

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

Towards Analyzing Semantic Robustness of Deep Neural Networks

- Supplementary Materials -

Anonymous ICCV submission

Paper ID 152

1. Detailed Derivations of the Update Directions of the Bounds

In our case we consider a more general case where we are interested in the $\mathbf{u} \in \Omega \subset \mathbb{R}^n$, a hidden latent parameter that generates the image and is passed to scene generator (e.g. a renderer function \mathbf{R}) that takes the parameter \mathbf{u} and an object shape \mathbf{S} of a class that is identified by classifier \mathbf{C} . Ω is the continuous semantic space for the parameters that we intend to study. The renderer creates the image $\mathbf{x} \in \mathbb{R}^d$, and then we study the behaviour of a classifier \mathbf{C} of that image across multiple shapes and multiple famous DNNs. Now this function of interest is defined as follows.

$$f(\mathbf{u}) = \mathbf{C}_z(\mathbf{R}(\mathbf{S}_z, \mathbf{u})) , \quad 0 \leq f(\mathbf{u}) \leq 1 \quad (1)$$

where z is class label of interest of study and we observe the network score for that class by rendering a shape \mathbf{S}_z of the same class. The shape and class labels are constants and only the parameters varies for f . The robust-region-finding operator is then defined as follows

$$\begin{aligned} \Phi_{\text{robust}}(f(\mathbf{u}), \mathbf{S}_z, \mathbf{u}_0) &= \mathbb{D} = \{\mathbf{u} : \mathbf{a} \leq \mathbf{u} \leq \mathbf{b}\} \\ \text{s.t. } \mathbb{E}_{\mathbf{u} \sim \mathbb{D}}[f(\mathbf{u})] &\geq 1 - \epsilon_m , \quad \mathbf{u}_0 \in \mathbb{D} , \quad \text{VAR}[f(\mathbf{u})] \leq \epsilon_v \end{aligned} \quad (2)$$

where the left and right bounds of \mathbb{D} are $\mathbf{a} = [a_1, a_2, \dots, a_n]$ and $\mathbf{b} = [b_1, b_2, \dots, b_n]$ respectively. The two small thresholds ϵ_m, ϵ_v are to insure high performance and low variance of the DNN network in that robust region. We can define the opposite operator which is to find adversarial regions like follows :

$$\begin{aligned} \Phi_{\text{adv}}(f(\mathbf{u}), \mathbf{S}_z, \mathbf{u}_0) &= \mathbb{D} = \{\mathbf{u} : \mathbf{a} \leq \mathbf{u} \leq \mathbf{b}\} \\ \text{s.t. } \mathbb{E}_{\mathbf{u} \sim \mathbb{D}}[f(\mathbf{u})] &\leq \epsilon_m , \quad \mathbf{u}_0 \in \mathbb{D} , \quad \text{VAR}[f(\mathbf{u})] \geq \epsilon_v \end{aligned} \quad (3)$$

We can show clearly that Φ_{adv} and Φ_{robust} are related as follows

$$\Phi_{\text{adv}}(f(\mathbf{u}), \mathbf{S}_z, \mathbf{u}_0) = \Phi_{\text{robust}}(1 - f(\mathbf{u}), \mathbf{S}_z, \mathbf{u}_0) \quad (4)$$

So we can just focus our attentions on Φ_{robust} to find robust regions , and the adversarial regions follows directly from Eq (4).

1.1. Divergence of the Bounds

To develop an algorithm for Φ , we deploy the idea by [7] which focus on maximizing the inner area of the function in the region and fitting the bounds to grow the region bounds. As we will show , maximizing the region by maximizing the integral can lead to divergence , as follows :

Lemma 1.1. *Let f be a continuous scalar function $f : \mathbb{R}^1 \rightarrow \mathbb{R}^1$, and let L be the function defining the definite integral of f in terms of the two integral bounds , i.e. $L(a, b) = \int_a^b f(u)du$. Then, to maximize L , $-f(a)$ and $f(b)$ are valid ascent directions for the two bounds a, b respectively.*

Proof. a direction \mathbf{p} is an ascent direction of objective L if it satisfies the inequality $\mathbf{p}^T \nabla L \geq 0$.[2].

To find $\frac{\partial L}{\partial a} = \frac{\partial}{\partial a} \int_a^b f(u)du$, we use Leibniz rule from the fundamental theorem of calculus which states that $\frac{d}{dx} \int_{a(x)}^{b(x)} f(u)du = f(b(x)) \frac{d}{dx} b(x) - f(a(x)) \frac{d}{dx} a(x)$

Therefore, $\frac{\partial}{\partial a} \int_a^b f(u)du = f(b) \times 0 - f(a) \times 1 = -f(a)$. Similarly, $\frac{\partial}{\partial b} \int_a^b f(u)du = f(b)$. By picking $\mathbf{p} = [-f(a), f(b)]^T$, then $\mathbf{p}^T \nabla L = f(a)^2 + f(b)^2 \geq 0$. This proves that \mathbf{p} is a valid ascent direction for objective L . \square

Theorem 1.2. *Let f be a positive continuous scalar function $f : \mathbb{R}^1 \rightarrow (0, 1)$, and let L be the function defining the definite integral of f in terms of the two integral bounds , i.e. $L(a, b) = \int_a^b f(u)du$. Then, following the ascent direction in Lemma 1.1 can diverge the bounds if followed in a gradient ascent technique with fixed learning rate .*

Proof. If we follow the direction $= [-f(a), f(b)]^T$, with a fixed learning rate η , then the update rules for a, b will be as follows. $a_k = a_{k-1} - \eta f(a)$, $b_k = b_{k-1} + \eta f(b)$. for initial points a_0, b_0 , then $a_k = a_0 - \eta \sum_{i=0}^k f(\text{Area}_{\text{in}})$, and $b_k = b_0 + \eta \sum_{i=0}^k f(b_i)$. We can see now if $f(u) = c$, $0 < c < 1$, then as $k \rightarrow \infty$, the bounds $a_k \rightarrow -\infty$, $b_k \rightarrow \infty$. This leads to the claimed divergence. \square

To solve the issue of bounds diverging we propose the following formulations for one dimensional bounds , and

108 then we extend them to n-dimensions , which allows for
 109 finding the n-dimensional semantic robust/adversarial re-
 110 gions that are similar to figure 2 . some of the formulations
 111 are black-box in nature (they dint need the gradient of the
 112 function f in order to update the current estimates of the
 113 bound) while others .
 114

115 1.2. Naive Approach

$$116 \quad L = -\text{Area}_{\text{in}} + \frac{\lambda}{2} \|b - a\|_2^2 \quad (5)$$

120 using Libeniz rule as in Lemma 1.1, we get the following
 121 update steps for the objective L :
 122

$$\begin{aligned} 123 \quad \frac{\partial L}{\partial a} &= f(a) - \lambda(b - a) \\ 124 \quad \frac{\partial L}{\partial b} &= -f(b) + \lambda(b - a) \end{aligned} \quad (6)$$

128 where λ is regularizing the boundaries not to extend too
 129 much in the case the function evaluation wa positive all the
 130 time.

131 **Extension to n-dimension:** Lets start by $n = 2$. Now, we
 132 have $f : \mathbb{R}^2 \rightarrow (0, 1)$, then we define the loss integral as a
 133 function of four bounds of a rectangular region as follows.
 134

$$\begin{aligned} 135 \quad L(a_1, a_2, b_1, b_2) &= \\ 136 \quad &- \int_{a_2}^{b_2} \int_{a_1}^{b_1} f(u, v) dv du + \frac{\lambda}{2} \|b_1 - a_1\|_2^2 + \frac{\lambda}{2} \|b_2 - a_2\|_2^2 \end{aligned} \quad (7)$$

140 We apply the trapezoidal approximation on the loss to ob-
 141 tain the following expression.

$$\begin{aligned} 143 \quad L(a_1, a_2, b_1, b_2) &\approx \frac{\lambda}{2} \|b_1 - a_1\|_2^2 + \frac{\lambda}{2} \|b_2 - a_2\|_2^2 \\ 144 \quad &- \frac{(b_1 - a_1)(b_2 - a_2)}{4} (f(a_1, a_2) + f(b_1, a_2) + \\ 145 \quad &\quad f(a_1, b_2) + f(b_1, b_2)) \end{aligned} \quad (8)$$

149 to find the update direction for the first bound a_1 by tak-
 150 ing the partial derivative of the function in Eq (7) we get
 151 the following update direction for a_1 along with its trape-
 152 zoidal approximation in order to able to compute it during
 153 the optimization:

$$\begin{aligned} 155 \quad \frac{\partial L}{\partial a_1} &= \int_{a_2}^{b_2} f(a_1, v) dv - \lambda(b_1 - a_1) \\ 156 \quad &\approx \frac{(b_2 - a_2)}{2} (f(a_1, a_2) + f(a_1, b_2)) - \lambda(b_1 - a_1) \end{aligned} \quad (9)$$

160 Doing similar steps to the first bound for the other
 161 three bounds we obtain the full update directions for

$$\begin{aligned} 162 \quad (a_1, a_2, b_1, b_2) \\ 163 \quad \frac{\partial L}{\partial a_1} &\approx \frac{(b_2 - a_2)}{2} (f(a_1, a_2) + f(a_1, b_2)) - \lambda(b_1 - a_1) \\ 164 \quad \frac{\partial L}{\partial b_1} &\approx -\frac{(b_2 - a_2)}{2} (f(b_1, a_2) + f(b_1, b_2)) + \lambda(b_1 - a_1) \\ 165 \quad \frac{\partial L}{\partial a_2} &\approx \frac{(b_1 - a_1)}{2} (f(a_1, a_2) + f(b_1, a_2)) - \lambda(b_2 - a_2) \\ 166 \quad \frac{\partial L}{\partial b_2} &\approx -\frac{(b_1 - a_1)}{2} (f(a_1, b_2) + f(b_1, b_2)) + \lambda(b_2 - a_2) \end{aligned} \quad (10)$$

168 Now, for $f : \mathbb{R}^n \rightarrow (0, 1)$, we define the inner region hyper-
 169 rectangle as before $\mathbb{D} = \{\mathbf{x} : \mathbf{a} \leq \mathbf{x} \leq \mathbf{b}\}$.Here , we
 170 assume the size of the region is positive at every dimension
 171 , i.e. $\mathbf{r} = \mathbf{b} - \mathbf{a} > 0$. The volume of the region \mathbb{D} normalized
 172 by exponent of dimension n is expressed as follows
 173

$$\text{volume}(\mathbb{D}) = \Delta = \frac{1}{2^n} \prod_{i=1}^n \mathbf{r}_i \quad (11)$$

177 The region \mathbb{D} can also be defined in terms of the matrix \mathbf{D}
 178 of all the corner points $\{\mathbf{d}^i\}_{i=1}^{2^n}$ as follows.

$$\begin{aligned} 182 \quad \text{corners}(\mathbb{D}) &= \mathbf{D}_{n \times 2^n} = [\mathbf{d}^1 | \mathbf{d}^2 | \dots | \mathbf{d}^{2^n}] \\ 183 \quad \mathbf{D} &= \mathbf{1}^T \mathbf{a} + \mathbf{M}^T \odot (\mathbf{1}^T \mathbf{r}) \end{aligned} \quad (12)$$

185 where $\mathbf{1}$ is the all-ones vector of size 2^n , \odot is the
 186 Hadamard product of matrices (elemnt-wise) , and \mathbf{M} is
 187 a constant masking matrix defined as the matrix of binary
 188 numbers of n bits that range from 0 to $2n - 1$ deined as
 189 follows.

$$\mathbf{M}_{n \times 2^n} = [\mathbf{m}^0 | \mathbf{m}^1 | \dots | \mathbf{m}^{2^n-1}] , \text{ where } \mathbf{m}^i = \text{binary}_n(i) \quad (13)$$

190 We define the function vector as the vector $\mathbf{f}_{\mathbb{D}}$ of all function
 191 evaluations at all corner points of \mathbb{D}

$$\mathbf{f}_{\mathbb{D}} = [f(\mathbf{d}^1), f(\mathbf{d}^2), \dots, f(\mathbf{d}^{2^n})]^T, \mathbf{d}^i = \mathbf{D}_{:,i} \quad (14)$$

192 We follow similar steps as in $n = 2$ and obtain the following
 193 loss expressions and update directions :

$$\begin{aligned} 194 \quad L(\mathbf{a}, \mathbf{b}) &= - \int \dots \int_{\mathbb{D}} f(u_1, \dots, u_n) du_1 \dots du_n + \frac{\lambda}{2} |\mathbf{r}|^2 \\ 195 \quad &\approx -\Delta \mathbf{1}^T \mathbf{f}_{\mathbb{D}} + \frac{\lambda}{2} |\mathbf{r}|^2 \\ 196 \quad \nabla_{\mathbf{a}} L &\approx 2\Delta \text{diag}^{-1}(\mathbf{r}) \bar{\mathbf{M}} \mathbf{f}_{\mathbb{D}} + \lambda \mathbf{r} \\ 197 \quad \nabla_{\mathbf{b}} L &\approx -2\Delta \text{diag}^{-1}(\mathbf{r}) \mathbf{M} \mathbf{f}_{\mathbb{D}} - \lambda \mathbf{r} \end{aligned} \quad (15)$$

198 1.3. Outer-Inner Ratio Loss (OIR)

199 We introduce an outer region A, B with bigger area that
 200 contains the small region (a, b) . We follow the following

216 assumption to insure that outer area is always positive.
 217

$$218 \quad A = a - \alpha \frac{b-a}{2}, B = b + \alpha \frac{b-a}{2} \quad (16)$$

219 where α is the small boundary factor of the outer area to inner area. We formulate the problem as a ratio of outer over inner area and we try to make this ratio as close as possible to 0. $L = \frac{\text{Area}_{\text{out}}}{\text{Area}_{\text{in}}}$ By using DencklBeck technique for solving non-linear fractional programming problems [5]. Using their formulation to transform L as follows.

$$220 \quad L = \frac{\text{Area}_{\text{out}}}{\text{Area}_{\text{in}}} = \text{Area}_{\text{out}} - \lambda \text{Area}_{\text{in}} \\ 221 \quad = \int_A^B f(a)du - \int_a^b f(a)du - \lambda \int_a^b f(a)du \quad (17)$$

222 where $\lambda^* = \frac{\text{Area}_{\text{out}}^*}{\text{Area}_{\text{in}}^*}$ is the DencklBeck factor and it is equal to the small objective best achieved.

