &= \|\nabla f_{\star}^{\top}\|_{\mathcal{D}, \mathbf{H}_{\star}}^2 - \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \sim \text{i.i.d.} \mathcal{D}} \sum_{i,j=1}^n \underbrace{\left(\Delta_{f_{\star} \rightarrow f_{\star}}^{(i)}(\mathbf{x}_1)^{\top} \mathbf{H}_{\star} \Delta_{f_{\star} \rightarrow f_{\star}}^{(j)}(\mathbf{x}_2) \right)}_{X_1} \cdot \underbrace{\left(\Delta_{f_{\star} \rightarrow f_{\diamond}}^{(i)}(\mathbf{x}_1)^{\top} \mathbf{J}^{\dagger} \Delta_{f_{\star} \rightarrow f_{\diamond}}^{(j)}(\mathbf{x}_2) \right)}_{X_2}. \quad (30)
 \end{aligned}$$

As the generalized first adversarial transferability \mathbf{A}_1 is about the magnitude of the output deviation (defined in equation 11), and we can separate the \mathbf{A}_1 out from the above equation. Then, what left should be about the directions about the output deviation, which we will put into the matrix \mathbf{A}_2 , *i.e.*, the generalized second adversarial transferability.

Recall that the generalized the first adversarial transferability is a n -dimensional vector $\mathbf{A}_1^{\star, \diamond}(\mathbf{x})$ including the adversarial losses of all of the generalized adversarial attacks, where the i^{th} element in the vector is

$$\mathbf{A}_1^{\star, \diamond}(\mathbf{x})^{(i)} = \frac{\|\Delta_{f_{\star} \rightarrow f_{\diamond}}^{(i)}(\mathbf{x})\|_{\mathbf{H}_{\diamond}}}{\|\nabla f_{\diamond}(\mathbf{x})\|_{\mathbf{H}_{\diamond}}}.$$

Moreover, to connect the magnitude of the output deviation to the generalized singular values (equation 9), we have

$$\|\Delta_{f_{\star} \rightarrow f_{\star}}^{(i)}(\mathbf{x})\|_{\mathbf{H}_{\star}} = \|\nabla f_{\star}(\mathbf{x})^{\top} \delta_{f_{\star}}^{(i)}(\mathbf{x})\|_{\mathbf{H}_{\star}} = \sigma_{f_{\star}, \mathbf{H}_{\star}}^{(i)}(\mathbf{x}),$$

and similarly,

$$\|\nabla f_{\diamond}(\mathbf{x})\|_{\mathbf{H}_{\diamond}} = \|\nabla f_{\diamond}(\mathbf{x}) \delta_{f_{\diamond}}^{(1)}(\mathbf{x})\|_{\mathbf{H}_{\diamond}} = \sigma_{f_{\diamond}, \mathbf{H}_{\diamond}}^{(1)}(\mathbf{x}).$$

Therefore, we can finally rewrite the X_1, X_2 in equation 30 as

$$\begin{aligned}
 X_1 &= \sigma_{f_{\star}, \mathbf{H}_{\star}}^{(i)}(\mathbf{x}_1) \sigma_{f_{\star}, \mathbf{H}_{\star}}^{(j)}(\mathbf{x}_2) \cdot \langle \widehat{\Delta_{f_{\star} \rightarrow f_{\star}}^{(i)}(\mathbf{x}_1)}|_{\mathbf{H}_{\star}}, \widehat{\Delta_{f_{\star} \rightarrow f_{\star}}^{(j)}(\mathbf{x}_2)}|_{\mathbf{H}_{\star}} \rangle \\
 X_2 &= \mathbf{A}_1^{\star, \diamond}(\mathbf{x}_1)^{(i)} \mathbf{A}_1^{\star, \diamond}(\mathbf{x}_2)^{(j)} \cdot \langle \widehat{\Delta_{f_{\star} \rightarrow f_{\diamond}}^{(i)}(\mathbf{x}_1)}|_{\mathbf{H}_{\diamond}}, \widehat{\Delta_{f_{\star} \rightarrow f_{\diamond}}^{(j)}(\mathbf{x}_2)}|_{\mathbf{H}_{\diamond}} \rangle_{\widehat{\mathbf{J}^{\dagger}}|_{\mathbf{H}_{\diamond}}} \cdot \sigma_{f_{\diamond}, \mathbf{H}_{\diamond}}^{(1)}(\mathbf{x}_1) \sigma_{f_{\diamond}, \mathbf{H}_{\diamond}}^{(1)}(\mathbf{x}_2) \|\mathbf{J}^{\dagger}\|_{\mathbf{H}_{\diamond}}.
 \end{aligned}$$

Recall the $(i, j)^{\text{th}}$ entry of the matrix \mathbf{A}_2 is

$$\mathbf{A}_2^{\star, \diamond}(\mathbf{x}_1, \mathbf{x}_2)^{(i,j)} = \langle \widehat{\Delta_{f_{\star} \rightarrow f_{\star}}^{(i)}(\mathbf{x}_1)}|_{\mathbf{H}_{\star}}, \widehat{\Delta_{f_{\star} \rightarrow f_{\star}}^{(j)}(\mathbf{x}_2)}|_{\mathbf{H}_{\star}} \rangle \cdot \langle \widehat{\Delta_{f_{\star} \rightarrow f_{\diamond}}^{(i)}(\mathbf{x}_1)}|_{\mathbf{H}_{\diamond}}, \widehat{\Delta_{f_{\star} \rightarrow f_{\diamond}}^{(j)}(\mathbf{x}_2)}|_{\mathbf{H}_{\diamond}} \rangle_{\widehat{\mathbf{J}^{\dagger}}|_{\mathbf{H}_{\diamond}}}.$$

We can write

$$X_1 X_2 = \sigma_{f_{\diamond}, \mathbf{H}_{\diamond}}^{(1)}(\mathbf{x}_1) \sigma_{f_{\star}, \mathbf{H}_{\star}}^{(i)}(\mathbf{x}_1) \mathbf{A}_1^{\star, \diamond}(\mathbf{x}_1)^{(i)} \cdot \mathbf{A}_2^{\star, \diamond}(\mathbf{x}_1, \mathbf{x}_2)^{(i,j)} \cdot \sigma_{f_{\diamond}, \mathbf{H}_{\diamond}}^{(1)}(\mathbf{x}_2) \sigma_{f_{\star}, \mathbf{H}_{\star}}^{(j)}(\mathbf{x}_2) \mathbf{A}_1^{\star, \diamond}(\mathbf{x}_2)^{(j)} \|\mathbf{J}^{\dagger}\|_{\mathbf{H}_{\diamond}}.$$

Plugging the above into equation 30, and rearranging the double summation, we have

$$\begin{aligned}
 (30) &= \|\nabla f_{\star}^{\top}\|_{\mathcal{D}, \mathbf{H}_{\star}}^2 \\
 &- \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \sim \text{i.i.d.} \mathcal{D}} \left[(\sigma_{f_{\diamond}, \mathbf{H}_{\diamond}}^{(1)}(\mathbf{x}_1) \sigma_{f_{\star}, \mathbf{H}_{\star}}^{(i)}(\mathbf{x}_1) \odot \mathbf{A}_1^{\star, \diamond}(\mathbf{x}_1))^{\top} \mathbf{A}_2^{\star, \diamond}(\mathbf{x}_1, \mathbf{x}_2) (\sigma_{f_{\diamond}, \mathbf{H}_{\diamond}}^{(1)}(\mathbf{x}_2) \sigma_{f_{\star}, \mathbf{H}_{\star}}^{(j)}(\mathbf{x}_2) \odot \mathbf{A}_1^{\star, \diamond}(\mathbf{x}_2)) \right] \|\mathbf{J}^{\dagger}\|_{\mathbf{H}_{\diamond}}. \quad (31)
 \end{aligned}$$

Denoting

$$\mathbf{v}^{\star, \diamond}(\mathbf{x}) = \sigma_{f_{\diamond}, \mathbf{H}_{\diamond}}^{(1)}(\mathbf{x}) \sigma_{f_{\star}, \mathbf{H}_{\star}}(\mathbf{x}) \odot \mathbf{A}_1^{\star, \diamond}(\mathbf{x}),$$

and rearranging equation 31 give us the Theorem 3.1. □

D.2. Proof of Proposition 3.1

From the proof of Theorem 3.1 in the above subsection, we can see why this proposition holds.

Proposition D.1 (Proposition 3.1 Restated). *In Theorem 3.1,*

$$0 \leq \frac{\mathbb{E}[\mathbf{v}^{*,\diamond}(\mathbf{x}_1)^\top \mathbf{A}_2^{*,\diamond}(\mathbf{x}_1, \mathbf{x}_2) \mathbf{v}^{*,\diamond}(\mathbf{x}_2)]}{\|\nabla f_\star^\top\|_{\mathcal{D}, \mathbf{H}_\star}^2 \cdot \|\mathbf{J}^\dagger\|_{\mathbf{H}_\diamond}^{-1}} \leq 1.$$

Proof. Recall Theorem 3.1 states

$$\min_{g \in \mathbb{G}} \|\nabla f_\star^\top - \nabla(g \circ f_\diamond)^\top\|_{\mathcal{D}, \mathbf{H}_\star} = \left(1 - \frac{\mathbb{E}[\mathbf{v}^{*,\diamond}(\mathbf{x}_1)^\top \mathbf{A}_2^{*,\diamond}(\mathbf{x}_1, \mathbf{x}_2) \mathbf{v}^{*,\diamond}(\mathbf{x}_2)]}{\|\nabla f_\star^\top\|_{\mathcal{D}, \mathbf{H}_\star}^2 \cdot \|\mathbf{J}^\dagger\|_{\mathbf{H}_\diamond}^{-1}}\right)^{\frac{1}{2}} \|\nabla f_\star^\top\|_{\mathcal{D}, \mathbf{H}_\star}.$$

We can see that the ≤ 1 part stands, since $\min_{g \in \mathbb{G}} \|\nabla f_\star^\top - \nabla(g \circ f_\diamond)^\top\|_{\mathcal{D}, \mathbf{H}_\star}$ is always non-negative.

