

Supplementary Materials

The Attack Generator: A Systematic Approach Towards Constructing Adversarial Attacks

In this document, we provide additional observations and results in support of the primary manuscript. The supplementary materials are divided into two parts: Appendix A is a comprehensive table which further illustrates the applicability of the attack generator concept. Here, we analyze existing adversarial attacks with respect to their building blocks, in particular their specificity, perturbation scope and perturbation imperceptibility measures. Appendix B shows further results of the two new semantic segmentation attacks. The first attack tries to enlarge the pedestrian class (confusion) while being imperceptible with respect to the L^∞ -norm. The second attack erases the pedestrian class (dynamic target) by utilizing flow field perturbations (attention-based imperceptibility).

Due to the size of the table, Appendix A starts on the following page.

A. Analysis Existing Attacks

Attack Name	Opt. Problem	Opt. Method	Specificity	Scope	Imperceptibility
FGSM [?]	$\min_{\delta} -l(F(x + \delta), y^x)$ $s.t. \ \delta\ _\infty \leq \epsilon,$ with y^x true label of image x and $l(\cdot, \cdot)$ loss function of neural network.	<i>First Order Method</i> Use scaled sign of gradient. $-l(F(x + \delta), y^x)$	<i>Untargeted</i> Measure: $M_{sp}(\delta)$	<i>Individual Scope</i> Only considering single image. Measure: $M_{sc}(\delta)$	$M_{min}(\delta)$ <i>L^p-based Imperceptibility Measure:</i> $\ \delta\ _\infty$
L-BFGS [?]	$\min_{\delta} l(F(x + \delta), y_{tar}^x) + \ \delta\ _2$ $s.t. x + \delta \in [0, 1]^m$ with y_{tar}^x adversarial target.	<i>Second Order Method</i> Box-constrained L-BFGS is a quasi-Newton method.	<i>Targeted</i> y_{tar}^x is desired outcome. Measure: $l(F(x + \delta), y_{tar}^x)$	<i>Individual Scope</i> Only considering single image. Measure: $M_{sp}(\delta)$	<i>L^p-based Imperceptibility Measure:</i> $\ \delta\ _2$
PGD [?]	$\min_{\delta} -l(F(x + \delta), y^x)$ $s.t. \ \delta\ _\infty \leq \epsilon,$ with y^x true label of x and $l(\cdot, \cdot)$ loss function of neural network.	<i>First Order Method</i> Multi-step gradient descent combined with projections.	<i>Untargeted</i> Measure: $-l(F(x + \delta), y^x)$	<i>Individual Scope</i> Only considering single image. Measure: $M_{sp}(\delta)$	<i>L^p-based Imperceptibility Measure:</i> $\ \delta\ _\infty$
Boundary Attack [?]	$\min_{\delta} \ \delta\ _2^2$ $s.t. c(\delta) = 1$ with $c(\cdot) \in \{0, 1\}$ adversarial criterion for data point x .	<i>Evolution & Random Sampling</i> Method is initialized from an adversarial point and then random walk along decision boundary.	<i>Untargeted / Targeted</i> Depends on choice of adversarial criterion (and initialization point). Measure: $c(x + \delta)$	<i>Individual Scope</i> Only considering single image. Measure: $M_{sp}(\delta)$	<i>L^p-based Imperceptibility Measure:</i> $\ \delta\ _2^2$
EOT [?]	$\min_{\delta} -\mathbb{E}_{t \sim T} [\log P(y_{tar}^x \mid t(x + \delta))]$ $s.t. \mathbb{E}_{t \sim T} [d(t(x), t(x + \delta))] \leq \epsilon,$ $x \in [0, 1]^m$ with distribution T of transformation functions, distance function $d(\cdot, \cdot)$ and y_{tar}^x target.	<i>First Order Method</i> Use projected gradient descent.	<i>Targeted</i> Measure also takes transformation t as input. Measure: $-\log P(y_{tar}^x \mid t(x + \delta))$	<i>Contextual</i> Consider one image x and relevant transformations of x . Measure: $\mathbb{E}_{t \sim T} [\mathcal{M}_{sp}(t, \delta)]$	<i>Attention-based / L^p-based Imperceptibility Measure:</i> Depends on choice of distance. Measure: $d(\cdot, \cdot)$

Table 1. Building block analysis of existing adversarial attacks.