223 Black-Box (OIR_B).

224 Here we set $\lambda = 1$ to simplify the problem. This yields the following expression of the loss

$$225 \quad L = \text{Area}_{\text{out}} - \text{Area}_{\text{in}} \\ 226 \quad = \int_A^a f(u)du + \int_b^B f(u)du - \int_a^b f(u)du \\ 227 \quad = \int_A^B f(u)du - 2 \int_a^b f(u)du \quad (18) \\ 228 \quad = \int_{a-\alpha \frac{b-a}{2}}^{b+\alpha \frac{b-a}{2}} f(u)du - 2 \int_a^b f(u)du$$

229 using Libeniz ruloe as in Lemma 1.1, we get the following update steps for the objective L :

$$230 \quad \frac{\partial L}{\partial a} = -(1 + \frac{\alpha}{2})f(A) - \frac{\alpha}{2}f(B) + 2f(a) \quad (19) \\ 231 \quad \frac{\partial L}{\partial b} = (1 + \frac{\alpha}{2})f(B) + \frac{\alpha}{2}f(A) - 2f(b)$$

232 **Extension to n-dimension:** Lets start by $n = 2$. Now, we have $f : \mathbb{R}^2 \rightarrow (0, 1)$, and with the following constrains on the outer region.

$$233 \quad A_1 = a_1 - \frac{b_1 - a_1}{2}, \quad B_1 = b_1 + \frac{b_1 - a_1}{2} \quad (20) \\ 234 \quad A_2 = a_2 - \frac{b_2 - a_2}{2}, \quad B_1 = b_2 + \frac{b_2 - a_2}{2}$$

235 we define the loss integral as a function of four bounds of a rectangular region as follows.

$$236 \quad L(a_1, a_2, b_1, b_2) = \\ 237 \quad \int_{A_2}^{B_2} \int_{A_1}^{B_1} f(u, v)dvdu - 2 \int_{a_2}^{b_2} \int_{a_1}^{b_1} f(u, v)dvdu \quad (21)$$

238 We apply the trapezoidal approximation on the loss to obtain the following expression.

$$239 \quad L(a_1, a_2, b_1, b_2) \approx \\ 240 \quad \frac{(B_1 - A_1)(B_2 - A_2)}{4} (f(A_1, A_2) + f(B_1, A_2) + \\ 241 \quad f(A_1, B_2) + f(B_1, B_2)) \\ 242 \quad - \frac{(b_1 - a_1)(b_2 - a_2)}{2} (f(a_1, a_2) + f(b_1, a_2) + \\ 243 \quad f(a_1, b_2) + f(b_1, b_2)) \\ 244 \quad = \frac{(b_1 - a_1)(b_2 - a_2)}{4} (\\ 245 \quad (1 + \alpha)^2 (f(A_1, A_2) + f(B_1, A_2) + \\ 246 \quad f(A_1, B_2) + f(B_1, B_2)) \\ 247 \quad - 2(f(a_1, a_2) + f(b_1, a_2) + f(a_1, b_2) + f(b_1, b_2))) \quad (22)$$

248 to find the update direction for the first bound a_1 by taking the partial derivative of the function in Eq (21) we get the following update direction for a_1 along with its trapezoidal approximation in order to able to compute it during the optimization:

$$249 \quad \frac{\partial L}{\partial a_1} = -(1 + \frac{\alpha}{2}) \int_{A_2}^{B_2} \left(f(A_1, v) + \frac{\alpha}{2} f(B_1, v) \right) dv \\ 250 \quad + 2 \int_{a_2}^{b_2} f(a_1, v) dv \\ 251 \quad \approx \frac{(b_2 - a_2)}{2} (\\ 252 \quad - (1 + \alpha) \left((1 + \frac{\alpha}{2})(f(A_1, A_2) + f(A_1, B_2)) \right. \\ 253 \quad \left. + \frac{\alpha}{2} (f(B_1, A_2) + f(B_1, B_2)) \right) \\ 254 \quad + 2 \left(f(a_1, a_2) + f(a_1, b_2) \right) \quad (23)$$

255 Doing similar steps to the first bound for the other three bounds we obtain the full update directions for (a_1, a_2, b_1, b_2)

$$256 \quad \frac{\partial L}{\partial a_1} \approx \frac{(b_2 - a_2)}{2} (\\ 257 \quad - (1 + \alpha) \left((1 + \frac{\alpha}{2})(f(A_1, A_2) + f(A_1, B_2)) \right. \\ 258 \quad \left. + \frac{\alpha}{2} (f(B_1, A_2) + f(B_1, B_2)) \right) \\ 259 \quad + 2 \left(f(a_1, a_2) + f(a_1, b_2) \right) \quad (24)$$

$$\begin{aligned} \frac{\partial L}{\partial b_1} &\approx \frac{(b_2 - a_2)}{2} \left((1 + \alpha) \left((1 + \frac{\alpha}{2}) (f(B_1, A_2) + f(B_1, B_2)) \right) \right. \\ &\quad \left. + \frac{\alpha}{2} (f(A_1, A_2) + f(A_1, B_2)) \right. \\ &\quad \left. - 2(f(b_1, a_2) + f(b_1, b_2)) \right) \end{aligned} \quad (25)$$

$$\begin{aligned} \frac{\partial L}{\partial a_2} &\approx \frac{(b_1 - a_1)}{2} \left(-(1 + \alpha) \left((1 + \frac{\alpha}{2}) (f(A_1, A_2) + f(B_1, A_2)) \right) \right. \\ &\quad \left. + \frac{\alpha}{2} (f(A_1, B_2) + f(B_1, B_2)) \right) \\ &\quad + 2(f(a_1, a_2) + f(b_1, a_2)) \end{aligned} \quad (26)$$

$$\begin{aligned} \frac{\partial L}{\partial b_2} &\approx \frac{(b_1 - a_1)}{2} \left((1 + \alpha) \left((1 + \frac{\alpha}{2}) (f(A_1, B_2) + f(B_1, B_2)) \right) \right. \\ &\quad \left. + \frac{\alpha}{2} (f(A_1, A_2) + f(B_1, A_2)) \right) \\ &\quad - 2(f(a_1, b_2) + f(b_1, b_2)) \end{aligned} \quad (27)$$

Now, for $f : \mathbb{R}^n \rightarrow (0, 1)$, we define the inner region hyper-rectangle as before $\mathbb{D} = \{\mathbf{x} : \mathbf{a} \leq \mathbf{x} \leq \mathbf{b}\}$, but now define an outer bigger region \mathbb{Q} that include the smaller region \mathbb{D} and defined as follows : $\mathbb{Q} = \{\mathbf{x} : \mathbf{a} - \frac{\alpha}{2}\mathbf{r} \leq \mathbf{x} \leq \mathbf{b} + \frac{\alpha}{2}\mathbf{r}\}$, where $\mathbf{a}, \mathbf{b}, \mathbf{r}$ are defined as before, while α is defined as the boundary factor of the outer region in all the dimensions equivalently. The inner and outer regions can also be defined in terms of the corner points as follows.

$$\begin{aligned} \text{corners}(\mathbb{D}) &= \mathbf{D}_{n \times 2^n} = [\mathbf{d}^1 | \mathbf{d}^2 | \dots | \mathbf{d}^{2^n}] \\ \mathbf{D} &= \mathbf{1}^T \mathbf{a} + \mathbf{M}^T \odot (\mathbf{1}^T \mathbf{r}) \\ \text{corners}(\mathbb{Q}) &= \mathbf{Q}_{n \times 2^n} = [\mathbf{q}^1 | \mathbf{q}^2 | \dots | \mathbf{q}^{2^n}] \\ \mathbf{Q} &= \mathbf{1}^T (\mathbf{a} - \frac{\alpha}{2}\mathbf{r}) + (1 + \alpha)\mathbf{M}^T \odot (\mathbf{1}^T \mathbf{r}) \end{aligned} \quad (28)$$

where $\mathbf{1}$ is the all-ones vector of size 2^n , \odot is the Hadamard product of matrices (element-wise), and \mathbf{M} is a constant masking matrix defined in Eq (60). Now we define two function vectors evaluated at all possible inner and outer corner points respectively.

$$\begin{aligned} \mathbf{f}_{\mathbb{D}} &= [f(\mathbf{d}^1), f(\mathbf{d}^2), \dots, f(\mathbf{d}^{2^n})]^T, \mathbf{d}^i = \mathbf{D}_{:, i} \\ \mathbf{f}_{\mathbb{Q}} &= [f(\mathbf{q}^1), f(\mathbf{q}^2), \dots, f(\mathbf{q}^{2^n})]^T, \mathbf{c}^i = \mathbf{Q}_{:, i} \end{aligned} \quad (29)$$

Now the loss and update directions for the n-dimensional

case becomes as follows .

$$\begin{aligned} L(\mathbf{a}, \mathbf{b}) &= \int \dots \int_{\mathbb{Q}} f(u_1, \dots, u_n) du_1 \dots du_n \\ &\quad - 2 \int \dots \int_{\mathbb{D}} f(u_1, \dots, u_n) du_1 \dots du_n \\ &\approx \Delta \left((1 + \alpha)^n \mathbf{1}^T \mathbf{f}_{\mathbb{Q}} - 2 \mathbf{1}^T \mathbf{f}_{\mathbb{D}} \right) \\ \nabla_{\mathbf{a}} L &\approx 2 \Delta \text{diag}^{-1}(\mathbf{r}) \left(2 \bar{\mathbf{M}} \mathbf{f}_{\mathbb{D}} - \bar{\mathbf{M}}_{\mathbb{Q}} \mathbf{f}_{\mathbb{Q}} \right) \\ \nabla_{\mathbf{b}} L &\approx 2 \Delta \text{diag}^{-1}(\mathbf{r}) \left(-2 \mathbf{M} \mathbf{f}_{\mathbb{D}} + \mathbf{M}_{\mathbb{Q}} \mathbf{f}_{\mathbb{Q}} \right) \end{aligned} \quad (30)$$

wheere $\bar{\mathbf{M}}_{\mathbb{Q}}$ is the outer region mask defined as follows.

$$\begin{aligned} \bar{\mathbf{M}}_{\mathbb{Q}} &= (1 + \alpha)^{n-1} \left((1 + \frac{\alpha}{2}) \bar{\mathbf{M}} + \frac{\alpha}{2} \mathbf{M} \right) \\ \mathbf{M}_{\mathbb{Q}} &= (1 + \alpha)^{n-1} \left((1 + \frac{\alpha}{2}) \mathbf{M} + \frac{\alpha}{2} \bar{\mathbf{M}} \right) \end{aligned} \quad (31)$$

White-Box OIR (OIR_W.) The following formulation is white-box in nature (it need the gradient of the function f in order to update the current estimates of the bound) this is useful when the function in hand is differntiable (e.g. DNN), to obtain more intelligent regions, rather then the regions surrounded by near 0 values of the function f . We set $\lambda = \frac{\alpha}{\beta}$ in Eq (17), where α is the small boundary factor of the outer area, β is the emphasis factor (we will show later how it determines the emphasis on the function vs the gradient). Hence, the objective in Eq (17) becomes :

$$\begin{aligned} \arg \min_{a, b} L &= \arg \min_{a, b} \text{Area}_{\text{out}} - \lambda \text{Area}_{\text{in}} \\ &= \arg \min_{a, b} \int_A^a f(u) du + \int_b^B f(u) du - \frac{\alpha}{\beta} \int_a^b f(u) du \\ &= \arg \min_{a, b} \frac{\beta}{\alpha} \int_{a - \alpha \frac{b-a}{2}}^{b + \alpha \frac{b-a}{2}} f(u) du - (1 + \frac{\beta}{\alpha}) \int_a^b f(u) du \end{aligned} \quad (32)$$

using Libeniz ruloe as in Lemma 1.1, we get the following derivatives of the bound a :

$$\begin{aligned} \frac{\partial L}{\partial a} &= \frac{\beta}{\alpha} \left(f(a) - f \left(a - \alpha \frac{b-a}{2} \right) \right) \\ &\quad - \frac{\beta}{2} f \left(b + \alpha \frac{b-a}{2} \right) - \frac{\beta}{2} f \left(a - \alpha \frac{b-a}{2} \right) + f(a) \end{aligned} \quad (33)$$

now since λ^* should be small for the optimal objective , then as $\lambda \rightarrow 0$, $\alpha \rightarrow 0$ and hence the derivative in Eq (33) becomes the following.

$$\begin{aligned} \lim_{\alpha \rightarrow 0} \frac{\partial L}{\partial a} &= \lim_{\alpha \rightarrow 0} \beta \frac{\left(f(a) - f \left(a - \alpha \frac{b-a}{2} \right) \right)}{\alpha} \\ &\quad - \frac{\beta}{2} f(b) - \frac{\beta}{2} f(a) + f(a) \end{aligned} \quad (34)$$

432 We can see that the first term is proportional to the derivative
 433 of f at a , and hence the expression becomes :

$$\lim_{\alpha \rightarrow 0} \frac{\partial L}{\partial a} = \frac{\beta}{2} \left((b-a)f'(a) + f(b) \right) + \left(1 - \frac{\beta}{2} \right) f(a) \quad (35)$$

439 we can see that the update rule depends on the function
 440 value and the derivative of f at the boundary a with β
 441 controlling the dependence. Similarly for the boundary b we
 442 can see the following direction

$$\lim_{\alpha \rightarrow 0} \frac{\partial L}{\partial b} = \frac{\beta}{2} \left((b-a)f'(b) + f(a) \right) - \left(1 - \frac{\beta}{2} \right) f(b) \quad (36)$$

443 If $\beta \rightarrow 0$, the update directions in Eq (35,36) collapse to
 444 the unregularized naive update in Eq (6).