The ≥ 0 part can be proved by observing

$$\begin{aligned} \left(1 - \frac{\mathbb{E}[\mathbf{v}^{*,\diamond}(\mathbf{x}_1)^\top \mathbf{A}_2^{*,\diamond}(\mathbf{x}_1, \mathbf{x}_2) \mathbf{v}^{*,\diamond}(\mathbf{x}_2)]}{\|\nabla f_\star^\top\|_{\mathcal{D}, \mathbf{H}_\star}^2 \cdot \|\mathbf{J}^\dagger\|_{\mathbf{H}_\diamond}^{-1}}\right)^{\frac{1}{2}} \|\nabla f_\star^\top\|_{\mathcal{D}, \mathbf{H}_\star} &= \min_{g \in \mathbb{G}} \|\nabla f_\star^\top - \nabla(g \circ f_\diamond)^\top\|_{\mathcal{D}, \mathbf{H}_\star} \\ &\leq \|\nabla f_\star^\top - \nabla(0 \circ f_\diamond)^\top\|_{\mathcal{D}, \mathbf{H}_\star} = \|\nabla f_\star^\top\|_{\mathcal{D}, \mathbf{H}_\star} \end{aligned}$$

□

D.3. Proof of Theorem 3.2

We introduce two lemmas before proving Theorem 3.2.

Lemma D.3. *Assume that function $h(\cdot)$ satisfies the β -smoothness under $\|\cdot\|_{\mathbf{H}_\star}$ norm (Assumption 1), and assume there is a vector \mathbf{x}_0 in the same space as $\mathbf{x} \sim \mathcal{D}$ such that $h(\mathbf{x}_0) = 0$. Given $\tau > 0$, there exists \mathbf{x}' as a function of \mathbf{x} such that $\|\mathbf{x} - \mathbf{x}'\|_2 \leq \tau$, and*

$$\|h(\mathbf{x})\|_{\mathbf{H}_\star}^2 \leq 2 \left(\|\nabla h(\mathbf{x}')^\top\|_{\mathbf{H}_\star}^2 + \beta^2 (\|\mathbf{x} - \mathbf{x}_0\|_2 - \tau)_+^2 \right) \cdot \|\mathbf{x} - \mathbf{x}_0\|_2^2,$$

where the $(\cdot)_+$ is an operator defined by $\forall x \in \mathbb{R}: (x)_+ = x$ if $x \geq 0$ and $(x)_+ = 0$ otherwise.

Proof. To begin with, we note that the assumption of $h(\mathbf{x}_0) = 0$ is only used for this lemma, and the assumption will be naturally guaranteed when we invoke this lemma in the proof of Theorem 3.2.

With the smoothness assumption, we know that $h(\cdot)$ has continuous gradient. Thus, we have

$$\|h(\mathbf{x})\|_{\mathbf{H}_\star} = \|h(\mathbf{x}) - h(\mathbf{x}_0)\|_{\mathbf{H}_\star} = \|\nabla h(\mathbf{x}_0 + \xi(\mathbf{x} - \mathbf{x}_0))^\top (\mathbf{x} - \mathbf{x}_0)\|_{\mathbf{H}_\star},$$

where the last equation is by mean value theorem and thus $\xi \in (0, 1)$.

Then, noting that $\|\cdot\|_{\mathbf{H}_\star}$ and $\|\cdot\|_2$ are compatible (Lemma E.1), we have

$$\|\nabla h(\mathbf{x}_0 + \xi(\mathbf{x} - \mathbf{x}_0))^\top (\mathbf{x} - \mathbf{x}_0)\|_{\mathbf{H}_\star} \leq \|\nabla h(\mathbf{x}_0 + \xi(\mathbf{x} - \mathbf{x}_0))^\top\|_{\mathbf{H}_\star} \cdot \|\mathbf{x} - \mathbf{x}_0\|_2.$$

Now we discuss two cases to define a random variable \mathbf{x}' as a function of \mathbf{x} .

If $(1 - \xi)\|\mathbf{x} - \mathbf{x}_0\|_2 \leq \tau$, we define \mathbf{x}' as

$$\mathbf{x}' = \mathbf{x}_0 + \xi(\mathbf{x} - \mathbf{x}_0),$$

and we can see that $\|\mathbf{x}' - \mathbf{x}\|_2 \leq \tau$.

Otherwise, *i.e.*, $(1 - \xi)\|\mathbf{x} - \mathbf{x}_0\|_2 > \tau$, we apply triangle inequality to derive

$$\begin{aligned} & \|\nabla h(\mathbf{x}_0 + \xi(\mathbf{x} - \mathbf{x}_0))^\top (\mathbf{x} - \mathbf{x}_0)\|_{\mathbf{H}_*} \\ &= \|\nabla h(\mathbf{x}_0 + \xi(\mathbf{x} - \mathbf{x}_0))^\top - \nabla h(\mathbf{x} - \tau(\widehat{\mathbf{x}} - \mathbf{x}_0))^\top + \nabla h(\mathbf{x} - \tau(\widehat{\mathbf{x}} - \mathbf{x}_0))^\top\|_{\mathbf{H}_*} \\ &\leq \underbrace{\|\nabla h(\mathbf{x}_0 + \xi(\mathbf{x} - \mathbf{x}_0))^\top - \nabla h(\mathbf{x} - \tau(\widehat{\mathbf{x}} - \mathbf{x}_0))^\top\|_{\mathbf{H}_*}}_X + \|\nabla h(\mathbf{x} - \tau(\widehat{\mathbf{x}} - \mathbf{x}_0))^\top\|_{\mathbf{H}_*}, \end{aligned}$$

where we define

$$\mathbf{x}' = \mathbf{x} - \tau(\widehat{\mathbf{x}} - \mathbf{x}_0).$$

By definition, in this case $\|\mathbf{x}' - \mathbf{x}\|_2 \leq \tau$ as well. We then treat X : it can be bounded using β -smoothness, *i.e.*,

$$\begin{aligned} X &\leq \beta\|\mathbf{x}_0 + \xi(\mathbf{x} - \mathbf{x}_0) - \mathbf{x} + \tau(\widehat{\mathbf{x}} - \mathbf{x}_0)\|_2 \\ &= \beta\|\tau(\widehat{\mathbf{x}} - \mathbf{x}_0) - (1 - \xi)(\mathbf{x} - \mathbf{x}_0)\|_2 \\ &= \beta|\tau - (1 - \xi) \cdot \|(\mathbf{x} - \mathbf{x}_0)\|_2| \\ &= \beta((1 - \xi) \cdot \|(\mathbf{x} - \mathbf{x}_0)\|_2 - \tau), \end{aligned}$$

where the last step is because we are exactly considering the case of $(1 - \xi) \cdot \|(\mathbf{x} - \mathbf{x}_0)\|_2 > \tau$.

Therefore, combining the two cases together, we can write

$$\|\nabla h(\mathbf{x}_0 + \xi(\mathbf{x} - \mathbf{x}_0))^\top\|_{\mathbf{H}_*} \leq \beta((1 - \xi) \cdot \|(\mathbf{x} - \mathbf{x}_0)\|_2 - \tau)_+ + \|\nabla h(\mathbf{x}')^\top\|_{\mathbf{H}_*},$$

where $\|\mathbf{x} - \mathbf{x}'\|_2 \leq \tau$.