B. Sample Results

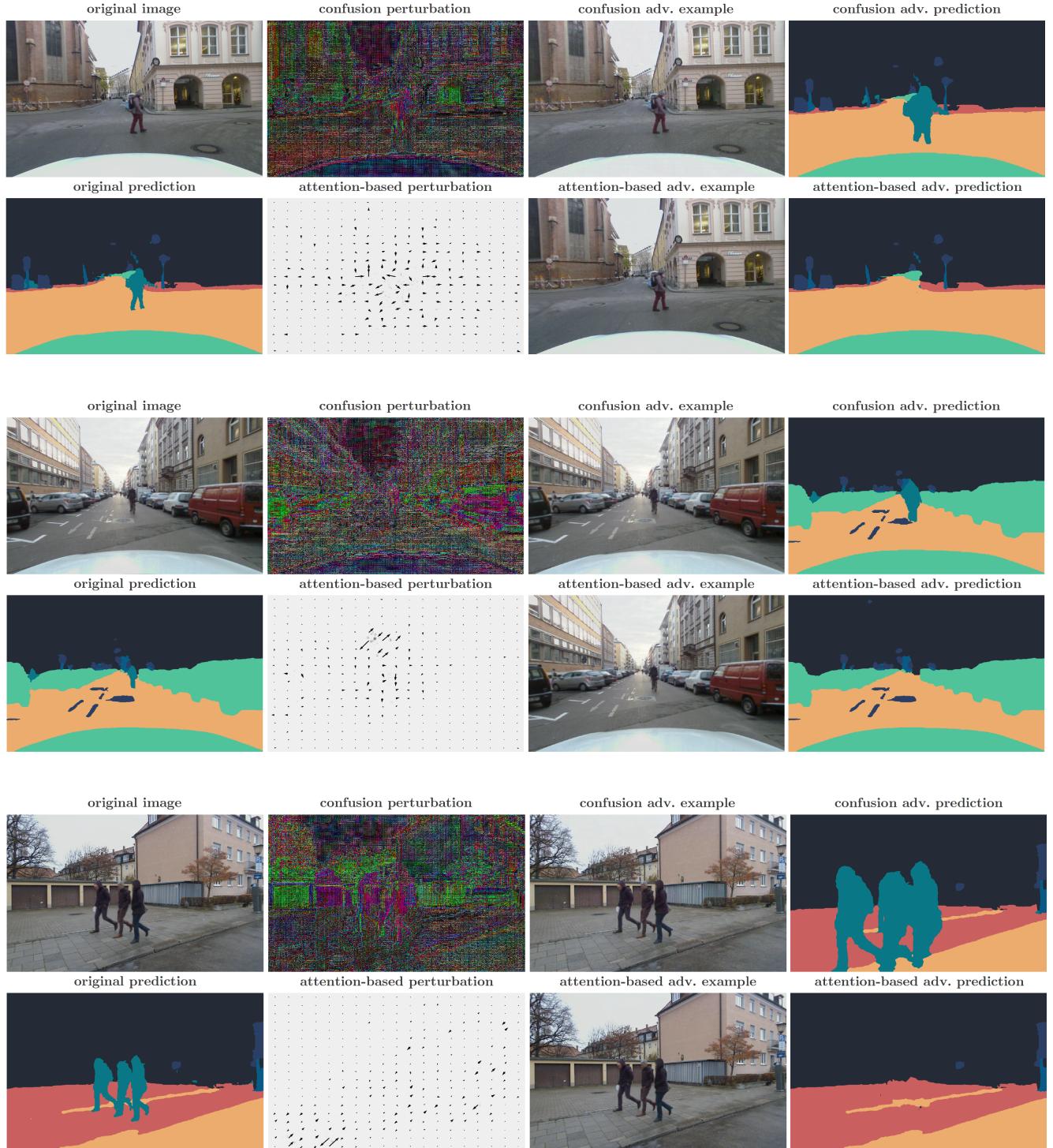


Figure 1. Sample results of adapted semantic segmentation attacks with confusion target or attention-based imperceptibility. Attack input is a self-trained ICNet and a single image ($N=1$, $\epsilon = 15$).