451 **Extension to n-dimension.** Lets start by $n = 2$. Now, we
 452 have $f : \mathbb{R}^2 \rightarrow (0, 1)$, and with the following constrains on
 453 the outer region.

$$\begin{aligned} A_1 &= a_1 - \frac{b_1 - a_1}{2}, \quad B_1 = b_1 + \frac{b_1 - a_1}{2} \\ A_2 &= a_2 - \frac{b_2 - a_2}{2}, \quad B_2 = b_2 + \frac{b_2 - a_2}{2} \end{aligned} \quad (37)$$

464 we follow similar approach as in Eq (32) to obtain the fol-
 465 lowing expression.

$$\begin{aligned} L(a_1, a_2, b_1, b_2) &= \\ &\frac{\beta}{\alpha} \int_{A_2}^{B_2} \int_{A_1}^{B_1} f(u, v) dv du - \left(1 + \frac{\beta}{\alpha} \right) \int_{a_2}^{b_2} \int_{a_1}^{b_1} f(u, v) dv du \end{aligned} \quad (38)$$

473 We apply the trapezoidal approximation on the loss to ob-
 474 tain the following approximation.

$$\begin{aligned} L(a_1, a_2, b_1, b_2) &\approx -1 + (1 + \alpha)^2 \\ &\frac{(f(A_1, A_2) + f(B_1, A_2) + f(A_1, B_2) + f(B_1, B_2))}{(f(a_1, a_2) + f(b_1, a_2) + f(a_1, b_2) + f(b_1, b_2))} \end{aligned} \quad (39)$$

482 to find the update direction for the first bound a_1 by taking
 483 the partial derivative of the function in Eq (21) we get the
 484 following update direction for a_1 along with its trapezoidal
 485 approximation in order to able to compute it during the op-

timization:

$$\begin{aligned} \frac{\partial L}{\partial a_1} &= -\frac{\beta}{\alpha} \left(1 + \frac{\alpha}{2} \right) \int_{A_2}^{B_2} \left(f(A_1, v) + \frac{\alpha}{2} f(B_1, v) \right) dv \\ &\quad + \left(1 + \frac{\beta}{\alpha} \right) \int_{a_2}^{b_2} f(a_1, v) dv \\ &\approx \frac{(b_2 - a_2)}{2} \left(-\left(1 + \alpha \right) \left(\frac{\beta}{\alpha} + \frac{\beta}{2} \right) (f(A_1, A_2) + f(A_1, B_2)) \right. \\ &\quad \left. + \frac{\beta}{2} (f(B_1, A_2) + f(B_1, B_2)) \right) \\ &\quad + \left(1 + \frac{\beta}{\alpha} \right) \left(f(a_1, a_2) + f(a_1, b_2) \right) \end{aligned} \quad (40)$$

grouping the terms which are divided by α together and then
 taking the limit of $\alpha \rightarrow \infty$ (as explained in the 1-d case),
 we get the following expressions.

$$\begin{aligned} \lim_{\alpha \rightarrow 0} \frac{\partial L}{\partial a_1} &\approx \frac{(b_2 - a_2)}{2} \left(-\lim_{\alpha \rightarrow 0} \frac{\beta}{\alpha} ((f(a_1, a_2) + f(a_1, b_2)) - (f(A_1, A_2) + f(A_1, B_2))) \right. \\ &\quad \left. - \lim_{\alpha \rightarrow 0} \frac{3\beta}{2} (f(A_1, A_2) + f(A_1, B_2)) \right. \\ &\quad \left. - \lim_{\alpha \rightarrow 0} \frac{\beta}{2} (f(B_1, A_2) + f(B_1, B_2)) \right. \\ &\quad \left. + (f(a_1, a_2) + f(a_1, b_2)) \right) \end{aligned} \quad (41)$$

Noting that the first term is related to the directional derivatives of f , we get the following limit expression

$$\begin{aligned} \lim_{\alpha \rightarrow 0} \frac{\partial L}{\partial a_1} &\approx \frac{(b_2 - a_2)}{2} \left(\beta \left(\nabla f(a_1, a_2) \cdot \left(\frac{b_1 - a_1}{b_2 - a_2} \right) + \nabla f(a_1, b_2) \cdot \left(\frac{b_1 - a_1}{-b_2 + a_2} \right) \right) \right. \\ &\quad \left. + \left(1 - \frac{3\beta}{2} \right) (f(a_1, a_2) + f(a_1, b_2)) \right) \\ &\quad - \frac{\beta}{2} (f(b_1, a_2) + f(b_1, b_2)) \end{aligned} \quad (42)$$

Doing similar steps to the first bound for the other three bounds we obtain the full update directions for (a_1, a_2, b_1, b_2)

$$\begin{aligned} \lim_{\alpha \rightarrow 0} \frac{\partial L}{\partial a_1} &\approx \frac{(b_2 - a_2)}{2} \left(\beta \left(\nabla f(a_1, a_2) \cdot \left(\frac{(b_1 - a_1)}{(b_2 - a_2)} \right) + \nabla f(a_1, b_2) \cdot \left(\frac{(b_1 - a_1)}{-(b_2 - a_2)} \right) \right) \right. \\ &\quad \left. + \left(1 - \frac{3\beta}{2} \right) (f(a_1, a_2) + f(a_1, b_2)) \right) \\ &\quad - \frac{\beta}{2} (f(b_1, a_2) + f(b_1, b_2)) \end{aligned} \quad (43)$$

540 $\lim_{\alpha \rightarrow 0} \frac{\partial L}{\partial b_1} \approx \frac{(b_2 - a_2)}{2} ($
 541 $\beta (\nabla f(b_1, a_2) \cdot \binom{(b_1 - a_1)}{-(b_2 - a_2)} + \nabla f(b_1, b_2) \cdot \binom{(b_1 - a_1)}{(b_2 - a_2)}))$
 542 $- \left(1 - \frac{3\beta}{2}\right) (f(b_1, a_2) + f(b_1, b_2)))$
 543 $+ \frac{\beta}{2} (f(a_1, a_2) + f(a_1, b_2)))$
 544 $)$
 545 (44)

546 $k \in \{1, 2, \dots, n\}.$

547 $\mathbf{s}_k = \frac{1}{\mathbf{r}_k} \sum_{i=1, i \neq k}^n \mathbf{r}_i ((\bar{\mathbf{M}}_{i,:} - \mathbf{M}_{i,:}) \odot \bar{\mathbf{M}}_{k,:}) \mathbf{G}_{:,i}$
 548 $\bar{\mathbf{s}}_k = \frac{1}{\mathbf{r}_k} \sum_{i=1, i \neq k}^n \mathbf{r}_i ((\mathbf{M}_{i,:} - \bar{\mathbf{M}}_{i,:}) \odot \mathbf{M}_{k,:}) \mathbf{G}_{:,i}$
 549 $k \in \{1, 2, \dots, n\}$
 550 (49)

550 $\lim_{\alpha \rightarrow 0} \frac{\partial L}{\partial a_2} \approx \frac{(b_2 - a_2)}{2} ($
 551 $\beta (\nabla f(a_1, a_2) \cdot \binom{(b_1 - a_1)}{(b_2 - a_2)} + \nabla f(b_1, a_2) \cdot \binom{-(b_1 - a_1)}{(b_2 - a_2)}))$
 552 $+ \left(1 - \frac{3\beta}{2}\right) (f(a_1, a_2) + f(b_1, a_2)))$
 553 $- \frac{\beta}{2} (f(a_1, b_2) + f(b_1, b_2)))$
 554 $)$
 555 (45)

556 $\lim_{\alpha \rightarrow 0} \frac{\partial L}{\partial b_2} \approx \frac{(b_2 - a_2)}{2} ($
 557 $\beta (\nabla f(a_1, b_2) \cdot \binom{-(b_1 - a_1)}{(b_2 - a_2)} + \nabla f(a_2, b_2) \cdot \binom{(b_1 - a_1)}{(b_2 - a_2)}))$
 558 $- \left(1 - \frac{3\beta}{2}\right) (f(a_1, b_2) + f(b_1, b_2)))$
 559 $+ \frac{\beta}{2} (f(a_1, a_2) + f(b_1, a_2)))$
 560 $)$
 561 (46)

562 Now, for $f : \mathbb{R}^n \rightarrow (0, 1)$ Following previous notations we
 563 have the following expressions for the loss and update di-
 564 rections for the bound

565 $L(\mathbf{a}, \mathbf{b}) \approx \frac{(1 + \alpha)^n \mathbf{1}^T \mathbf{f}_{\mathbb{Q}}}{\mathbf{1}^T \mathbf{f}_{\mathbb{D}}} - 1$
 566 $\nabla_{\mathbf{a}} L \approx \Delta (\text{diag}^{-1}(\mathbf{r}) \bar{\mathbf{M}}_{\mathbb{D}} \mathbf{f}_{\mathbb{D}} + \beta \text{diag}(\bar{\mathbf{M}} \mathbf{G}_{\mathbb{D}}) + \beta \bar{\mathbf{s}})$
 567 $\nabla_{\mathbf{b}} L \approx \Delta (-\text{diag}^{-1}(\mathbf{r}) \mathbf{M}_{\mathbb{D}} \mathbf{f}_{\mathbb{D}} + \beta \text{diag}(\mathbf{M} \mathbf{G}_{\mathbb{D}}) + \beta \mathbf{s})$
 568 $)$
 569 (47)

570 where the mask is the special mask

571 $\bar{\mathbf{M}}_{\mathbb{D}} = \left(\gamma_n \bar{\mathbf{M}} - \beta \mathbf{M} \right)$
 572 $\mathbf{M}_{\mathbb{D}} = \left(\gamma_n \mathbf{M} - \beta \bar{\mathbf{M}} \right)$
 573 $\gamma_n = 2 - \beta(2n - 1)$
 574 (48)

575 Where $\text{diag}(\cdot)$ is the diagonal matrix of the vector argument or the diagonal vector of the matrix argument. \mathbf{s} is a weighted sum of the gradient from other dominions ($i \neq k$) contributing to the update direction of dimension k , where

576 $k \in \{1, 2, \dots, n\}.$
 577 $\mathbf{s}_k = \frac{1}{\mathbf{r}_k} \sum_{i=1, i \neq k}^n \mathbf{r}_i ((\bar{\mathbf{M}}_{i,:} - \mathbf{M}_{i,:}) \odot \bar{\mathbf{M}}_{k,:}) \mathbf{G}_{:,i}$
 578 $\bar{\mathbf{s}}_k = \frac{1}{\mathbf{r}_k} \sum_{i=1, i \neq k}^n \mathbf{r}_i ((\mathbf{M}_{i,:} - \bar{\mathbf{M}}_{i,:}) \odot \mathbf{M}_{k,:}) \mathbf{G}_{:,i}$
 579 $k \in \{1, 2, \dots, n\}$
 580 (49)

1.4. Trapezoidal Approximation Formulation

581 Here we use the trapezoidal approximation of the integral, a first-order approximation from Newton-Cortes formulas for numerical integration [8]. The rule states that a definite integral can be approximated as follows:

582 $\int_a^b f(u) du \approx (b - a) \frac{f(a) + f(b)}{2}$
 583 (50)

584 asymptotic error estimate is given by
 585 $-\frac{(b-a)^2}{48} [f'(b) - f'(a)] + \mathcal{O}(0.125)$. So as long the derivatives are bounded by some lipschitz constant \mathbb{L} , then the error becomes bounded by the following
 586 $|error| \leq \mathbb{L}(b - a)^2$. The regularized loss of interest in Eq
 587 (5) becomes the following .

588 $L = -\text{Area}_{\text{in}} + \lambda |b - a|_2^2$
 589 $\approx -(b - a) \frac{f(a) + f(b)}{2} + \lambda |b - a|_2^2$
 590 (51)

591 taking the derivative of L approximation directly with re-
 592 spect to these bounds , yields the following update direc-
 593 tions which are different from the expressions in Eq (6)

594 $\frac{\partial L}{\partial a} = -\frac{b - a}{2} f'(a) + \frac{f(a) + f(b)}{2} - \lambda(b - a)$
 595 $\frac{\partial L}{\partial b} = -\frac{b - a}{2} f'(b) - \frac{f(a) + f(b)}{2} + \lambda(b - a)$
 596 (52)

597 note that it needs the first derivative $f'(\cdot)$ of the function f
 598 evaluated at the bound to update that bound.

Extension to n-dimensions.

599 Lets start by $n = 2$. Now, we have $f : \mathbb{R}^2 \rightarrow (0, 1)$, and
 600 we define the loss integral as a function of four bounds of
 601 a rectangler region and apply the trapezoidal approximation

648

as follows.

$$\begin{aligned} L(a_1, a_2, b_1, b_2) &= - \int_{a_1}^{b_1} \int_{a_2}^{b_2} f(u, v) dv du \\ &\quad + \frac{\lambda}{2} |b_1 - a_1|_2^2 + \frac{\lambda}{2} |b_2 - a_2|_2^2 \\ &\approx \frac{\lambda}{2} |b_1 - a_1|_2^2 + \frac{\lambda}{2} |b_2 - a_2|_2^2 \\ &\quad - \frac{(b_1 - a_1)(b_2 - a_2)}{4} \left(f(a_1, a_2) + f(b_1, a_2) + \right. \\ &\quad \left. f(a_1, b_2) + f(b_1, b_2) \right) \end{aligned} \quad (53)$$

Then following similar steps as in the one-dimensional case we can obtain the following update directions for the four bounds

$$\begin{aligned} \frac{\partial L}{\partial a_1} &= (b_1 - a_1)(b_2 - a_2) \frac{f'_1(a_1, a_2) + f'_1(a_1, b_2)}{4} \\ &\quad - (b_2 - a_2) \frac{f(a_1, a_2) + f(a_1, b_2) + f(b_1, a_2) + f(b_1, b_2)}{4} \\ &\quad - \lambda(b_1 - a_1) \end{aligned} \quad (54)$$

$$\begin{aligned} \frac{\partial L}{\partial a_2} &= (b_1 - a_1)(b_2 - a_2) \frac{f'_2(b_1, a_2) + f'_2(b_1, b_2)}{4} \\ &\quad - (b_2 - a_2) \frac{f(a_1, a_2) + f(a_1, b_2) + f(b_1, a_2) + f(b_1, b_2)}{4} \\ &\quad - \lambda(b_2 - a_2) \end{aligned} \quad (55)$$

$$\begin{aligned} \frac{\partial L}{\partial b_1} &= (b_1 - a_1)(b_2 - a_2) \frac{f'_1(a_1, a_2) + f'_1(b_1, a_2)}{4} \\ &\quad + (b_1 - a_1) \frac{f(a_1, a_2) + f(a_1, b_2) + f(b_1, a_2) + f(b_1, b_2)}{4} \\ &\quad + \lambda(b_1 - a_1) \end{aligned} \quad (56)$$

$$\begin{aligned} \frac{\partial L}{\partial b_2} &= (b_1 - a_1)(b_2 - a_2) \frac{f'_2(a_1, b_2) + f'_2(b_1, b_2)}{4} \\ &\quad + (b_1 - a_1) \frac{f(a_1, a_2) + f(a_1, b_2) + f(b_1, a_2) + f(b_1, b_2)}{4} \\ &\quad + \lambda(b_2 - a_2) \end{aligned} \quad (57)$$