Combining the above, we have

$$\|h(\mathbf{x})\|_{\mathbf{H}_*} \leq (\|\nabla h(\mathbf{x}')^\top\|_{\mathbf{H}_*} + \beta(\|\mathbf{x} - \mathbf{x}_0\|_2 - \tau)_+) \cdot \|\mathbf{x} - \mathbf{x}_0\|_2.$$

Take the square on both sides, and apply the Cauchy-Schwarz inequality, we have the lemma proved.

$$\begin{aligned} \|h(\mathbf{x})\|_{\mathbf{H}_*}^2 &\leq (\|\nabla h(\mathbf{x}')^\top\|_{\mathbf{H}_*} + \beta(\|\mathbf{x} - \mathbf{x}_0\|_2 - \tau)_+)^2 \cdot \|\mathbf{x} - \mathbf{x}_0\|_2^2 \\ &\leq 2 \left(\|\nabla h(\mathbf{x}')^\top\|_{\mathbf{H}_*}^2 + \beta^2 (\|\mathbf{x} - \mathbf{x}_0\|_2 - \tau)_+^2 \right) \cdot \|\mathbf{x} - \mathbf{x}_0\|_2^2. \end{aligned}$$

□

Lemma D.4. Assume that function $h(\cdot)$ satisfies the β -smoothness under $\|\cdot\|_{\mathbf{H}_*}$ norm (Assumption 1). Given $\tau > 0$, there exists \mathbf{x}'_i as a function of \mathbf{x} for $\forall i \in [n]$ such that $\|\mathbf{x} - \mathbf{x}'_i\|_2 \leq \tau$, and

$$\tau^2 \cdot \|\nabla h(\mathbf{x})^\top\|_{\mathbf{H}_*}^2 \leq 3 \left(\sum_{i=1}^n \|h(\mathbf{x}'_i)\|_{\mathbf{H}_*}^2 + n\|h(\mathbf{x})\|_{\mathbf{H}_*}^2 + n\tau^4\beta^2 \right).$$

Proof. Denote the dimension of \mathbf{x} as n , and let \mathbf{U} be an orthogonal matrix in $\mathbb{R}^{n \times n}$, where we denote its column vectors as $\mathbf{u}_i \in \mathbb{R}^n$ for $i \in [n]$. Applying the mean value theorem, there exists $\xi_i \in (0, 1)$ such that

$$\begin{aligned} h(\mathbf{x} + \tau\mathbf{u}_i) - h(\mathbf{x}) &= \nabla h(\mathbf{x} + \tau\xi_i\mathbf{u}_i)^\top \tau\mathbf{u}_i \\ &= \tau (\nabla h(\mathbf{x})^\top \mathbf{u}_i + (\nabla h(\mathbf{x} + \tau\xi_i\mathbf{u}_i)^\top - \nabla h(\mathbf{x})^\top) \mathbf{u}_i). \end{aligned}$$

Rearranging the equality, we have

$$\nabla h(\mathbf{x})^\top \mathbf{u}_i = \frac{1}{\tau} \gamma_i,$$

where we denote

$$\gamma_i = h(\mathbf{x} + \tau \mathbf{u}_i) - h(\mathbf{x}) - \tau(\nabla h(\mathbf{x} + \tau \xi_i \mathbf{u}_i)^\top - \nabla h(\mathbf{x})^\top) \mathbf{u}_i.$$

Collecting each γ_i for $i \in [n]$ into a matrix $\mathbf{\Gamma} = [\gamma_1 \dots \gamma_n]$, we can re-formulate the above equality as

$$\begin{aligned} \tau \nabla h(\mathbf{x})^\top \mathbf{U} &= \mathbf{\Gamma} \\ \tau \nabla h(\mathbf{x})^\top &= \mathbf{\Gamma} \mathbf{U}^\top, \end{aligned}$$

where the last equality is because that \mathbf{U} is orthogonal.

Taking the $\|\cdot\|_{\mathbf{H}_*}^2$ on both sides, with some linear algebra manipulation we can derive

$$\begin{aligned} \tau^2 \cdot \|\nabla h(\mathbf{x})^\top\|_{\mathbf{H}_*}^2 &= \|\mathbf{\Gamma} \mathbf{U}^\top\|_{\mathbf{H}_*}^2 \\ &= \text{tr}(\mathbf{U} \mathbf{\Gamma}^\top \mathbf{H}_* \mathbf{\Gamma} \mathbf{U}^\top) = \text{tr}(\mathbf{\Gamma}^\top \mathbf{H}_* \mathbf{\Gamma}) = \text{tr}(\mathbf{H}_* \mathbf{\Gamma} \mathbf{\Gamma}^\top) \\ &= \text{tr}(\mathbf{H}_* \sum_{i=1}^n \gamma_i \gamma_i^\top) = \sum_{i=1}^n \text{tr}(\mathbf{H}_* \gamma_i \gamma_i^\top) = \sum_{i=1}^n \text{tr}(\gamma_i^\top \mathbf{H}_* \gamma_i) \\ &= \sum_{i=1}^n \|\gamma_i\|_{\mathbf{H}_*}^2. \end{aligned} \tag{32}$$

Taking $\|\gamma_i\|_{\mathbf{H}_*}$ to work on further, we can derive its upper bound as

$$\begin{aligned} \|\gamma_i\|_{\mathbf{H}_*} &= \|h(\mathbf{x} + \tau \mathbf{u}_i) - h(\mathbf{x}) - \tau(\nabla h(\mathbf{x} + \tau \xi_i \mathbf{u}_i)^\top - \nabla h(\mathbf{x})^\top) \mathbf{u}_i\|_{\mathbf{H}_*} \\ &\leq \|h(\mathbf{x} + \tau \mathbf{u}_i)\|_{\mathbf{H}_*} + \|h(\mathbf{x})\|_{\mathbf{H}_*} + \tau \|(\nabla h(\mathbf{x} + \tau \xi_i \mathbf{u}_i)^\top - \nabla h(\mathbf{x})^\top) \mathbf{u}_i\|_{\mathbf{H}_*} \\ &\leq \|h(\mathbf{x} + \tau \mathbf{u}_i)\|_{\mathbf{H}_*} + \|h(\mathbf{x})\|_{\mathbf{H}_*} + \tau \|\nabla h(\mathbf{x} + \tau \xi_i \mathbf{u}_i)^\top - \nabla h(\mathbf{x})^\top\|_{\mathbf{H}_*} \\ &\leq \|h(\mathbf{x} + \tau \mathbf{u}_i)\|_{\mathbf{H}_*} + \|h(\mathbf{x})\|_{\mathbf{H}_*} + \tau^2 \beta \xi_i \\ &\leq \|h(\mathbf{x} + \tau \mathbf{u}_i)\|_{\mathbf{H}_*} + \|h(\mathbf{x})\|_{\mathbf{H}_*} + \tau^2 \beta, \end{aligned} \tag{33}$$

where the first inequality is by triangle inequality, the second inequality is by Lemma E.1 and the fact that $\|\mathbf{u}_i\|_2 = 1$, the third inequality is done by applying the β -smoothness assumption, and the last inequality is by the fact that $\xi_i \in (0, 1)$ from the mean value theorem.

Plugging the equation 33 into equation 32, we have

$$\begin{aligned} \tau^2 \cdot \|\nabla h(\mathbf{x})^\top\|_{\mathbf{H}_*}^2 &\leq \sum_{i=1}^n (\|h(\mathbf{x} + \tau \mathbf{u}_i)\|_{\mathbf{H}_*} + \|h(\mathbf{x})\|_{\mathbf{H}_*} + \tau^2 \beta)^2 \\ &\leq \sum_{i=1}^n 3 (\|h(\mathbf{x} + \tau \mathbf{u}_i)\|_{\mathbf{H}_*}^2 + \|h(\mathbf{x})\|_{\mathbf{H}_*}^2 + \tau^4 \beta^2) \\ &= 3 \sum_{i=1}^n \|h(\mathbf{x} + \tau \mathbf{u}_i)\|_{\mathbf{H}_*}^2 + 3n \|h(\mathbf{x})\|_{\mathbf{H}_*}^2 + 3n \tau^4 \beta^2, \end{aligned}$$

where the inequality is done Cauchy-Schwarz inequality.

Denoting $\mathbf{x}'_i = \mathbf{x} + \tau \mathbf{u}_i$, we have the lemma proved. \square

Theorem D.2 (Theorem 3.2 Restated). *Given a data distribution \mathcal{D} and $\tau > 0$, there exist distributions $\mathcal{D}_1, \mathcal{D}_2$ such that the type-1 Wasserstein distance $W_1(\mathcal{D}, \mathcal{D}_1) \leq \tau$ and $W_1(\mathcal{D}, \mathcal{D}_2) \leq \tau$ satisfying*

$$\begin{aligned} \frac{1}{2B^2} \|h_{*,\diamond}\|_{\mathcal{D}, \mathbf{H}_*}^2 &\leq \|\nabla h'_{*,\diamond}\|_{\mathcal{D}_1, \mathbf{H}_*}^2 + \beta^2 (B - \tau)_+^2 \\ \frac{1}{3n} \|\nabla h'_{*,\diamond}\|_{\mathcal{D}, \mathbf{H}_*}^2 &\leq \frac{2}{\tau^2} \|h_{*,\diamond}\|_{\mathcal{D}_2, \mathbf{H}_*}^2 + \beta^2 \tau^2, \end{aligned}$$

where n is the dimension of $\mathbf{x} \sim \mathcal{D}$, and $B = \inf_{\mathbf{x}_0 \in \mathbb{R}^n} \sup_{\mathbf{x} \in \text{supp}(\mathcal{D})} \|\mathbf{x} - \mathbf{x}_0\|_2$ is the radius of the $\text{supp}(\mathcal{D})$. The $(\cdot)_+$ is an operator defined by $\forall x \in \mathbb{R}: (x)_+ = x$ if $x \geq 0$ and $(x)_+ = 0$ otherwise.

Proof. Let us begin with recalling the definition of $h_{*,\diamond}$ and $h'_{*,\diamond}$.

The optimal affine transformation $g \in \mathbb{G}$ in the function matching distance (13) is \tilde{g} , and one of the optimal $g \in \mathbb{G}$ in the gradient matching distance is (14) \tilde{g}' . Accordingly, we denote

$$h_{*,\diamond} := f_* - \tilde{g} \circ f_\diamond \quad \text{and} \quad h'_{*,\diamond} := f_* - \tilde{g}' \circ f_\diamond,$$

and we can see that the gradient matching distance and the function matching distance can be written as

$$(13) = \|h_{*,\diamond}\|_{\mathcal{D}, \mathbf{H}_*} \quad \text{and} \quad (14) = \|\nabla h'_{*,\diamond}\|_{\mathcal{D}, \mathbf{H}_*}^\top. \quad (34)$$

The first inequality. Then, we can prove the first inequality using Lemma D.3.