Where $f'_1(\cdot) = \frac{\partial f(u, v)}{\partial u}$, $f'_2(\cdot) = \frac{\partial f(u, v)}{\partial v}$. extending the 2-dimensional to general n-dimensions is straight forward. For $f : \mathbb{R}^n \rightarrow (0, 1)$, we define the following. Let the left bound vector be $\mathbf{a} = [a_1, a_2, \dots, a_n]$ and the right bound vector $\mathbf{b} = [b_1, b_2, \dots, b_n]$ define the n-dimensional hyper-rectangle region of interest. The region is then defined as follows : $\mathbb{D} = \{\mathbf{x} : \mathbf{a} \leq \mathbf{x} \leq \mathbf{b}\}$, Herer , we assume the size of the region is positve at ervy dimension , i.e. $\mathbf{r} = \mathbf{b} - \mathbf{a} > 0$. The volume of the region \mathbb{D} normalized by exponent of dimension n is expressed as follows

$$\text{volume}(\mathbb{D}) = \Delta = \frac{1}{2^n} \prod_{i=1}^n \mathbf{r}_i \quad (58)$$

The region \mathbb{D} can also be defined in terms of the matrix \mathbf{D} of all the corner points $\{\mathbf{d}^i\}_{i=1}^{2^n}$ as follows.

$$\begin{aligned} \text{corners}(\mathbb{D}) &= \mathbf{D}_{n \times 2^n} = [\mathbf{d}^1 | \mathbf{d}^2 | \dots | \mathbf{d}^{2^n}] \\ \mathbf{D} &= \mathbf{1}^T \mathbf{a} + \mathbf{M}^T \odot (\mathbf{1}^T \mathbf{r}) \end{aligned} \quad (59)$$

where $\mathbf{1}$ is the all-ones vector of size 2^n , \odot is the Hadamard product of matrices (elemnt-wise) , and \mathbf{M} is a constant masking matrix defined as the matrix of binary numbers of n bits that range from 0 to $2n - 1$

$$\mathbf{M}_{n \times 2^n} = [\mathbf{m}^0 | \mathbf{m}^1 | \dots | \mathbf{m}^{2^n-1}] , \text{ where } \mathbf{m}^i = \text{binary}_n(i) \quad (60)$$

We define the function vector as the vector $\mathbf{f}_{\mathbb{D}}$ of all function evaluations at all corner points of \mathbb{D}

$$\mathbf{f}_{\mathbb{D}} = [f(\mathbf{d}^1), f(\mathbf{d}^2), \dots, f(\mathbf{d}^{2^n})]^T, \mathbf{d}^i = \mathbf{D}_{:,i} \quad (61)$$

now we can see that the loss integral in Eq (53) becomes as follows .

$$\begin{aligned} L(\mathbf{a}, \mathbf{b}) &= \int \dots \int_{\mathbb{D}} f(u_1, \dots, u_n) du_1 \dots du_n + \frac{\lambda}{2} |\mathbf{r}|^2 \\ &\approx \Delta \mathbf{1}^T \mathbf{f}_{\mathbb{D}} + \frac{\lambda}{2} |\mathbf{r}|^2 \end{aligned} \quad (62)$$

Similarly to Eq (61), we define the gradient matrix $\mathbf{G}_{\mathbb{D}}$ as the matrix of all gradient vectors evaluated at all corner points of \mathbb{D}

$$\mathbf{G}_{\mathbb{D}} = [\nabla f(\mathbf{d}^1) | \nabla f(\mathbf{d}^2) | \dots | \nabla f(\mathbf{d}^{2^n})]^T \quad (63)$$

The update directions for the left bound \mathbf{a} and the right bound \mathbf{b} becomes as follows by the trapezoid approximation

$$\begin{aligned} \nabla_{\mathbf{a}} L &\approx \Delta \left(\text{diag}(\overline{\mathbf{M}} \mathbf{G}_{\mathbb{D}}) + \mathbf{1}^T \mathbf{f}_{\mathbb{D}} \text{diag}^{-1}(\mathbf{r}) \mathbf{1} \right) + \lambda \mathbf{r} \\ \nabla_{\mathbf{b}} L &\approx \Delta \left(\text{diag}(\mathbf{M} \mathbf{G}_{\mathbb{D}}) - \mathbf{1}^T \mathbf{f}_{\mathbb{D}} \text{diag}^{-1}(\mathbf{r}) \mathbf{1} \right) - \lambda \mathbf{r} \end{aligned} \quad (64)$$

Where $\overline{\mathbf{M}}$ is the complement of the binary mask matrix \mathbf{M} .

756	2. Analyzing Deep Neural Networks	810
757		811
758	Here we visualize different Network semantic maps generated during our analysis. In the 1D case we fix the elevation of the camera to a nominal angle of 35° and rotate around the object. In the 2D case, we change both the elevation and azimuth around the object. These maps can be generated to any type of semantic parameters that affect the generation of the image , and not viewing angle.	812
760		813
761		814
762		815
763		816
764		817
765		818
766		819
767	2.1. Networks Semantic Maps (1D)	820
768	In Figure 1,2 we visualize the 1D semantic maps of rotating around the object and recording different DNNs performance and averaging the profile over 10 different shapes per class.	821
769		822
770		823
771		824
772	2.2. Networks Semantic Maps (2D)	825
773	In Figure 3,4 we visualize the 2D semantic maps of elevation angles and rotating around the object and recording different DNNs performance and averaging the maps over 10 different shapes per class.	826
774		827
775		828
776		829
777		830
778	2.3. Convergence of the Region Finding Algorithms	831
779	Here we show how when we apply the region detection algorithms, the naive detect the smallest region while the OIR formulations detect bigger more comprehensive robust region. This happens even with different initial points, they always converge to the same bounds of that robust region of the semantic maps. Figure 5 show 4 different initializations for 1D case along with predicted regions. In Figure 6,7 shows the bounds evolving during the optimization of the three algorithms (naive , OIR_B and OIR_W) for 500 steps.	832
780		833
781		834
782		835
783		836
784		837
785		838
786		839
787		840
788		841
789		842
790	2.4. Examples of Found Regions (with Example Renderings)	843
791	In Figure 8,9 we provide examples of 2D regions found with the three algorithms along with renderings of the shapes from the robust regions detected.	844
792		845
793		846
794		847
795		848
796	2.5. Analyzing Semantic Data Bias in ImageNet	849
797		850
798	In Figure 10,11, we visualizing semantic data bias in the common training dataset (<i>i.e.</i> ImageNet [6]) by averaging the Networks Semantic Maps (NSM) of different networks and on different shapes, Different classes have different semantic bias in ImageNet as clearly shown in the maps above. These places of high confidence probably reveal the areas where abundance of examples exists in ImageNet, while holes conveys scarcity of such examples in ImageNet corresponding class.	851
799		852
800		853
801		854
802		855
803		856
804		857
805		858
806		859
807		860
808		861
809		862
		863

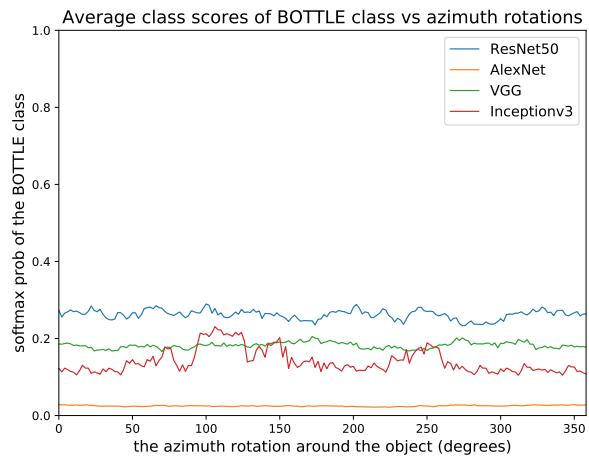
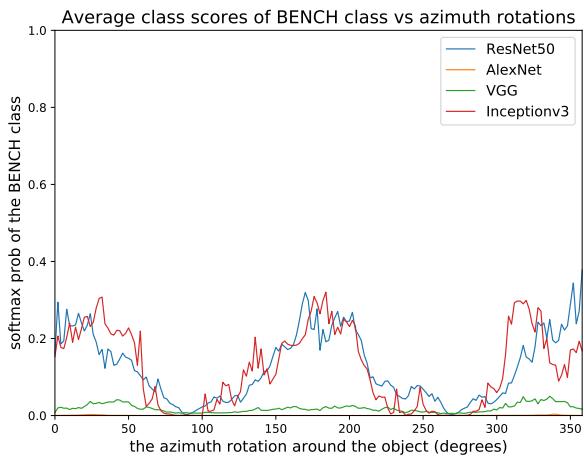
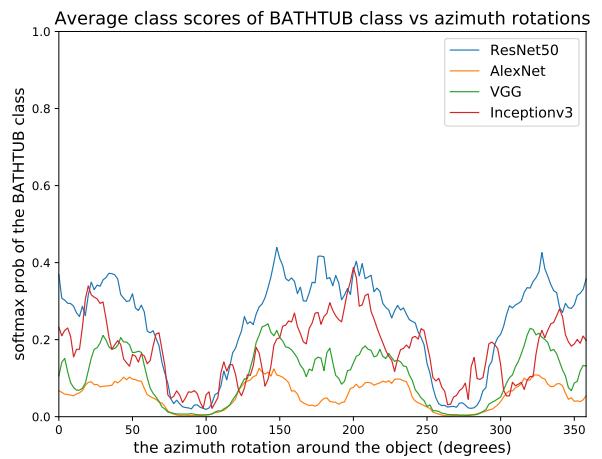
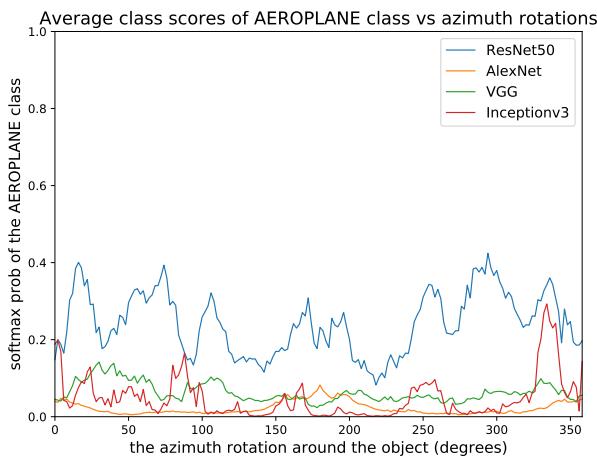


Figure 1: **1D Network Semantic Maps NMS-I.** Visualizing 1D Semantic Robustness profile for different networks averaged over 10 different shapes. Observe that different DNNs profiles differ depending on the training, accuracy, and network architectures that all result in a unique "signatures" for the DNN on that class. The correlation between the DNN profiles is due to the common data bias in ImageNet.

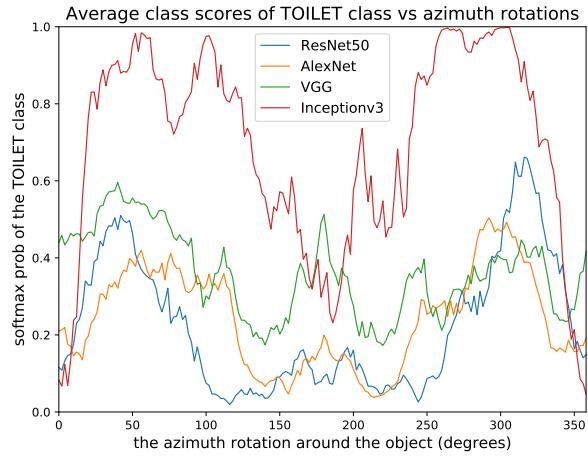
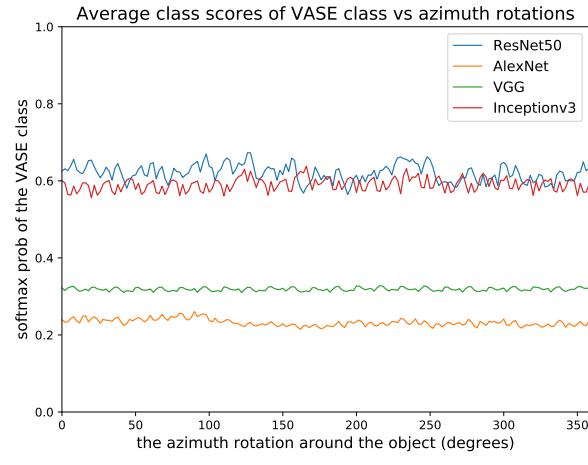
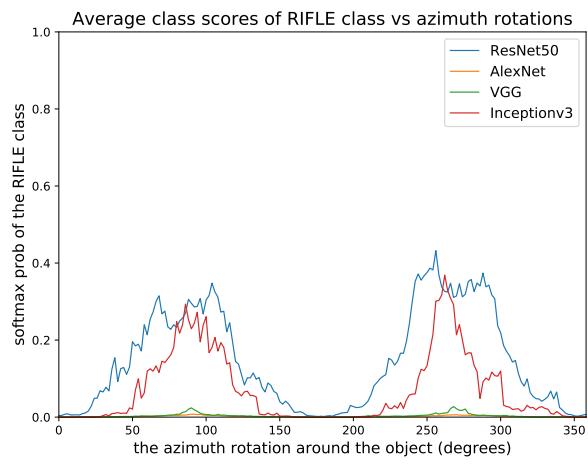
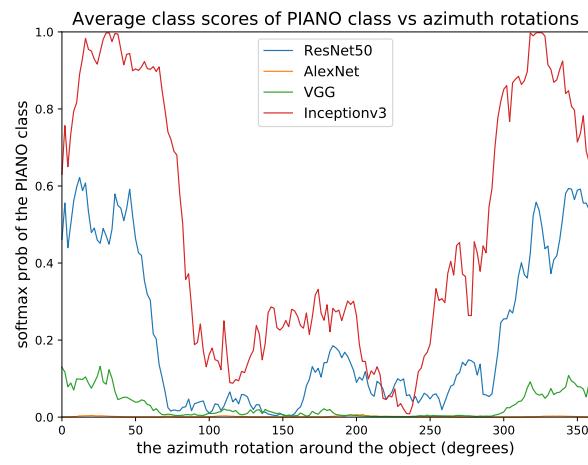
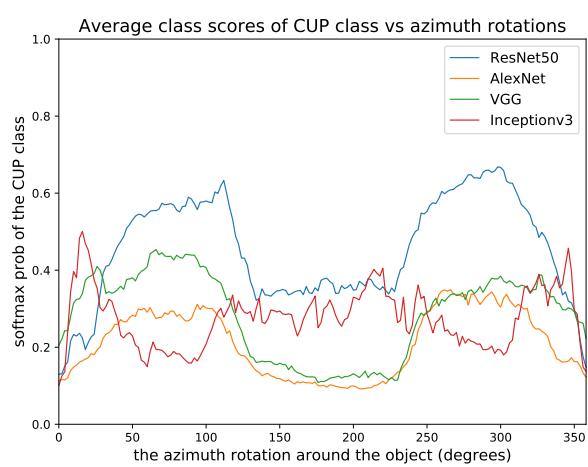
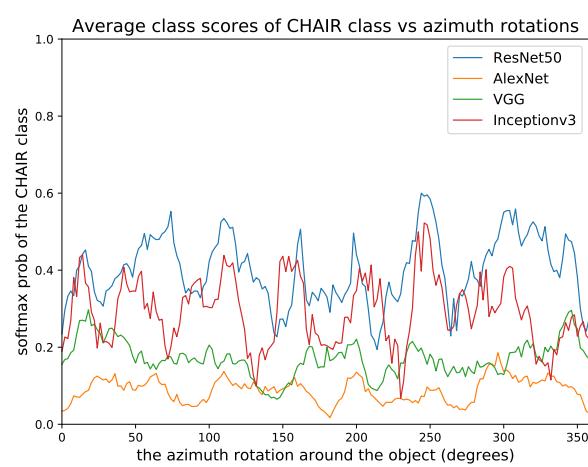


Figure 2: **1D Network Semantic Maps NMS-II**. Visualizing 1D Semantic Robustness profile for different networks averaged over 10 different shapes. Observe that different DNNs profiles differ depending on the training, accuracy, and network architectures that all result in a unique "signatures" for the DNN on that class. The correlation between the DNN profiles is due to the common data bias in ImageNet.

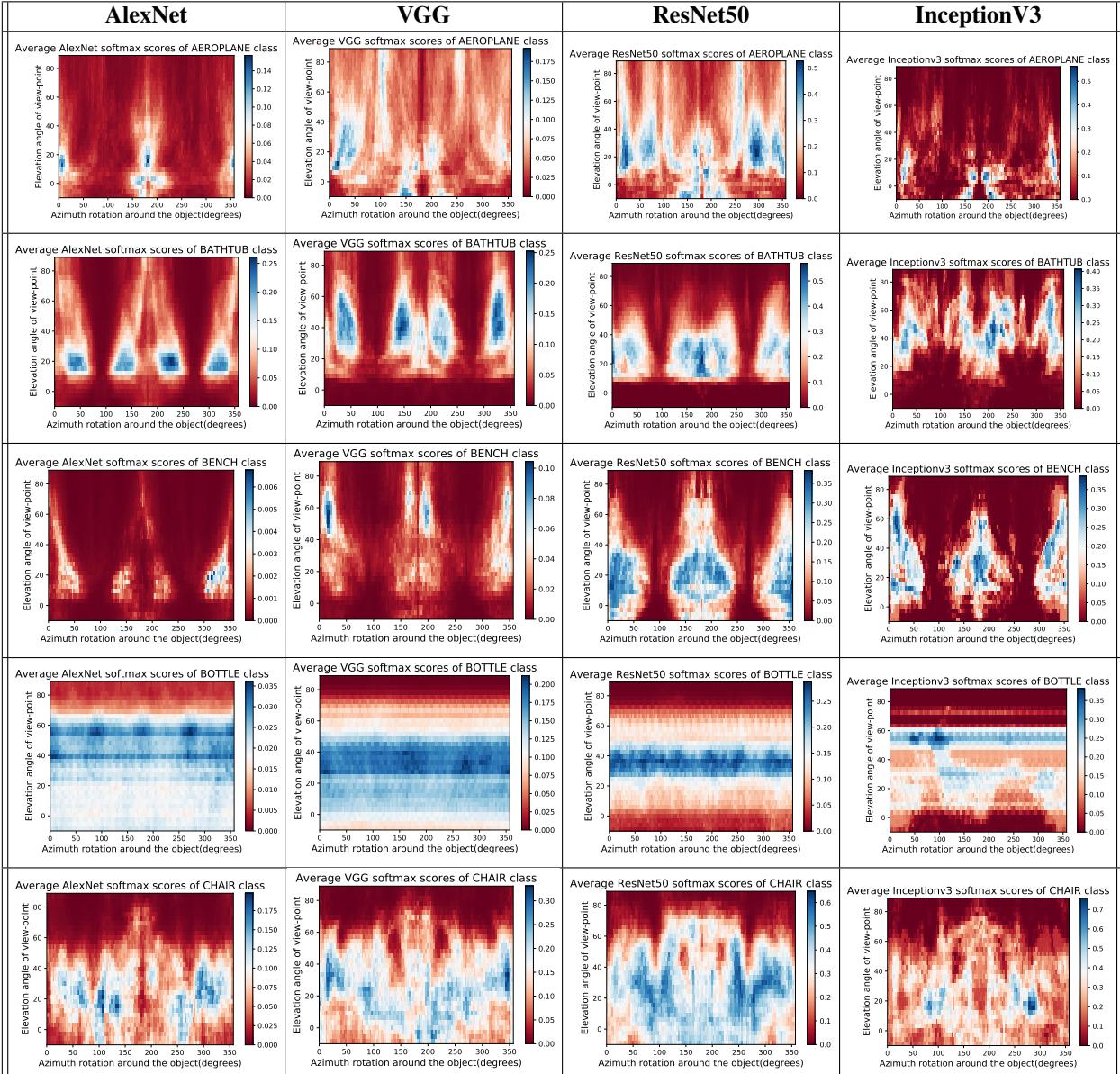
1080
1081
1082
1083
1084
10851134
1135
1136
1137
1138
11391086
10871140
1141

Figure 3: **2D Network Semantic Maps NMS-I.** Visualizing 2D Semantic Robustness profile for different networks averaged over 10 different shapes. Every row is different class. observe that differnt DNNs profiles differ depending on the training , accuracy , and network architectures that all result in a unique "signatures" for the DNN on that class. Also observe InceptionV3 [9] creates traps of very low score inside a region of high confidence

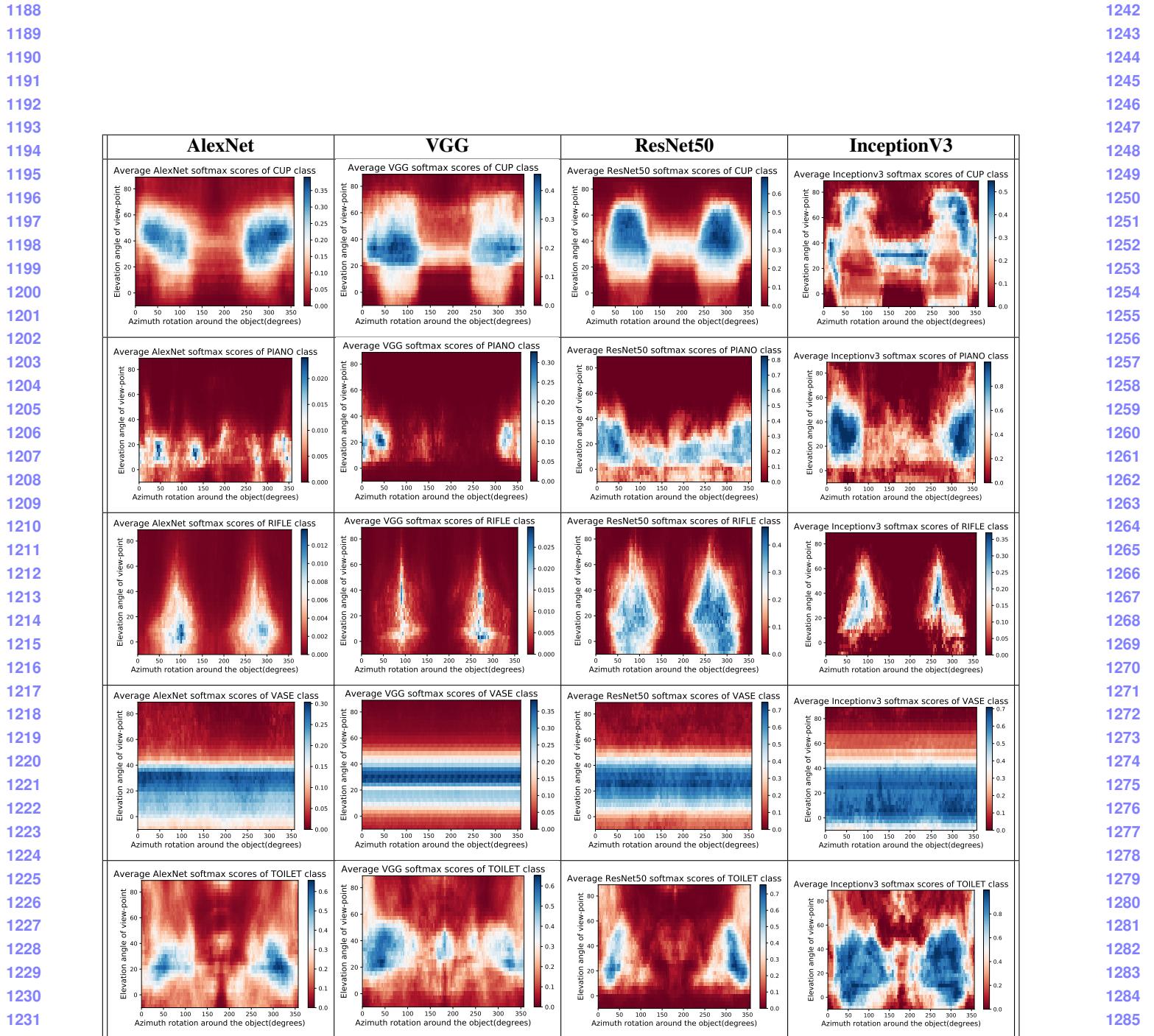


Figure 4: **2D Network Semantic Maps NMS-II**. Visualizing 2D Semantic Robustness profile for different networks averaged over 10 different shapes. Every row is different class. observe that differnt DNNs profiles differ depending on the training , accuracy , and network architectures that all result in a unique "signatures" for the DNN on that class. Also observe InceptionV3 [9] creates traps of very low score inside a region of high confidence

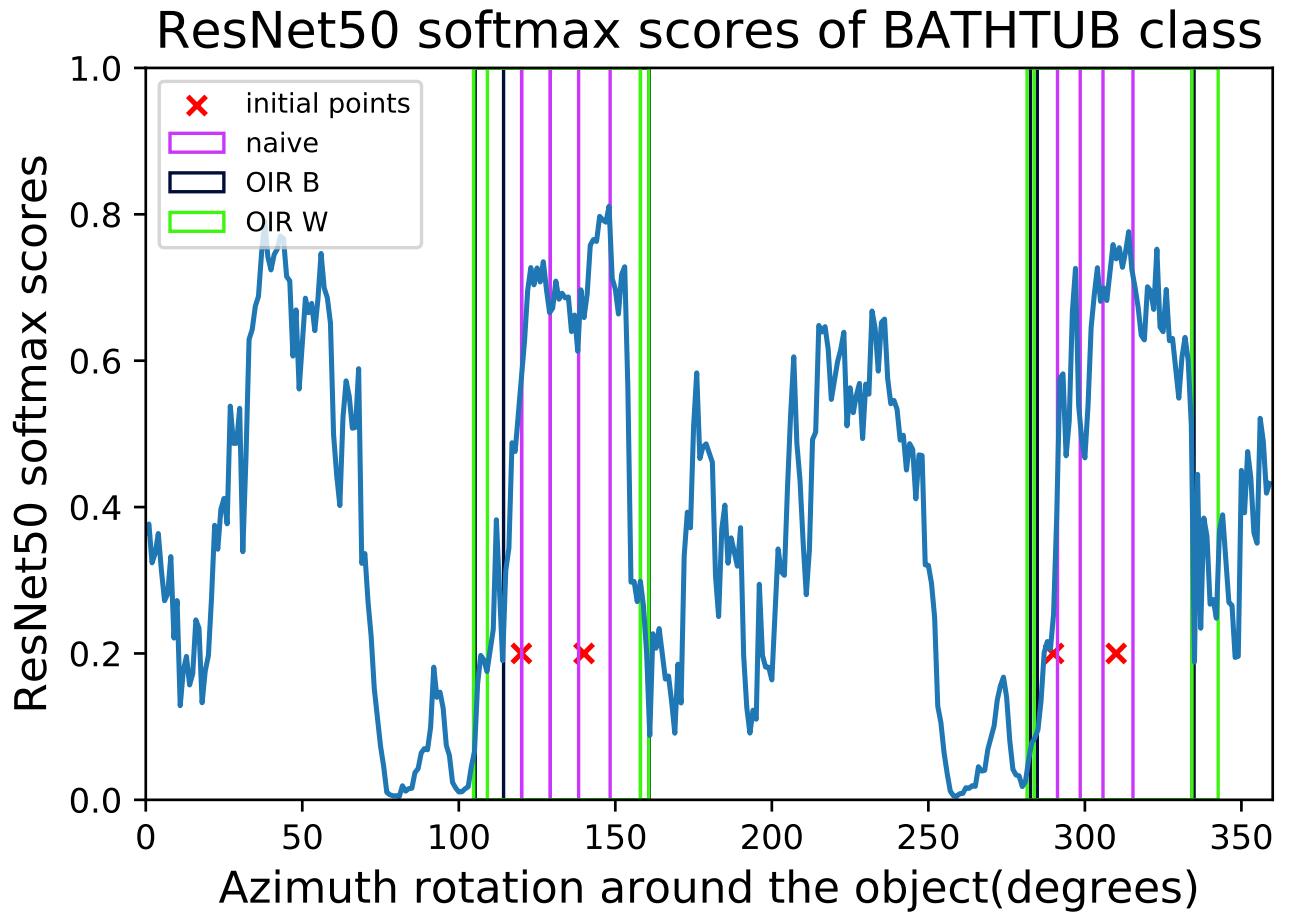


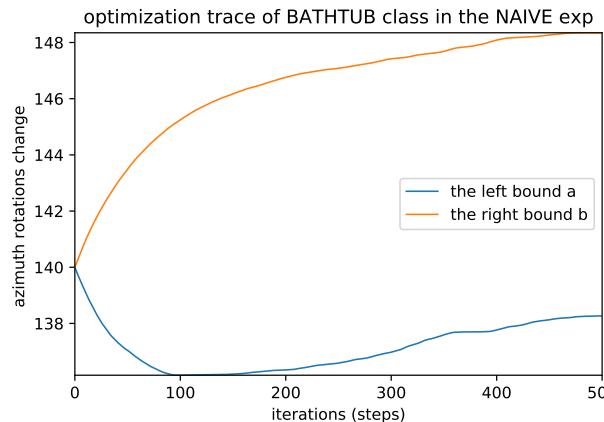
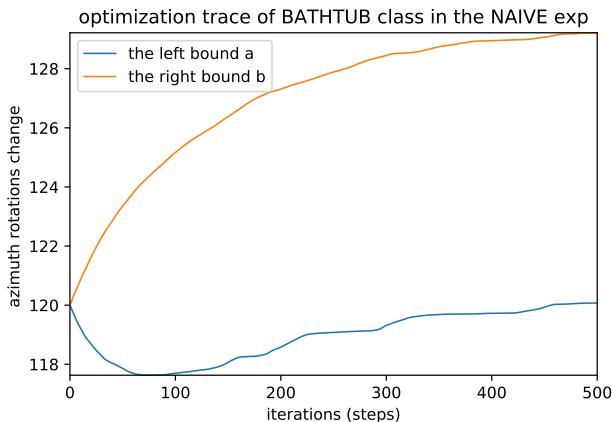
Figure 5: **Robust Region Detection with different initializations** Visualizing the Robust Regions Bounds found by the three algorithms for 4 initial points. we can see that the naive produce different bounds of the same region for different initializations while OIR detect the same region regardless of initialization.