Let $\mathbf{x}_0 \in \mathbb{R}^n$ be a free variable, and then set $\mathbf{b} = h'_{*,\diamond}(\mathbf{x}_0)$. Noting that $\|h_{*,\diamond}\|_{\mathcal{D}, \mathbf{H}_*}^2$ by definition is the minimum of this function distance, we have

$$\|h_{*,\diamond}\|_{\mathcal{D}, \mathbf{H}_*}^2 \leq \|h'_{*,\diamond} - \mathbf{b}\|_{\mathcal{D}, \mathbf{H}_*}^2. \quad (35)$$

Denoting $h := h'_{*,\diamond} - \mathbf{b}$, we can see $h(\mathbf{x}_0) = 0$. Therefore, h can be used to invoke Lemma D.3. That is, there exists \mathbf{x}' as a function of \mathbf{x} such that $\|\mathbf{x} - \mathbf{x}'\|_2 \leq \tau$, and

$$\|h(\mathbf{x})\|_{\mathbf{H}_*}^2 \leq 2 \left(\|\nabla h(\mathbf{x}')^\top\|_{\mathbf{H}_*}^2 + \beta^2 (\|\mathbf{x} - \mathbf{x}_0\|_2 - \tau)_+^2 \right) \cdot \|\mathbf{x} - \mathbf{x}_0\|_2^2,$$

Taking the expectation of $\mathbf{x} \sim \mathcal{D}$ of the both sides, and denote the induced distribution for \mathbf{x}' as \mathcal{D}_1 , we have

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \|h(\mathbf{x})\|_{\mathbf{H}_*}^2 \leq 2 \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left(\|\nabla h(\mathbf{x}')^\top\|_{\mathbf{H}_*}^2 + \beta^2 (\|\mathbf{x} - \mathbf{x}_0\|_2 - \tau)_+^2 \right) \cdot \|\mathbf{x} - \mathbf{x}_0\|_2^2.$$

Recall that \mathbf{x}_0 is a free variable, we can tighten the bound by

$$\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \|h(\mathbf{x})\|_{\mathbf{H}_*}^2 \leq \inf_{\mathbf{x}_0 \in \mathbb{R}^n} 2 \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left(\|\nabla h(\mathbf{x}')^\top\|_{\mathbf{H}_*}^2 + \beta^2 (\|\mathbf{x} - \mathbf{x}_0\|_2 - \tau)_+^2 \right) \cdot \|\mathbf{x} - \mathbf{x}_0\|_2^2. \quad (36)$$

Note that we can have tighter but similar results if we keep the $\inf_{\mathbf{x}_0 \in \mathbb{R}^n}$. However, by plugging in the radius

$$B = \inf_{\mathbf{x}_0 \in \mathbb{R}^n} \sup_{\mathbf{x} \in \text{supp}(\mathcal{D})} \|\mathbf{x} - \mathbf{x}_0\|_2$$

we can make the presentation much more simplified without losing its core messages.

That is,

$$(36) \leq 2 \left(\mathbb{E}_{\mathbf{x}' \sim \mathcal{D}_1} \|\nabla h(\mathbf{x}')^\top\|_{\mathbf{H}_*}^2 + \beta^2 (B - \tau)_+^2 \right) B^2.$$

Combining the above inequality and equation 35, and noting that

$$\begin{aligned} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \|h(\mathbf{x})\|_{\mathbf{H}_*}^2 &= \|h\|_{\mathcal{D}, \mathbf{H}_*}^2 \\ \mathbb{E}_{\mathbf{x}' \sim \mathcal{D}_1} \|\nabla h(\mathbf{x}')^\top\|_{\mathbf{H}_*}^2 &= \|\nabla h^\top\|_{\mathcal{D}_1, \mathbf{H}_*}^2, \end{aligned}$$

we have

$$\begin{aligned} \|h_{*,\diamond}\|_{\mathcal{D}, \mathbf{H}_*}^2 &\leq \|h'_{*,\diamond} - \mathbf{b}\|_{\mathcal{D}, \mathbf{H}_*}^2 = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \|h(\mathbf{x})\|_{\mathbf{H}_*}^2 \\ &\leq 2 \left(\mathbb{E}_{\mathbf{x}' \sim \mathcal{D}_1} \|\nabla h(\mathbf{x}')^\top\|_{\mathbf{H}_*}^2 + \beta^2 (B - \tau)_+^2 \right) B^2 \\ &= 2 \left(\|\nabla h^\top\|_{\mathcal{D}_1, \mathbf{H}_*}^2 + \beta^2 (B - \tau)_+^2 \right) B^2 \end{aligned}$$

Noting that h and $h'_{*,\diamond}$ only differs by a constant shift \mathbf{b} , we can see $\nabla h = \nabla h'_{*,\diamond}$. Therefore, by replacing ∇h^\top by $\nabla h'_{*,\diamond}^\top$ we finally have the first inequality in Theorem 3.2

$$\|h_{*,\diamond}\|_{\mathcal{D}, \mathbf{H}_*}^2 \leq 2 \left(\|\nabla h'_{*,\diamond}^\top\|_{\mathcal{D}_1, \mathbf{H}_*}^2 + \beta^2 (B - \tau)_+^2 \right) B^2.$$

It remains to show the Wasserstein distance between \mathcal{D}_1 and \mathcal{D} . As \mathbf{x}' is a function of the random variable $\mathbf{x} \sim \mathcal{D}$ with $\|\mathbf{x}' - \mathbf{x}\|_2 \leq \tau$, and \mathcal{D}_1 is the induced distribution of \mathbf{x}' as a function of \mathbf{x} , we can see that by the definition of type-1 Wasserstein distance between \mathcal{D} and \mathcal{D}_1 is bounded by τ .

Denote $\mathbb{J}(\mathcal{D}, \mathcal{D}')$ as the set of all joint distributions that have marginals \mathcal{D} and \mathcal{D}' , and recall the definition of type-1 Wasserstein distance is

$$W_1(\mathcal{D}, \mathcal{D}_1) = \inf_{\mathcal{J} \in \mathbb{J}(\mathcal{D}, \mathcal{D}_1)} \int \|\mathbf{x} - \mathbf{x}'\|_2 d\mathcal{J}(\mathbf{x}, \mathbf{x}').$$

Denote \mathcal{J}_0 as the joint distribution such that in $(\mathbf{x}, \mathbf{x}') \sim \mathcal{J}$ we always have \mathbf{x}' being a function of \mathbf{x} as how \mathbf{x}' is defined. We can see that

$$\begin{aligned} W_1(\mathcal{D}, \mathcal{D}_1) &= \inf_{\mathcal{J} \in \mathbb{J}(\mathcal{D}, \mathcal{D}_1)} \int \|\mathbf{x} - \mathbf{x}'\|_2 d\mathcal{J}(\mathbf{x}, \mathbf{x}') \leq \int \|\mathbf{x} - \mathbf{x}'\|_2 d\mathcal{J}_0(\mathbf{x}, \mathbf{x}') \leq \int \tau d\mathcal{J}_0(\mathbf{x}, \mathbf{x}') \\ &= \tau. \end{aligned} \quad (37)$$

Therefore, we have the first inequality in the theorem proved .

The second inequality. Invoking Lemma D.4 with $h_{\star, \diamond}$, and rearranging the inequality, we have

$$\frac{1}{3n} \|\nabla h_{\star, \diamond}(\mathbf{x})^\top\|_{\mathbf{H}_\star}^2 \leq \frac{2}{\tau^2} \left(\sum_{i=1}^n \frac{1}{2n} \|h_{\star, \diamond}(\mathbf{x}'_i)\|_{\mathbf{H}_\star}^2 + \frac{1}{2} \|h_{\star, \diamond}(\mathbf{x})\|_{\mathbf{H}_\star}^2 \right) + \tau^2 \beta^2.$$

Taking the expectation on both sides, we have

$$\frac{1}{3n} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \|\nabla h_{\star, \diamond}(\mathbf{x})^\top\|_{\mathbf{H}_\star}^2 \leq \frac{2}{\tau^2} \underbrace{\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left(\sum_{i=1}^n \frac{1}{2n} \|h_{\star, \diamond}(\mathbf{x}'_i)\|_{\mathbf{H}_\star}^2 + \frac{1}{2} \|h_{\star, \diamond}(\mathbf{x})\|_{\mathbf{H}_\star}^2 \right)}_X + \tau^2 \beta^2. \quad (38)$$

Note that X can be reformulated to be the expectation of an induced distribution from $\mathbf{x} \sim \mathcal{D}$, since \mathbf{x}'_i is a pre-defined function of \mathbf{x} . Denote \mathcal{D}_2 as the distribution induced by the following sampling process: first, sample $\mathbf{x} \sim \mathcal{D}$; then,

$$\begin{aligned} \mathbf{x}' &= \mathbf{x} && \text{with probability } \frac{1}{2} \\ \mathbf{x}' &= \mathbf{x}'_i && \text{with probability } \frac{1}{2n} \text{ for } \forall i \in [n]. \end{aligned}$$

Therefore, we can write X as

$$X = \|h_{\star, \diamond}\|_{\mathcal{D}_2, \mathbf{H}_\star}^2. \quad (39)$$

Similarly to equation 37, it also holds that $W_1(\mathcal{D}, \mathcal{D}_2) \leq \tau$.