1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
14211422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
14371438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
14571458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511

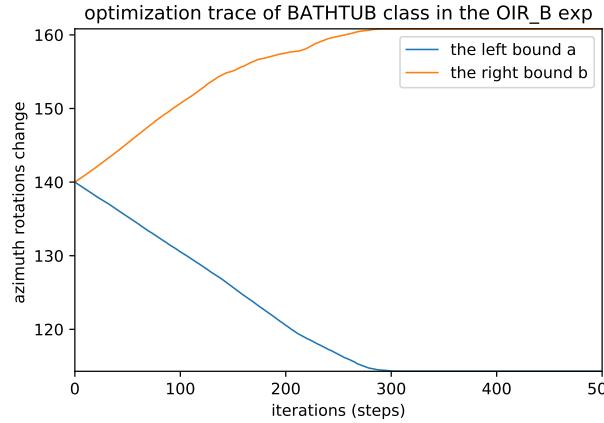
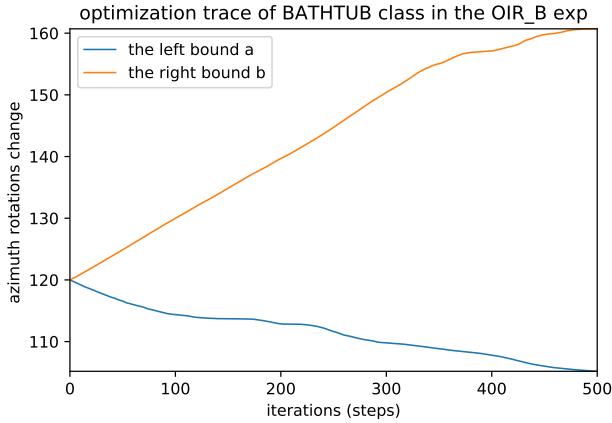
initialization = 120

initialization = 140

Naive algorithm



Black-Box OIR algorithm



White-Box OIR algorithm

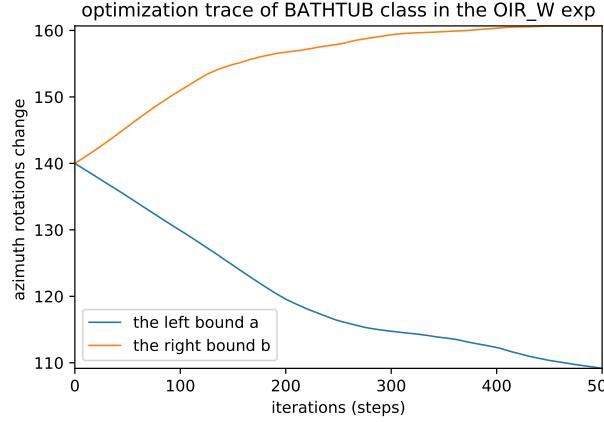
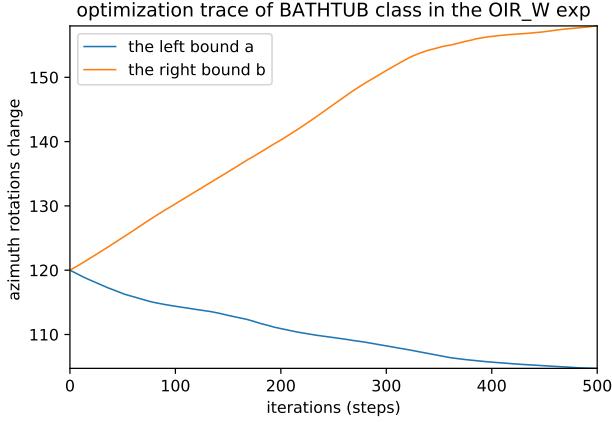


Figure 6: **Robust Region Bounds Growing I** Visualizing the bounds growing Using different algorithms from two different initializations (120 and 140) in Figure 5. We can see that OIR formulations converge to the same bounds of the robust region regardless of the initialization, which indicates effectiveness in detecting these regions, unlike the naive approach which can stop in a local optimum.

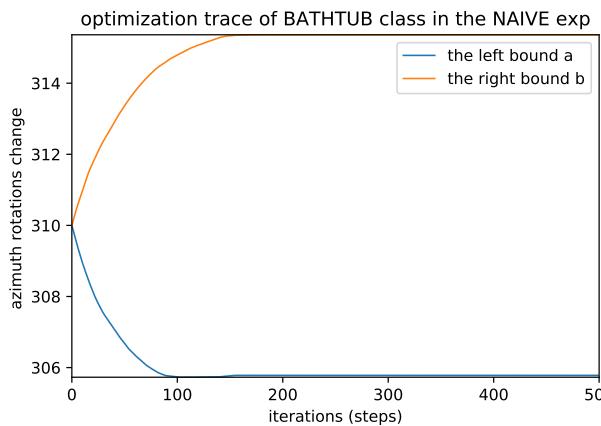
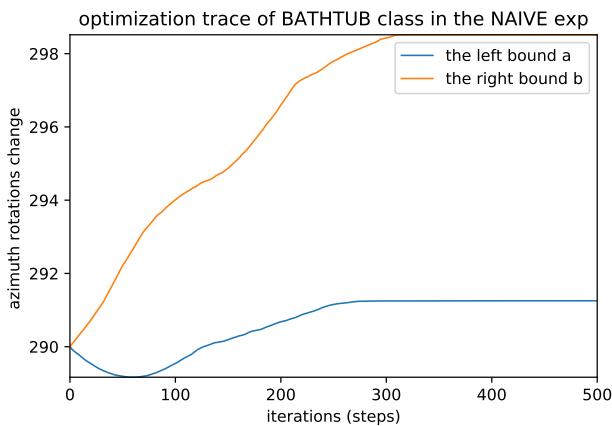
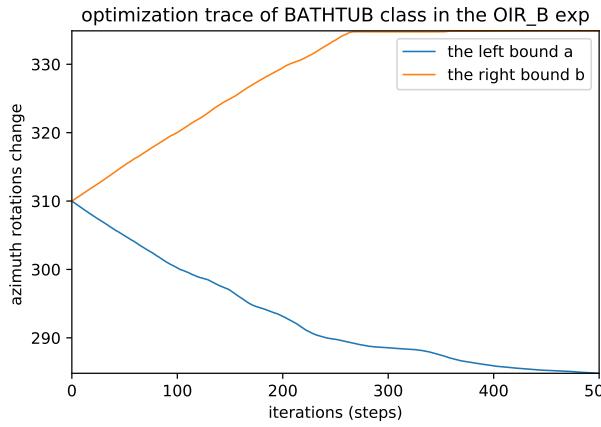
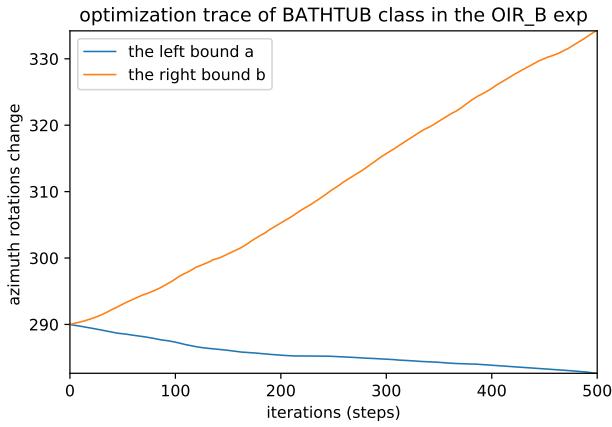
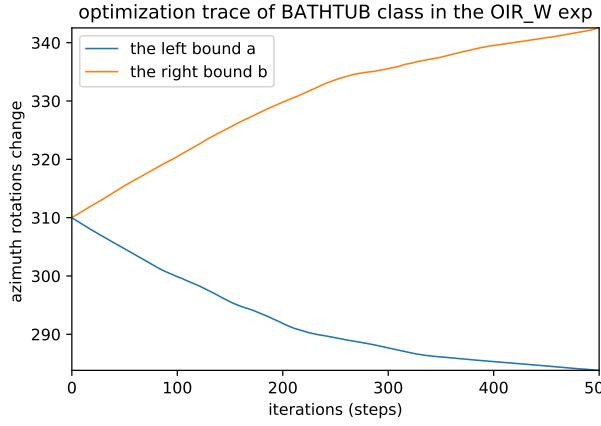
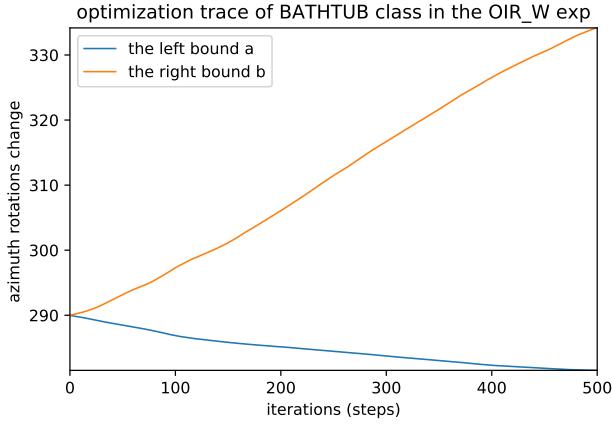
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
15291566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583**initialization = 290****initialization = 310****Naive algorithm****Black-Box OIR algorithm****White-Box OIR algorithm**

Figure 7: Robust Region Bounds Growing I Visualizing the bounds growing Using different algorithms from two different initializations (290, 310) in Figure 5. We can see that OIR formulations converge to the same bounds of the robust region regardless of the initialization, which indicates effectiveness in detecting these regions, unlike the naive approach which can stop in a local optimum.

1561

1615

1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673

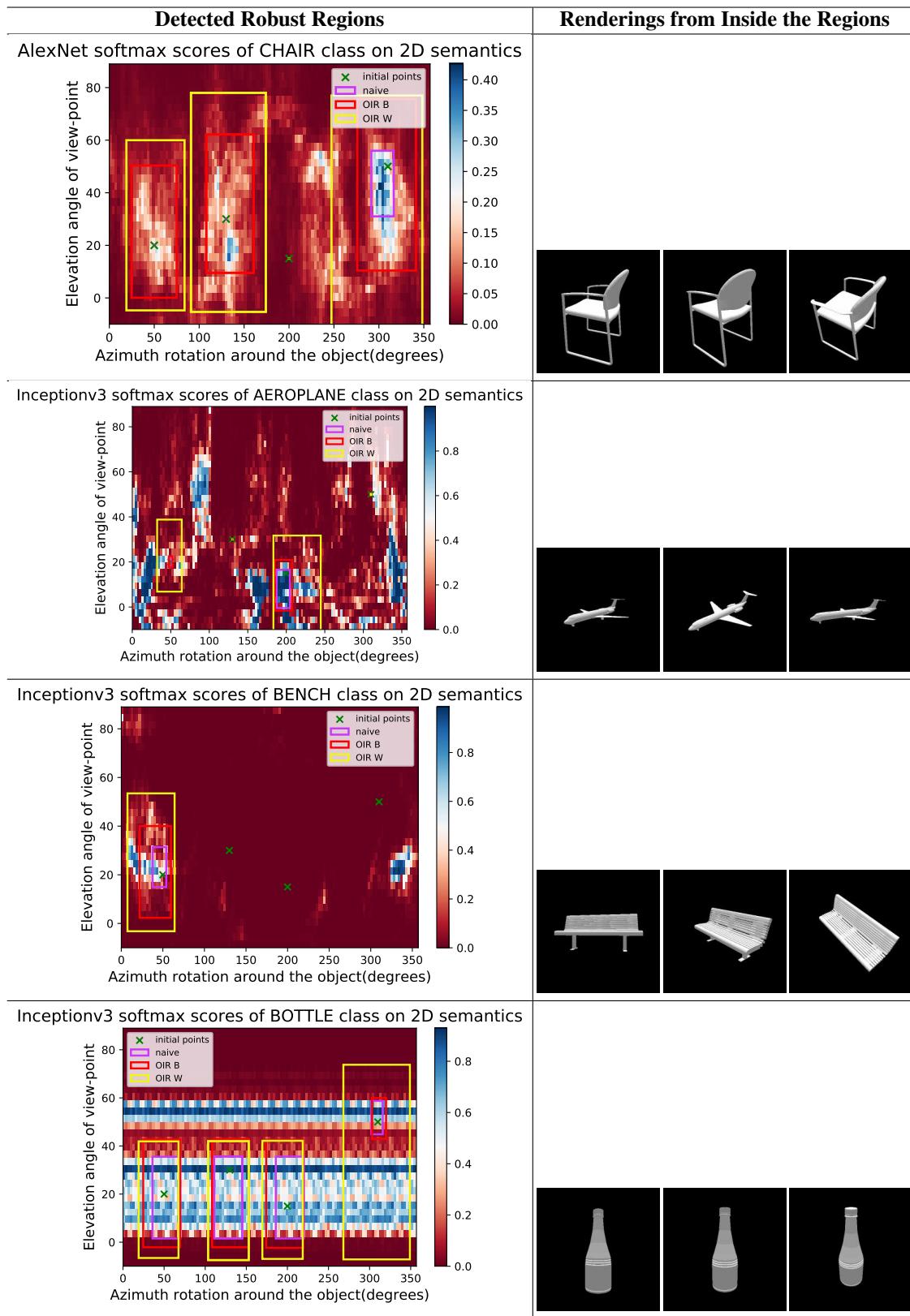


Figure 8: Qualitative Examples of Robust Regions I. Different runs of the algorithm to find robust regions along with different renderings from inside these regions for those specific shapes used in the experiments.

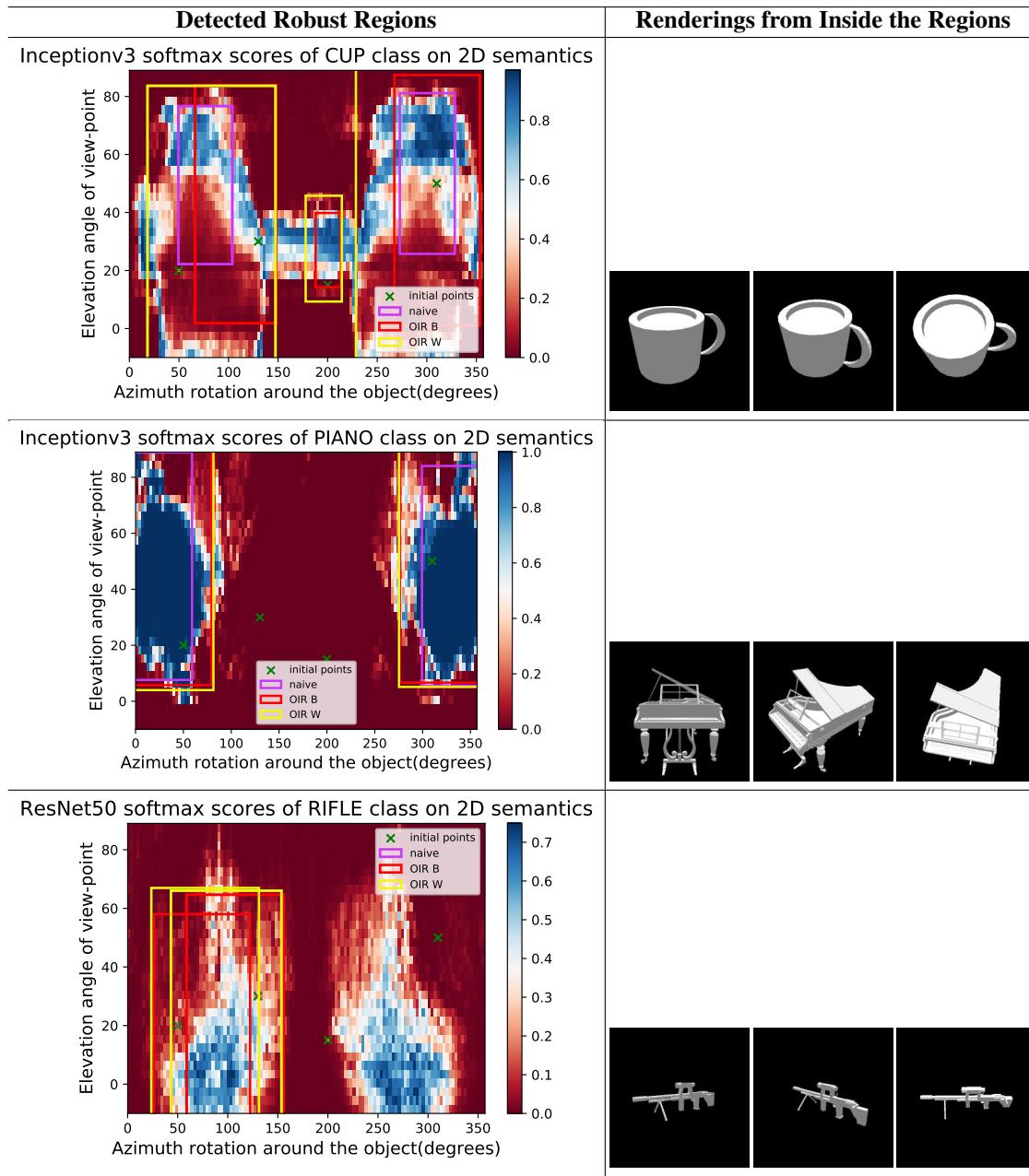
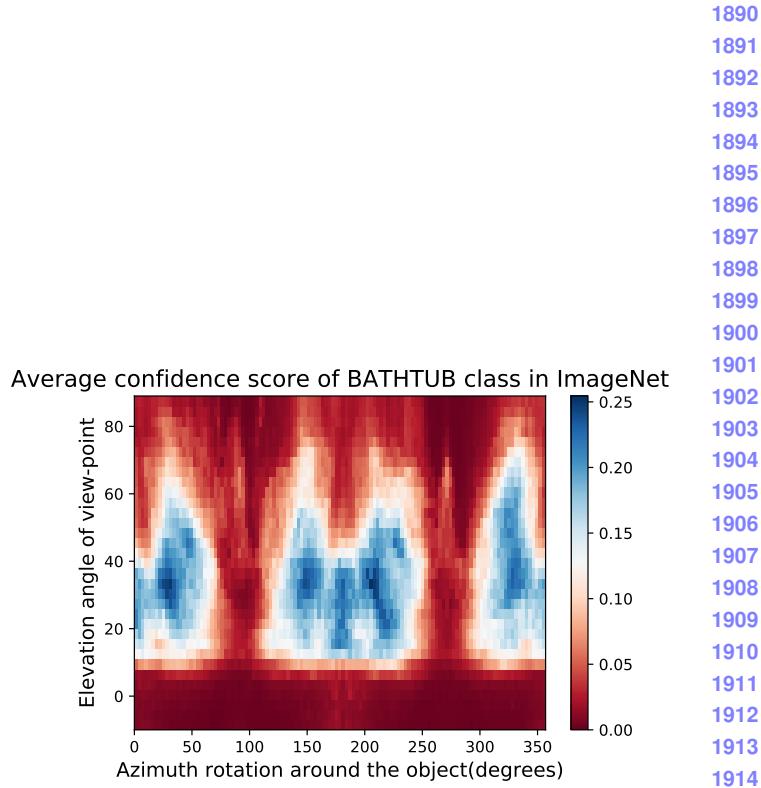
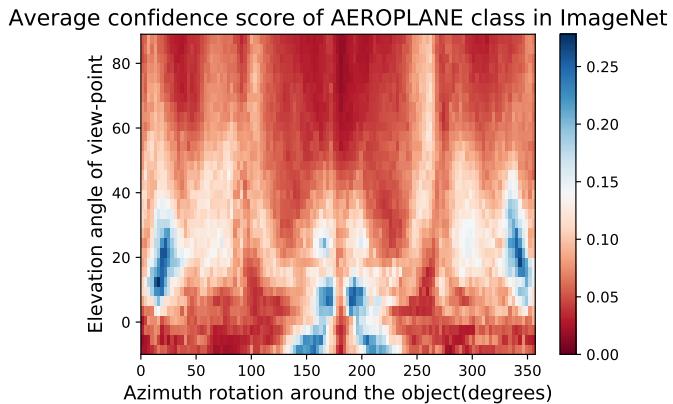
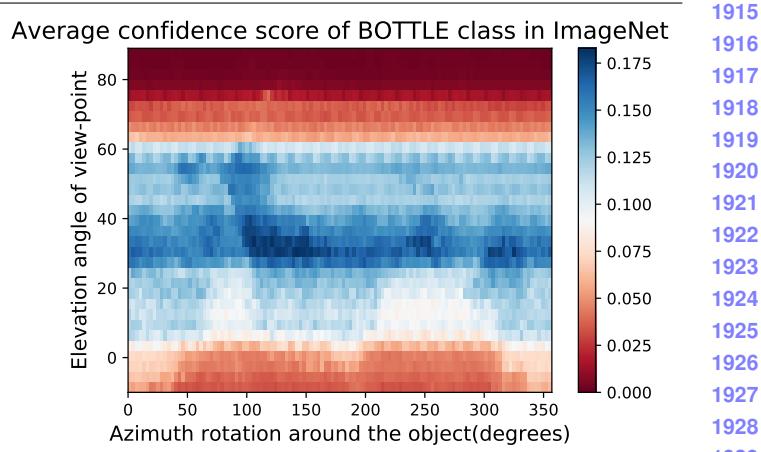
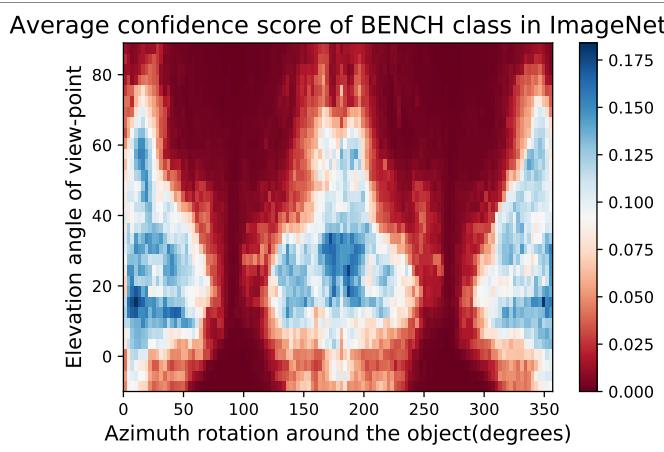


Figure 9: **Qualitative Examples of Robust Regions II** Different runs of the algorithm to find robust regions along with different renderings from inside these regions for those specific shapes used in the experiments.

1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860



1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875



1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889

Figure 10: **Data Semantic Maps DSM-I**. Visualizing Semantic Data Bias in the common training dataset (*i.e.* ImageNet [6]) by averaging the Networks Semantic Maps (NSM) of different networks and on different shapes, Different classes have different semantic bias in ImageNet as clearly shown in the maps above. The symmetry in the maps are attributed to the 3D symmetry of the objects.

1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943

1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961

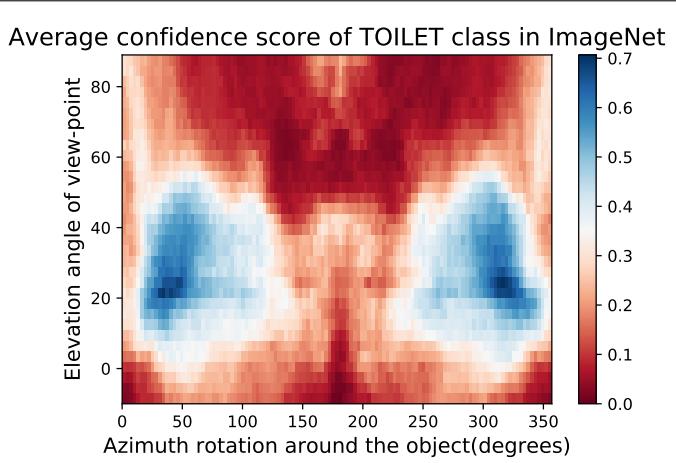
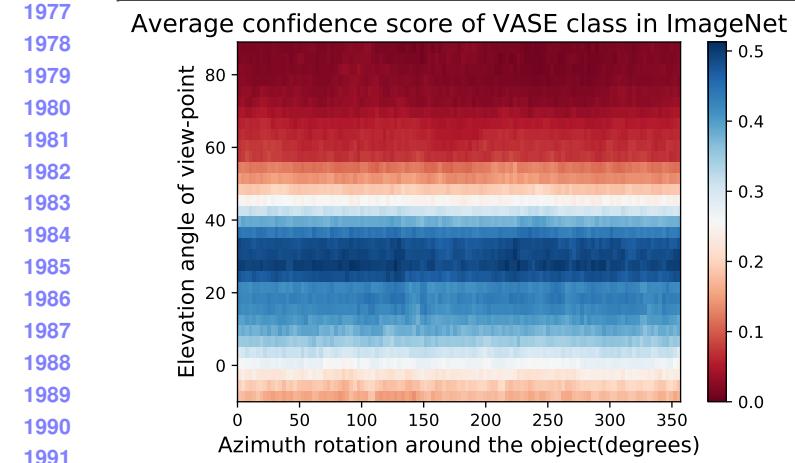
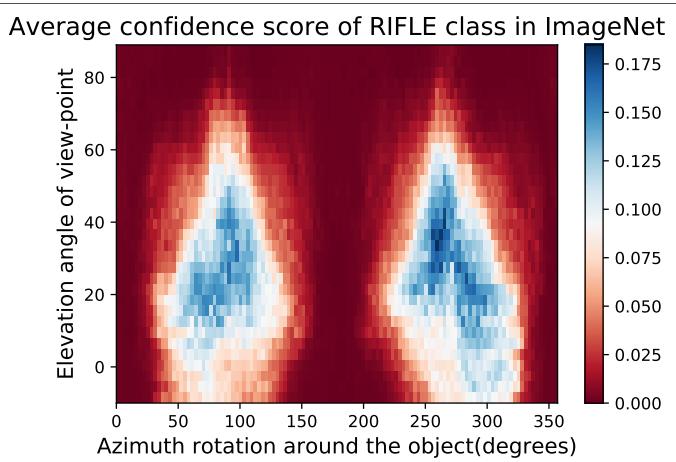
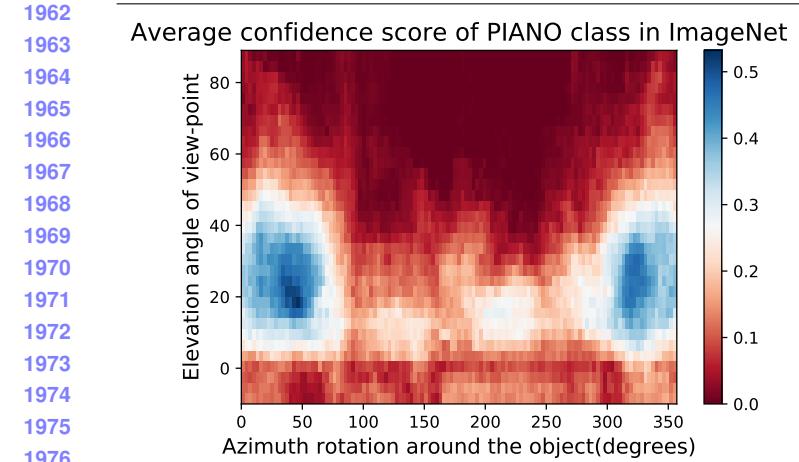
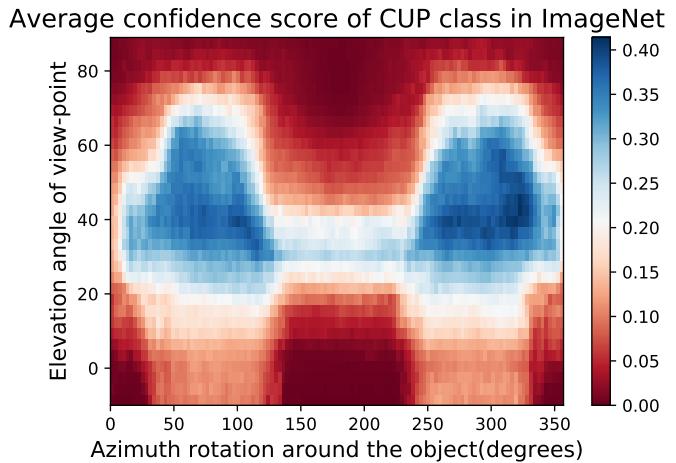
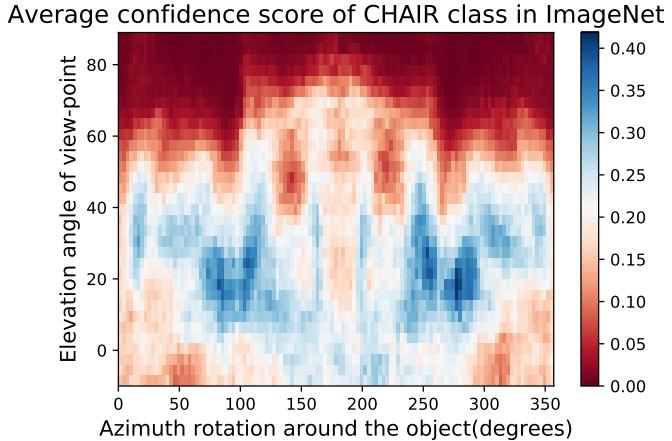