To finally complete the proof, noting that $\|\nabla h_{\star, \diamond}^\top\|_{\mathcal{D}, \mathbf{H}_\star}^2$ is the minimum of this gradient distance (equation 34), we have

$$\|\nabla h_{\star, \diamond}^\top\|_{\mathcal{D}, \mathbf{H}_\star}^2 \leq \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \|\nabla h_{\star, \diamond}(\mathbf{x})^\top\|_{\mathbf{H}_\star}^2. \quad (40)$$

Combining equation 38, equation 39 and equation 40, we have the second inequality proved.

Hence, we have proved Theorem 3.2. \square

D.4. Proof of Theorem 3.3

Theorem D.3 (Theorem 3.3 Restated). *The surrogate transfer loss (16) and the true transfer loss (17) are close, with an error of $\|f_T - y\|_{\mathcal{D}, \mathbf{H}_T}$.*

$$-\|f_T - y\|_{\mathcal{D}, \mathbf{H}_T} \leq (17) - (16) \leq \|f_T - y\|_{\mathcal{D}, \mathbf{H}_T}$$

Proof. Let us begin by recall the definition of the surrogate transfer loss (16) and the true transfer loss (17).

$$(16) := \min_{g \in \mathbb{G}} \|f_T - g \circ f_S\|_{\mathcal{D}, \mathbf{H}_T}$$

$$(17) := \min_{g \in \mathbb{G}} \|y - g \circ f_S\|_{\mathcal{D}, \mathbf{H}_T}.$$

Denote

$$\tilde{g}' := \arg \min_{g \in \mathbb{G}} \|f_T - g \circ f_S\|_{\mathcal{D}, \mathbf{H}_T}$$

$$\tilde{g} := \arg \min_{g \in \mathbb{G}} \|y - g \circ f_S\|_{\mathcal{D}, \mathbf{H}_T}.$$

First, we show an upper bound for (16).

$$(16) \leq \|f_T - \tilde{g} \circ f_S\|_{\mathcal{D}, \mathbf{H}_T} \leq \|y - \tilde{g} \circ f_S\|_{\mathcal{D}, \mathbf{H}_T} + \|f_T - y\|_{\mathcal{D}, \mathbf{H}_T} = (17) + \|f_T - y\|_{\mathcal{D}, \mathbf{H}_T}, \quad (41)$$

where the last inequality is by triangle inequality.

Similarly, we can derive its lower bound.

$$\begin{aligned} (16) &= \|f_T - \tilde{g}' \circ f_S\|_{\mathcal{D}, \mathbf{H}_T} \geq \|y - \tilde{g}' \circ f_S\|_{\mathcal{D}, \mathbf{H}_T} - \|f_T - y\|_{\mathcal{D}, \mathbf{H}_T} \\ &\geq \min_{g \in \mathbb{G}} \|y - g \circ f_S\|_{\mathcal{D}, \mathbf{H}_T} - \|f_T - y\|_{\mathcal{D}, \mathbf{H}_T} = (17) - \|f_T - y\|_{\mathcal{D}, \mathbf{H}_T}, \end{aligned} \quad (42)$$

where the first inequality is by triangle inequality.

Combining equation 41 and equation 42, we have the proposition proved. \square

D.5. Proof of Theorem B.1

Theorem D.4 (Theorem B.1 Restated). *Denote $\tilde{g}_{T,S} : \mathbb{R}^m \rightarrow \mathbb{R}^d$ as the optimal solution of equation 16, and $\tilde{g}_{S,T} : \mathbb{R}^d \rightarrow \mathbb{R}^m$ as the optimal solution of equation 18. Suppose the two optimal affine maps $\tilde{g}_{T,S}, \tilde{g}_{S,T}$ are both full-rank. For $\mathbf{v} \in \mathbb{R}^m$, denote the matrix representation of $\tilde{g}_{T,S}$ as $\tilde{g}_{T,S}(\mathbf{v}) = \tilde{\mathbf{W}}_{T,S} \mathbf{v} + \tilde{\mathbf{b}}_{T,S}$. Similarly, for $\mathbf{w} \in \mathbb{R}^d$, denote the matrix representation of $\tilde{g}_{S,T}$ as $\tilde{g}_{S,T}(\mathbf{w}) = \tilde{\mathbf{W}}_{S,T} \mathbf{w} + \tilde{\mathbf{b}}_{S,T}$. We have the following statements.*

If $d < m$, then $\tilde{g}_{S,T}$ is injective, and we have:

$$\|f_T - \tilde{g}_{T,S} \circ f_S\|_{\mathcal{D}, \mathbf{H}_T} \leq \sqrt{\|(\tilde{\mathbf{W}}_{S,T}^\top \tilde{\mathbf{W}}_{S,T})^{-1}\|_F \cdot \|\mathbf{H}_T\|_F} \cdot \|f_S - \tilde{g}_{S,T} \circ f_T\|_{\mathcal{D}}. \quad (19)$$

If $d > m$, then $\tilde{g}_{T,S}$ is injective, and we have:

$$\|f_S - \tilde{g}_{S,T} \circ f_T\|_{\mathcal{D}, \mathbf{H}_S} \leq \sqrt{\|(\tilde{\mathbf{W}}_{T,S}^\top \tilde{\mathbf{W}}_{T,S})^{-1}\|_F \cdot \|\mathbf{H}_S\|_F} \cdot \|f_T - \tilde{g}_{T,S} \circ f_S\|_{\mathcal{D}}. \quad (20)$$

If $d = m$, then both $\tilde{g}_{S,T}$ and $\tilde{g}_{T,S}$ are bijective, and we have both (19) and (20) stand.

Proof. Observing the symmetry, we only need to prove the following claim.

Claim. *For $\star, \diamond \in \{S, T\}$ and $\star \neq \diamond$, if $\tilde{g}_{\star, \diamond}$ is injective, then*

$$\|f_\diamond - \tilde{g}_{\diamond, \star} \circ f_\star\|_{\mathcal{D}, \mathbf{H}_\diamond}^2 \leq \|(\tilde{\mathbf{W}}_{\star, \diamond}^\top \tilde{\mathbf{W}}_{\star, \diamond})^{-1}\|_F \cdot \|\mathbf{H}_\diamond\|_F \cdot \|f_\star - \tilde{g}_{\star, \diamond} \circ f_\diamond\|_{\mathcal{D}}^2.$$

Proof of the Claim. We have mostly done with this claim with Lemma E.3. Noting that $\tilde{g}_{\diamond, \star}$ is the minimizer of $\min_{g \in \mathbb{G}} \|f_\diamond - g \circ f_\star\|_{\mathcal{D}, \mathbf{H}_\diamond}^2$, we have

$$\begin{aligned} \|f_\diamond - \tilde{g}_{\diamond, \star} \circ f_\star\|_{\mathcal{D}, \mathbf{H}_\diamond}^2 &\leq \|f_\diamond - \tilde{g}_{\star, \diamond}^{-1} \circ f_\star\|_{\mathcal{D}, \mathbf{H}_\diamond}^2 = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\|f_\diamond(\mathbf{x}) - \tilde{g}_{\star, \diamond}^{-1}(f_\star(\mathbf{x}))\|_{\mathbf{H}_\diamond}^2] \\ &\leq \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\|(\tilde{\mathbf{W}}_{\star, \diamond}^\top \tilde{\mathbf{W}}_{\star, \diamond})^{-1}\|_F \cdot \|\mathbf{H}_\diamond\|_F \cdot \|f_\star(\mathbf{x}) - \tilde{g}_{\star, \diamond}(f_\diamond(\mathbf{x}))\|_2^2] \\ &= \|(\tilde{\mathbf{W}}_{\star, \diamond}^\top \tilde{\mathbf{W}}_{\star, \diamond})^{-1}\|_F \cdot \|\mathbf{H}_\diamond\|_F \cdot \|f_\star - \tilde{g}_{\star, \diamond} \circ f_\diamond\|_{\mathcal{D}}^2, \end{aligned}$$

where the second inequality is by invoking Lemma E.3. \square

Taking the square root of this claim, and applying ($\diamond = T, \star = S$) or ($\diamond = S, \star = T$), we immediately have the first two statements about the case of $d < m$ or $d > m$. Finally, noting that when $m = d$, both $\tilde{g}_{S,T}$ and $\tilde{g}_{T,S}$ are bijective and thus also injective, we can see that both (19) and (20) stand. \square

E. Auxiliary Lemmas

Lemma E.1 (Compatibility of $\|\cdot\|_H$ and $\|\cdot\|_2$). *Let $H \in \mathbb{R}^{m \times m}$ be a positive semi-definite matrix, and denote $H = T^\top T$ as its symmetric decomposition with $T \in \mathbb{R}^{m \times m}$. For $W \in \mathbb{R}^{m \times n}$ and $v \in \mathbb{R}^n$, we have*

$$\|Wv\|_H \leq \|W\|_H \cdot \|v\|_2.$$

Proof.

$$\begin{aligned} \|Wv\|_H^2 &= v^\top W^\top T^\top T W v = \|TWv\|_2^2 \\ &\leq \|TW\|_2^2 \cdot \|v\|_2^2 \leq \|TW\|_F^2 \cdot \|v\|_2^2, \end{aligned}$$

where $\|\cdot\|_F$ is the Frobenius norm. Then, we can continue as

$$\|TW\|_F^2 = \text{tr}(W^\top T^\top T W) = \text{tr}(W^\top H W) = \|W\|_H^2.$$

Combining the above two parts, we have the lemma proved. \square

Lemma E.2 (Expectation Preserves the Inclusion Relationship Between Linear Spaces). *Given a distribution $x \sim \mathcal{D}$ in \mathbb{R}^n , we denote the associated probability measure as μ . Given linear maps $M_x : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and $N_x : \mathbb{R}^n \rightarrow \mathbb{R}^d$, noting that they are both functions of x , we have the following statement.*

$$\ker(\mathbb{E}_{x \sim \mathcal{D}} M_x^\top M_x) \subseteq \ker(\mathbb{E}_{x \sim \mathcal{D}} N_x^\top M_x),$$

where $\ker(\cdot)$ denotes the kernel space of a given liner map.