Figure 11: **Data Semantic Maps DSM-II**. Visualizing Semantic Data Bias in the common training dataset (*i.e.* ImageNet [6]) By averaging the Networks Semantic Maps (NSM) of different networks and on different shapes, Different classes have different semantic bias in ImageNet as clearly shown in the maps above. The symmetry in the maps are attributed to the 3D symmetry of the objects.

1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051

2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105

3. Analysis

3.1. Detected Robust Regions

We applies the algorithms on 2 semantic parameters (the azimuth angle of the camera around the object , and the elevation angle around the object) that are commonly used in the literature [4, 1]. When we use 1 parameter (the azimuth) , we fix the elevation to 35° . We used 2 instead of larger numbers because it is easier to verify and visualize 2D unlike higher dimensions. Also, the complexity of sampling increase exponentially with dimensionality for those algorithms (albeit being much better than grid sampling , see Table 1). The regions in Figure 8,9 look vertical rectangles most of the time. This is because the scale of the horizontal-axis (0,360) is much smaller than the vertical axis (-10,90), so most regions are squares but looks rectangles because of figure scales.

3.2. hyper-parameters

How to select the hyper parameters in all the above algorithms ? The answer is not straight forward. For the λ in the naive approach , it is set experimentally by trying different values and using the one which detect some regions that known to behave robustly. The values we found for this are $\lambda = 0.07 \sim 0.1$. The learning rate η is easier to detect with observing the convergence and definitely depends on the full range of study . A rule of thumb is to make 0.001of the full range. for the OIR formulations, we have the boundary factor α which we set to 0.05 . A rule of thumb is to set it to be between $0.5 \sim 1/N$, where N is the number of samples needed for that dimensions to adequately characterize the space. In our case $N = 180$, so $1/180 \approx 0.005$. The only hyperparameter with mathematically established bound is the emphasis factor of the OIR_W formulation β . The bound shown in Table 1 which is $0 \leq \beta \leq \frac{2}{2n-1}$ can be shown as follows. We start from Eq (48). This is the actual expression for the special masks (we apologize in the typos in the main paper). As we can see , the most important term is γ_n . IT dictates how the function at the boundaries determine the next move of the bounds. Here γ_n should always be positive to insure the correct direction of the move for a positive function evaluation.

$$\begin{aligned} \gamma_n &> 0 \\ 2 - (2n - 1)\beta &> 0 \\ \beta &< \frac{2}{2n - 2} \end{aligned} \quad (65)$$

3.3. Detest

The data set used is collected from ShapeNet [3] and consists of 10 classes and 10 shapes eaach that are all identified by at least ResNet50 trained on ImagNet. This criteria is important to obtain valid NSM

and DSM. The classes are [’aeroplane’,”bathtub”,’bench’, ’bottle’,’chair’,”cup”,”piano”,’rifle’,’vase’,”toilet”]. Part of the dataset is shown in Figure 6,7. 4 shapes faced difficulty of rendering during the SRVR experiments , therefore , they were replaced by another shapes from the same class.

3.4. Possible Future Directions

We analyse DNNs from a semantic lens and show how more confident networks tend to create adversarial semantic regions inside highly confident regions. We developed a bottom-up approach to analyse the networks semantically by growing adversarial regions, which scales well with dimensionality and we use it to benchmark the semantic robustness of DNNs. We aim to investigate how to use the insights we gain from this work to develop and train semantically robust networks from the start while maintaining the accuracy. Another direct extension of our work is to develop large scale semantic robustness challenge where we label these robust/adversarial regions in the semantic space and release some of them to allow for training and then we test the trained models on hidden test set to measure robustness while reporting the accuracy on ImageNet validation set to make sure that the features of the model did not get affected by the robust training.

Algorithm 1: Robust n-dimensional Region Finding for Black-Box DNNs by Outer-Inner Ratios

Requires: Semantic Function of a DNN $f(\mathbf{u})$ in Eq (1), initial semantic parameter \mathbf{u}_0 , number of iterations T , learning rate η , object shape \mathbf{S}_z of class label z , boundary factor α , smallll ϵ

Form constant binary matrices $\mathbf{M}, \overline{\mathbf{M}}, \mathbf{M}_{\mathbb{Q}}, \overline{\mathbf{M}}_{\mathbb{Q}}, \mathbf{M}_{\mathbb{D}}, \overline{\mathbf{M}}_{\mathbb{D}}$

Initialize bounds $\mathbf{a}_0 \leftarrow \mathbf{u}_0 - \epsilon \mathbf{1}$, $\mathbf{b}_0 \leftarrow \mathbf{u}_0 + \epsilon \mathbf{1}$

$\mathbf{r}_0 \leftarrow \mathbf{a}_0 - \mathbf{b}_0$, update region volume Δ_0 as in Eq (58)

for $t \leftarrow 1$ **to** T **do**

- form the all-corners function vectors $f_{\mathbb{D}}, f_{\mathbb{Q}}$ as in Eq (29)
- $\nabla_{\mathbf{a}} L \leftarrow 2\Delta_{t-1}\text{diag}^{-1}(\mathbf{r}_{t-1}) \left(2\overline{\mathbf{M}}f_{\mathbb{D}} - \overline{\mathbf{M}}_{\mathbb{Q}}f_{\mathbb{Q}} \right)$
- $\nabla_{\mathbf{b}} L \leftarrow 2\Delta_{t-1}\text{diag}^{-1}(\mathbf{r}_{t-1}) \left(-2\mathbf{M}f_{\mathbb{D}} + \mathbf{M}_{\mathbb{Q}}f_{\mathbb{Q}} \right)$
- update bounds: $\mathbf{a}_t \leftarrow \mathbf{a}_{t-1} - \eta \nabla_{\mathbf{a}} L$,
- $\mathbf{b}_t \leftarrow \mathbf{b}_{t-1} - \eta \nabla_{\mathbf{b}} L$
- $\mathbf{r}_t \leftarrow \mathbf{a}_t - \mathbf{b}_t$, update region volume Δ_t as in Eq (58)

end

Returns: robust region bounds: $\mathbf{a}_T, \mathbf{b}_T$.

2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159

2160	Analysis Approach	Paradigm	Total Sampling complexity	Black -box Functions	Forward pass /step	Backward pass /step	Identification Capabality	Hyper-parameters
2161	Grid Sampling	top-down	$\mathcal{O}(N^n)$ $N \gg 2$	✓	-	-	Fully identifies the semantic map of DNN	no hyper-parameters
2162	Naive	bottom-up	$\mathcal{O}(2^n)$	✓	2^n	0	finds strong robust regions only around \mathbf{u}_0	λ , experimentally determined
2163	OIR_B	bottom-up	$\mathcal{O}(2^{n+1})$	✓	2^{n+1}	0	finds strong and week robust regions around \mathbf{u}_0	α , experimentally determined
2164	OIR_W	bottom-up	$\mathcal{O}(2^n)$	✗	2^n	2^n	finds strong and week robust regions around \mathbf{u}_0	$0 \leq \beta < \frac{2}{2^n-1}$ dependeds on n and Lipschitz constant \mathbb{L}

Table 1: **Semantic Analysis Techniques:** Comparing different approaches to analyse the semantic robustness of DNN.**Algorithm 2:** Robust n-dimensional Region Finding for White-Box DNNs by Outer-Inner Ratios

Requires: Semantic Function of a DNN $f(\mathbf{u})$ in Eq (1), initial semantic parameter \mathbf{u}_0 , learning rate η , object shape \mathbf{S}_z of class label z , emphasis factor β , small ϵ . Form constant binary matrices $\mathbf{M}, \bar{\mathbf{M}}, \mathbf{M}_{\mathbb{D}}, \bar{\mathbf{M}}_{\mathbb{D}}$. Initialize bounds $\mathbf{a}_0 \leftarrow \mathbf{u}_0 - \epsilon \mathbf{1}$, $\mathbf{b}_0 \leftarrow \mathbf{u}_0 + \epsilon \mathbf{1}$. $\mathbf{r}_0 \leftarrow \mathbf{a}_0 - \mathbf{b}_0$, update region volume Δ_0 as in Eq (58).

for $t \leftarrow 1$ **to** T **do**

- form the all-corners function vector $f_{\mathbb{D}}$ as in Eq (29)
- form the all-corners gradients matrix $\mathbf{G}_{\mathbb{D}}$ as in Eq (63)
- form the gradient selection vectors $\mathbf{s}, \bar{\mathbf{s}}$ as in Eq (49)
- $\nabla_{\mathbf{a}} L \leftarrow \Delta_{t-1} (\text{diag}^{-1}(\mathbf{r}_{t-1}) \bar{\mathbf{M}}_{\mathbb{D}} \mathbf{f}_{\mathbb{D}} + \beta \text{diag}(\bar{\mathbf{M}} \mathbf{G}_{\mathbb{D}} + \beta \bar{\mathbf{s}}))$
- $\nabla_{\mathbf{b}} L \leftarrow \Delta_{t-1} (-\text{diag}^{-1}(\mathbf{r}_{t-1}) \mathbf{M}_{\mathbb{D}} \mathbf{f}_{\mathbb{D}} + \beta \text{diag}(\mathbf{M} \mathbf{G}_{\mathbb{D}}) + \beta \mathbf{s})$
- update bounds: $\mathbf{a}_t \leftarrow \mathbf{a}_{t-1} - \eta \nabla_{\mathbf{a}} L$,
- $\mathbf{b}_t \leftarrow \mathbf{b}_{t-1} - \eta \nabla_{\mathbf{b}} L$
- $\mathbf{r}_t \leftarrow \mathbf{a}_t - \mathbf{b}_t$, update region volume Δ_t as in Eq (58)

end

Returns: robust region bounds: $\mathbf{a}_T, \mathbf{b}_T$.

2268	References	2322
2269		2323
2270	[1] M. A. Alcorn, Q. Li, Z. Gong, C. Wang, L. Mai, W. Ku, and A. Nguyen. Strike (with) a pose: Neural networks are easily fooled by strange poses of familiar objects. <i>CoRR</i> , abs/1811.11553, 2018. 20	2324
2271		2325
2272		2326
2273		2327
2274	[2] S. Boyd and L. Vandenberghe. <i>Convex Optimization</i> . Cam- bridge University Press, New York, NY, USA, 2004. 1	2328
2275		2329
2276	[3] A. X. Chang, T. Funkhouser, L. Guibas, P. Hanrahan, Q. Huang, Z. Li, S. Savarese, M. Savva, S. Song, H. Su, J. Xiao, L. Yi, and F. Yu. ShapeNet: An Information-Rich 3D Model Repository. Technical Report arXiv:1512.03012 [cs.GR], Stanford University — Princeton University — Toy- ota Technological Institute at Chicago, 2015. 20	2330
2277		2331
2278		2332
2279		2333
2280		2334
2281	[4] A. Hamdi, M. Müller, and B. Ghanem. SADA: semantic adversarial diagnostic attacks for autonomous applications. <i>CoRR</i> , abs/1812.02132, 2018. 20	2335
2282		2336
2283		2337
2284	[5] R. G. Ródenas, M. L. López, and D. Verastegui. Extensions of dinkelbach’s algorithm for solving non-linear fractional pro- gramming problems. <i>Top</i> , 7(1):33–70, Jun 1999. 3	2338
2285		2339
2286		2340
2287	[6] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. S. Bernstein, A. C. Berg, and F. Li. Imagenet large scale visual recognition chal- lenge. <i>CoRR</i> , abs/1409.0575, 2014. 8 , 18 , 19	2341
2288		2342
2289		2343
2290		2344
2291	[7] Z. Shou, H. Gao, L. Zhang, K. Miyazawa, and S.-F. Chang. Autoloc: weakly-supervised temporal action localization in untrimmed videos. In <i>Proceedings of the European Confer- ence on Computer Vision (ECCV)</i> , pages 154–171, 2018. 1	2345
2292		2346
2293		2347
2294	[8] A. H. Stroud. Methods of numerical integration (philip j. davis and philip rabinowitz). <i>SIAM Review</i> , 18(3):528–2, 07 1976. Copyright - Copyright] 1976 Society for Industrial and Applied Mathematics; Last updated - 2012-03-05; CODEN - SIREAD. 6	2348
2295		2349
2296		2350
2297		2351
2298		2352
2299	[9] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna. Rethinking the inception architecture for computer vision. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pages 2818–2826, 2016. 11 , 12	2353
2300		2354
2301		2355
2302		2356
2303		2357
2304		2358
2305		2359
2306		2360
2307		2361
2308		2362
2309		2363
2310		2364
2311		2365
2312		2366
2313		2367
2314		2368
2315		2369
2316		2370
2317		2371
2318		2372
2319		2373
2320		2374
2321		2375