Proof. It suffice to show for $\forall v \in \ker(\mathbb{E}_{x \sim \mathcal{D}} M_x^\top M_x)$, we also have $v \in \ker(\mathbb{E}_{x \sim \mathcal{D}} N_x^\top M_x)$.

Denote $P := \mathbb{E}_{x \sim \mathcal{D}} M_x^\top M_x$, and let $v \in \ker(P)$, we have

$$Pv = 0.$$

Noting that P is positive semi-definite, we have the following equivalent statements.

$$v \in \ker(P) \iff v^\top Pv = 0,$$

where the ' \implies ' direction is trivial, and the ' \impliedby ' direction can be proved by decomposing $P = T^\top T$ as two matrices and noting that

$$v^\top T^\top T v = 0 \implies \|Tv\|_2^2 = 0 \implies Tv = 0 \implies T^\top Tv = 0 \implies Pv = 0.$$

Therefore, we have

$$\begin{aligned} &v^\top Pv = 0 \\ \implies &\mathbb{E}_{x \sim \mathcal{D}} [v^\top M_x^\top M_x v] = 0 \\ \implies &\mathbb{E}_{x \sim \mathcal{D}} [\|M_x v\|_2^2] = 0 \\ \implies &\int \|M_x v\|_2^2 d\mu = 0, \end{aligned}$$

which implies $M_x v = 0$ almost everywhere w.r.t. μ .

Therefore, applying \mathbf{v} to $\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[\mathbf{N}_x^\top \mathbf{M}_x]$ and we have

$$\begin{aligned} \mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[\mathbf{N}_x^\top \mathbf{M}_x] \mathbf{v} &= \int \mathbf{N}_x^\top \mathbf{M}_x \mathbf{v} \, d\mu \\ &= \int_{a.e.} \mathbf{N}_x^\top \mathbf{0} \, d\mu \\ &= \mathbf{0}, \end{aligned}$$

which means $\mathbf{v} \in \ker(\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \mathbf{N}_x^\top \mathbf{M}_x)$. □

Lemma E.3 (Inverse an Injective Linear Map). *Given a full-rank injective affine transformation $g : \mathbb{R}^m \rightarrow \mathbb{R}^d$, we denote its matrix representation as $g(\mathbf{v}) = \mathbf{W}\mathbf{v} + \mathbf{b}$ where $\mathbf{v} \in \mathbb{R}^m$, $\mathbf{W} \in \mathbb{R}^{d \times m}$, $\mathbf{b} \in \mathbb{R}^d$. The inverse of g is $g^{-1} : \mathbb{R}^d \rightarrow \mathbb{R}^m$ defined by $g^{-1}(\mathbf{w}) := \mathbf{W}^\dagger \mathbf{w} - \mathbf{W}^\dagger \mathbf{b}$ for $\mathbf{w} \in \mathbb{R}^d$, i.e., $g^{-1} \circ g$ is the identity function. Moreover, given a positive semi-definite matrix \mathbf{H} , for $\forall \mathbf{v} \in \mathbb{R}^m$ and $\forall \mathbf{w} \in \mathbb{R}^d$, we have*

$$\sqrt{\|(\mathbf{W}^\top \mathbf{W})^{-1}\|_F \cdot \|\mathbf{H}\|_F} \cdot \|\mathbf{w} - g(\mathbf{v})\|_2 \geq \|\mathbf{v} - g^{-1}(\mathbf{w})\|_{\mathbf{H}}.$$

Proof. First, let us verify that $g^{-1} \circ g$ is the identity function. The conditions of g being full-rank and injective are equivalent to \mathbf{W} being full-rank and $d \geq m$. That is being said, $\mathbf{W}^\top \mathbf{W}$ is invertible and $\mathbf{W}^\dagger = (\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top$. Therefore, for $\forall \mathbf{v} \in \mathbb{R}^m$, we have

$$\begin{aligned} g^{-1} \circ g(\mathbf{v}) &= \mathbf{W}^\dagger (\mathbf{W}\mathbf{v} + \mathbf{b}) - \mathbf{W}^\dagger \mathbf{b} = \mathbf{W}^\dagger \mathbf{W}\mathbf{v} \\ &= (\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{W}\mathbf{v} = \mathbf{v}. \end{aligned}$$

That is, $g^{-1} \circ g$ is indeed the identity function.

Next, to prove the inequality, let us start from the right-hand-side of the inequality.

$$\begin{aligned} \|\mathbf{v} - g^{-1}(\mathbf{w})\|_{\mathbf{H}} &= \|g^{-1} \circ g(\mathbf{v}) - g^{-1}(\mathbf{w})\|_{\mathbf{H}} \\ &= \|\mathbf{W}^\dagger (g(\mathbf{v}) - \mathbf{w})\|_{\mathbf{H}} \\ &\leq \|\mathbf{W}^\dagger\|_{\mathbf{H}} \cdot \|g(\mathbf{v}) - \mathbf{w}\|_2, \end{aligned} \tag{43}$$

where the inequality is done by applying Lemma E.1.

To complete the prove, we can see that

$$\begin{aligned} \|\mathbf{W}^\dagger\|_{\mathbf{H}}^2 &= \|(\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top\|_{\mathbf{H}}^2 = \text{tr}(\mathbf{W}(\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{H}(\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top) \\ &= \text{tr}((\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{H}(\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{W}) = \text{tr}((\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{H}) \\ &= \langle (\mathbf{W}^\top \mathbf{W})^{-1}, \mathbf{H} \rangle \\ &\leq \|(\mathbf{W}^\top \mathbf{W})^{-1}\|_F \cdot \|\mathbf{H}\|_F. \end{aligned} \tag{44}$$

Plugging the square root of equation 44 into equation 43, we have the lemma proved. □

F. Additional Details of Synthetic Experiments

In this section, we complete the description of the settings and methods used in the synthetic experiments. Moreover, we report two additional sets of results in cross-architecture scenarios.

In the main paper (section 4), the synthetic experiments are done on the setting where source models have the same architecture as the target model, i.e., all the models are one-hidden-layer neural networks with width $m = 100$. A natural question is what would the results be if using different architectures? That is, the architecture of the source models are different from the target model. To answer this question, we present two additional sets of synthetic experiments where the width of the source models is $m = 50$ or $m = 200$, different from the target model (width $m = 100$).

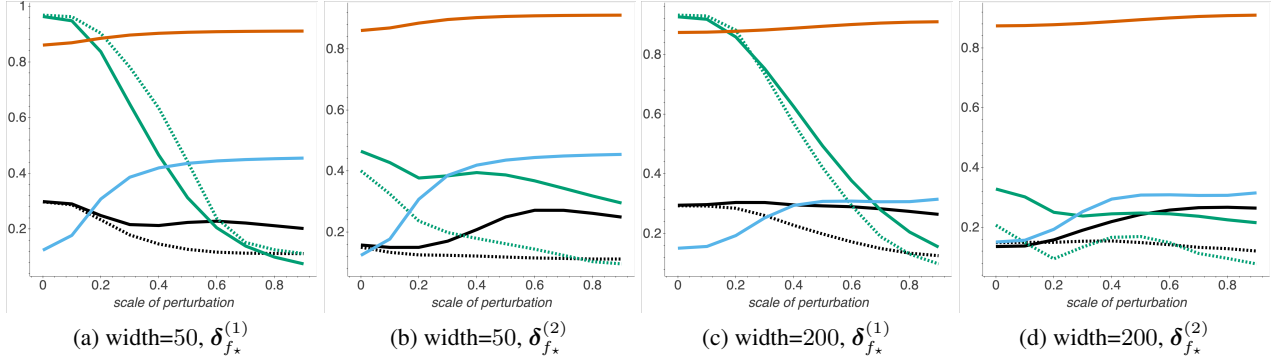


Figure 3. In this figure, 'width' is the width of the source models (one-hidden-layer neural networks). As defined in equation 9, $\delta_{f_*}^{(1)}$ corresponds to the regular adversarial attacks, while $\delta_{f_*}^{(2)}$ the secondary adversarial attack. That is, $\delta_{f_*}^{(2)}$ represents the other information in the adversarial transferring process compared with the first. The x-axis shows the scale of perturbation $t \in [0, 1]$ that controls how much the source model deviates from its corresponding reference source model. There are in total 6 quantities reported. Specifically, $\alpha_1^{f_T \rightarrow f_S}$ is **black solid**; $\alpha_1^{f_S \rightarrow f_T}$ is **black dotted**; $\alpha_2^{f_T \rightarrow f_S}$ is **green solid**; $\alpha_2^{f_S \rightarrow f_T}$ is **green dotted**; the gradient matching loss is **red solid**; and the knowledge transferability distance is **blue solid**.

As we have presented in the main paper about the description of the methods and models used in this experiment, here we present the detailed description of the settings and the datasets being used.

Settings. We follow the small- ϵ setting used in the theory, *i.e.*, the adversarial attack are constrained to a small magnitude, so that we can use its first-order Talyor approximation.

Dataset. Denote a radial basis function as $\phi_i(\mathbf{x}) = e^{-\|\mathbf{x} - \boldsymbol{\mu}_i\|_2^2 / (\sigma_i)^2}$, and for each input data we form its corresponding M -dimensional feature vector as $\boldsymbol{\phi}(\mathbf{x}) = [\phi_1(\mathbf{x}), \dots, \phi_M(\mathbf{x})]^\top$. We set the dimension of \mathbf{x} to be 50. For each radial basis function $\phi_i(\mathbf{x})$, $i \in [M]$, $\boldsymbol{\mu}_i$ is sampled from $U(-0.5, 0.5)^{50}$, and σ_i^2 is sampled from $U(0, 100)$. We use $M = 100$ radial basis functions so that the feature vector is 100-dimensional. Then, we set the target ground truth to be $y(\mathbf{x}) = \mathbf{W}\boldsymbol{\phi}(\mathbf{x}) + \mathbf{b}$ where $\mathbf{W} \in \mathbb{R}^{10 \times 100}$, $\mathbf{b} \in \mathbb{R}^{10}$ are sampled from $U(-0.5, 0.5)$ element-wise. We generate $N = 5000$ samples of \mathbf{x} from a Gaussian mixture formed by 10 Gaussians with different centers but the same covariance matrix $\boldsymbol{\Sigma} = \mathbf{I}$. The centers are sampled randomly from $U(-0.5, 0.5)^{50}$. That is, the dataset $D = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$ consists of $N = 5000$ sample from the distribution, where \mathbf{x}_i is 50-dimensional, \mathbf{y}_i is 10-dimensional. The ground truth target \mathbf{y}_i are computed using the ground truth target function $y(\mathbf{x}_i)$. That is, we want our neural networks to approximate $y(\cdot)$ on the Gaussian mixture.

Methods of Additional Experiments. Note that we have provided the detailed description of the methods used in the main paper synthetic experiments. Here, we present the methods for two additional sets of synthetic experiments, using the same dataset and settings, but different architectures. In the main paper, the source model and the target model are of the same architecture, and the source models are perturbed target model. Here, we use the same target model f_T (width $m = 100$) trained on the dataset D , but two different architectures for source models. That is, the source models and the target model are of different width.

To derive the source models, we first train two reference source models on D with width $m = 50$ and $m = 200$. For each of the reference models, denoting the weights of the model as \mathbf{W} , we randomly sample a direction \mathbf{V} where each entry of \mathbf{V} is sampled from $U(-0.5, 0.5)$, and choose a scale $t \in [0, 1]$. Subsequently, we perturb the model weights of the clean source model as $\mathbf{W}' := \mathbf{W} + t\mathbf{V}$, and define the source model f_S to be a one-hidden-layer neural network with weights \mathbf{W}' . Then, we compute each of the quantities we care about, including α_1, α_2 from both $f_S \rightarrow f_T$ and $f_T \rightarrow f_S$, the gradient matching distance (equation 7), and the actual knowledge transfer distance (equation 17). We use the standard ℓ_2 loss as the adversarial loss function.

Results. We present four sets of result in Figure 3. The indication relations between adversarial transferability and knowledge transferability can be observed in the cross-architecture setting. Moreover: 1. the metrics α_1, α_2 are more meaningful if using the regular attacks; 2. the gradient matching distance tracks the actual knowledge transferability loss; 3. the directions of $f_T \rightarrow f_S$ and $f_S \rightarrow f_T$ are similar.

G. Details of the Empirical Experiments

All experiments are run on a single GTX2080Ti.

G.1. Datasets

G.1.1. IMAGE DATASETS

- **CIFAR10**¹: it consists of 60000 32×32 colour images in 10 classes, with 6000 images per class. There are 50000 training images and 10000 test images.
- **STL10**²: it consists of 13000 labeled 96×96 colour images in 10 classes, with 1300 images per class. There are 5000 training images and 8000 test images. 500 training images (10 pre-defined folds), 800 test images per class.

G.1.2. NLP DATASETS

- **IMDB**³: Document-level sentiment classification on positive and negative movie reviews. We use this dataset to train the target model.
- **AG’s News (AG)**: Sentence-level classification with regard to four news topics: World, Sports, Business, and Science/Technology. Following [Zhang et al. \(2015\)](#), we concatenate the title and description fields for each news article. We use this dataset to train the source model.
- **Fake News Detection (Fake)**: Document-level classification on whether a news article is fake or not. The dataset comes from the Kaggle Fake News Challenge⁴. We concatenate the title and news body of each article. We use this dataset to train the source model.
- **Yelp**: Document-level sentiment classification on positive and negative reviews ([Zhang et al., 2015](#)). Reviews with a rating of 1 and 2 are labeled negative and 4 and 5 positive. We use this dataset to train the source model.

G.2. Adversarial Transferability Indicating Knowledge Transferability

G.2.1. IMAGE

For all the models, both source and target, in the Cifar10 to STL10 experiment, we train them by SGD with momentum and learning rate 0.1 for 100 epochs. For knowledge transferability, we randomly reinitialize and train the source models’ last layer for 10 epochs on STL10. Then we generate adversarial examples with the target model on the validation set and measure the adversarial transferability by feeding these adversarial examples to the source models. We employ two adversarial attacks in this experiment and show that they achieve the same purpose in practice: First, we generate adversarial examples by 50 steps of projected gradient descent and epsilon 0.1 (Results shown in Table 1). Then, we generate adversarial examples by the more efficient FGSM with epsilon 0.1 (Results shown in Table 6) and show that we can efficiently identify candidate models without the expensive PGD attacks.

To further visualize the averaged relation presented in Table 1 and 6, we plot scatter plots Figure 5 and Figure 4 with per sample α_1 as x axis and per sample transfer loss as y axis. Transfer loss is the cross entropy loss predicted by the source model with last layer fine-tuned on STL10. The Pearson score indicates strong correlation between adversarial transferability and knowledge transferability.

We note that in the figures where we report per-sample α_1 , although ideally $\alpha_1 \in [0, 1]$, we can observe that for some samples they have $\alpha_1 > 1$ due to the attacking algorithm is not ideal in practice. However, the introduced sample-level noise does not affect the overall results, *e.g.*, see the averaged results in our tables, or the overall correlation in these figures.

G.2.2. NLP

In the NLP experiments, to train source and target models, we finetune BERT-base models on different datasets for 3 epochs with learning rate equal to $5e - 5$ and warm-up steps equal to the 10% of the total training steps. For knowledge

¹<https://www.cs.toronto.edu/~kriz/cifar.html>

²<https://cs.stanford.edu/~acoates/stl10/>

³<https://datasets.imdbws.com/>

⁴<https://www.kaggle.com/c/fake-news/data>

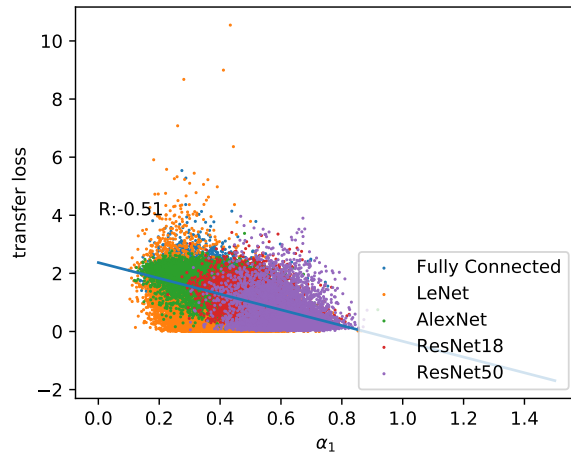


Figure 4. Distribution of per sample knowledge transfer loss and α_1 . The adversarial samples are generated by PGD. The Pearson score shows strong negative correlation between α_1 and the knowledge transfer loss. The higher the transfer loss is, the lower the knowledge transferability is, and the lower the α_1 is.

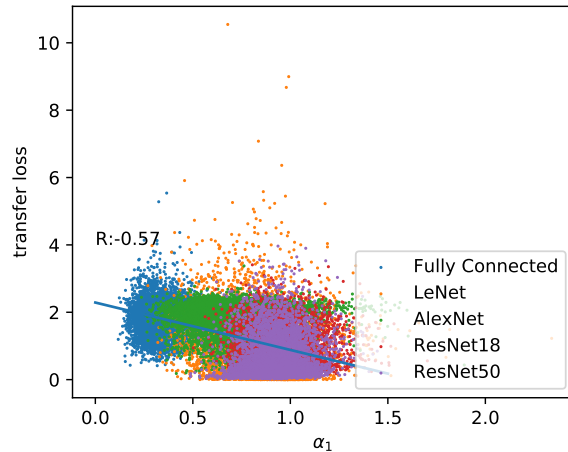


Figure 5. Distribution of per-sample knowledge transfer loss and α_1 . The adversarial samples are generated by FGSM. The Pearson score shows negative strong correlation between α_1 and transfer loss. The higher the transfer loss is, the lower the knowledge transferability is, the lower the α_1 should be.

transferability, we random initialize the last layer of source models and fine-tune all layers of BERT for 1 epoch on the targeted dataset (IMDB). Based on the test data from the target model, we generate 1,000 textual adversarial examples via the state-of-the-art adversarial attacks T3 (Wang et al., 2020) with adversarial learning rate equal to 0.2, maximum iteration steps equal to 100, and $c = \kappa = 100$.

G.2.3. ABLATION STUDIES ON CONTROLLING ADVERSARIAL TRANSFERABILITY

We conduct series of experiments on controlling adversarial transferability between source models and target model by promoting their Loss Gradient Diversity. Demontis et al. (2019) shows that for two models f_S and f_T , the cosine similarity between their loss gradient vectors $\nabla_{x^l} \ell_{f_S}$ and $\nabla_{x^l} \ell_{f_T}$ could be a significant indicator measuring two models’ adversarial transferability. Moreover, Kariyappa & Qureshi (2019) claims that adversarial transferability between two models could be well controlled by regularizing the cosine similarity between their loss gradient vectors. Inspired by this, we train several

Model	Knowledge Trans.	α_1	α_2	$\alpha_1 * \alpha_2$
Fully Connected	28.30	0.279	0.117	0.0103
AlexNet	45.65	0.614	0.208	0.0863
LeNet	55.09	0.803	0.298	0.205
ResNet18	76.60	1.000	0.405	0.410
ResNet50	77.92	0.962	0.392	0.368

Table 6. Knowledge transferability (Knowledge Trans.) among different model architectures. Adversarial examples are generated using FGSM attacks. Our correlation analysis shows Pearson score of -0.57 between the transfer loss and α_1 . Lower transfer loss corresponds to higher transfer accuracy. More details can be found in Figure 5

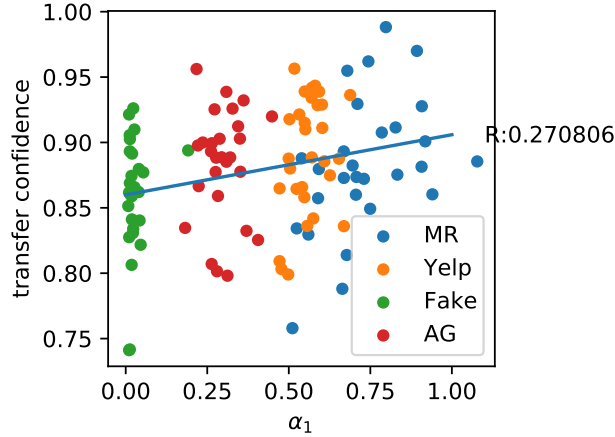


Figure 6. Distribution of per-batch knowledge transfer confidence and α_1 . The Pearson score shows positive correlation between α_1 and transfer confidence. The higher the confidence, the higher the knowledge transferability.

Table 7. Knowledge transferability (Knowledge Trans.) among different source models (controlling adversarial transferability by promoting Loss Gradient Diversity). Adversarial transferability is measured by using the adversarial examples generated against the Target Model to attack the Source Models and estimate α_1 and α_2 .

Model	Knowledge Trans.	α_1	α_2	$\alpha_1 * \alpha_2$
$\rho = 0.0$	73.91	0.394	0.239	0.103
$\rho = 0.5$	73.11	0.385	0.246	0.102
$\rho = 1.0$	72.47	0.371	0.244	0.100
$\rho = 2.0$	71.62	0.370	0.244	0.100
$\rho = 5.0$	72.16	0.378	0.240	0.098

source models f_S to one target model f_T with following training loss:

$$\mathcal{L}_{\text{train}} = \mathcal{L}_{\text{CE}}(f_S(\mathbf{x}), y) + \rho \cdot \mathcal{L}_{\text{cos}}(\nabla_{\hat{\mathbf{x}}} \ell_{f_S}, \nabla_{\hat{\mathbf{x}}} \ell_{f_T})$$

where \mathcal{L}_{CE} refers to cross-entropy loss and $\mathcal{L}_{\text{cos}}(\cdot, \cdot)$ the cosine similarity metric. \mathbf{x} presents *source domain* instances while $\hat{\mathbf{x}}$ presents *target domain* instances. We explore $\rho \in \{0.0, 0.5, 1.0, 2.0, 5.0\}$ and finetune each source model for 50 epochs with learning rate as 0.01. For knowledge transferability, we random initialize the last layer of each source model and finetune it on STL-10 for 10 epochs with learning rate as 0.01. During the adversarial example generation, we utilize standard ℓ_∞ PGD attack with perturbation scale $\epsilon = 0.1$ and 50 attack iterations with step size as $\epsilon/10$.

Table 7 shows the relationship between knowledge transferability and adversarial transferability of different source model trained by different ρ . With the increasing of ρ , the adversarial transferability between source model and target model decreases ($\alpha_1, \alpha_1 * \alpha_2$ become smaller), and the knowledge transferability also decreases. We also plot the α_1 with its corresponding transfer loss on each instance, as shown in Figure 7. The negative correlation between α_1 and transfer loss confirms our theoretical insights.

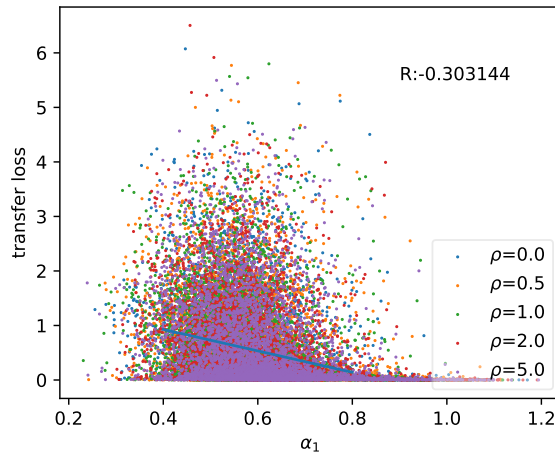


Figure 7. Distribution of per-sample knowledge transfer loss and α_1 . The Pearson score shows negative correlation between α_1 and transfer loss. The higher the loss is, the lower the knowledge transferability is, the lower the α_1 should be.

G.3. Knowledge Transferability Indicating Adversarial Transferability

G.3.1. IMAGE

We follow the same setup in the previous image experiment for source model training, transfer learning as well as generation of adversarial examples. However, there is one key difference: Instead of generating adversarial examples on the target model and measuring adversarial transferability on source models, we generate adversarial examples on each source model and measure the adversarial transferability by feeding these adversarial examples to the target model.

Similarly, we also visualize the results (Table 3) and compute the Pearson score. Due to the significant noise introduced by per-sample calculation, the R score is not as significant as figure 5, but the trend is still correct and valid, which shows that higher knowledge transferability indicates higher adversarial transferability.

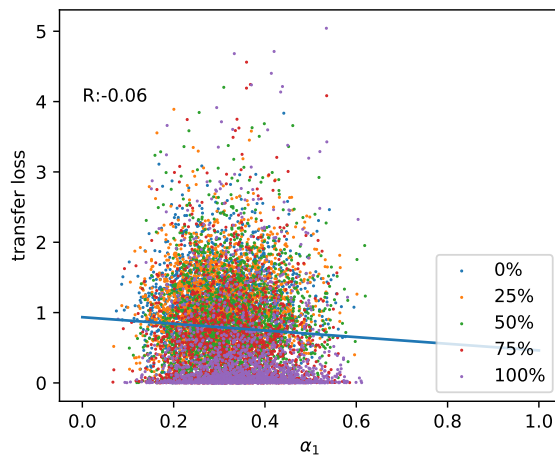


Figure 8. Distribution of per-sample knowledge transfer loss and α . The Pearson score shows negative strong correlation between α and transfer loss. The higher the loss is, the lower the knowledge transferability is, and the lower the α_1 is.

G.3.2. NLP

We follow the same setup to train the models and generate textual adversarial examples as §G.2 in the NLP experiments. We note that to measure the adversarial transferability, we generate 1,000 adversarial examples on each source model based on the test data from the target model, and measure the adversarial transferability by feeding these adversarial examples to the target model.

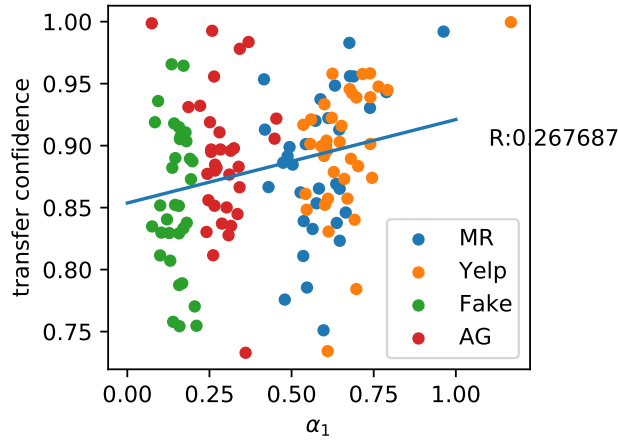


Figure 9. Distribution of per-batch knowledge transfer confidence and α_1 . The Pearson score shows positive correlation between α_1 and transfer confidence. The higher the confidence, the higher the knowledge transferability.