

# Verification of Non-Linear Specifications for Neural Networks

Hanqing Chao

# VERIFICATION OF NON-LINEAR SPECIFICATIONS FOR NEURAL NETWORKS

**Chongli Qin\*, Krishnamurthy (Dj) Dvijotham\*, Brendan O'Donoghue, Rudy Bunel, Robert Stanforth, Sven Gowal, Jonathan Uesato, Grzegorz Swirszcz, Pushmeet Kohli**

DeepMind  
London, N1C 4AG, UK  
correspondence: chongliqin@google.com

Hanqing Chao

# Adversarial Attack

- Loopholes exists on a well trained network
- The output cannot hold in some very proximity of the training examples



“panda”  
57.7% confidence

+ .007 ×



“nematode”  
8.2% confidence

=



“gibbon”  
99.3 % confidence

# Specification Verification

- Specification:

Some properties we expect a train network to have:

The output could hold in a very proximity of the training examples

$$y = f(x + r) \text{ are consistent for } \forall ||r|| < \epsilon, \epsilon > 0$$

# Specification Verification

- Specification:  $y = f(x + r)$  are consistent for  $\forall ||r|| < \epsilon, \epsilon > 0$
- $\epsilon$  is some kind measurement demonstrating the specification of the network (on point  $x$ ).

Complete verification:

Guaranteed to either find a proof that the specification is true  
Or find a counterexample proving that the specification is untrue.



$\epsilon = 0.5$ , True specification

$\epsilon = 0.3$ , Lower bound

Incomplete verification:

May not find a proof even if the specification is true.

However, if they do find a proof, the specification is guaranteed to be true.

# Specification Verification

- Specification:  $y = f(x + r)$  are consistent for  $\forall ||r|| < \epsilon, \epsilon > 0$
- $\epsilon$  is some kind measurement demonstrating the specification of the network (on point  $x$ ).

Complete verification:

Guaranteed to either find a proof that the specification is true  
Or find a counterexample proving that the specification is untrue.

—  $\epsilon = 0.5$ , True specification

—  $\epsilon = 0.3$ , Lower bound

Incomplete verification:

May not find a proof even if the specification is true.

However, if they do find a proof, the specification is guaranteed to be true.

# Problem Formulation

- A network mapping:  $f : \mathcal{X} \rightarrow \mathcal{Y}$ ,
- With specification  $F : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$
- $F$  is satisfied if

$$F(x, y) \leq 0 \quad \forall x \in \mathcal{S}_{\text{in}}, y = f(x).$$

- Where  $\mathcal{S}_{\text{in}}$  is the neighbor region of  $x^{\text{nom}}$ , for instance:

$$\mathcal{S}_{\text{in}} = \{x : \|x - x^{\text{nom}}\|_{\infty} \leq \delta\},$$

# Non-Linear Specification

$$F(x, y) \leq 0 \quad \forall x \in \mathcal{S}_{\text{in}}, y = f(x).$$

- Output consistency: linear specification

$$F(x, y) = c^T y + d \leq 0 \quad \forall x \in \mathcal{S}_{\text{in}}, y = f(x)$$

- Non-linear specification:
  - In a learned dynamics model of a physical system, law of conservation of energy should be met;
  - For a classifier, its output labels under adversarial perturbations (adversarial attack) should be semantically consistent;
  - In a system predicts the summation of handwritten digits, errors should be bounded.



# Conservation of Energy

$$F(x, y) \leq 0 \quad \forall x \in \mathcal{S}_{\text{in}}, y = f(x).$$

- Consider a simple pendulum with damping

$$E(w, h, \omega) = \underbrace{mgh}_{\text{potential}} + \underbrace{\frac{1}{2}ml^2\omega^2}_{\text{kinetic}},$$

- $w$  and  $h$  represent the horizontal and vertical coordinates of the pendulum respectively, and  $\omega$  refers to the angular velocity

# Conservation of Energy

$$F(x, y) \leq 0 \quad \forall x \in \mathcal{S}_{\text{in}}, y = f(x).$$

- Consider a simple pendulum with damping

$$E(w, h, \omega) = \underbrace{mgh}_{\text{potential}} + \underbrace{\frac{1}{2}ml^2\omega^2}_{\text{kinetic}},$$

- $w$  and  $h$  represent the horizontal and vertical coordinates of the pendulum respectively, and  $\omega$  refers to the angular velocity

$$E(w', h', \omega') \leq E(w, h, \omega)$$

# Semantic Consistency

$$F(x, y) \leq 0 \quad \forall x \in \mathcal{S}_{\text{in}}, y = f(x).$$

$$\mathbb{E} [d(i, j)] = \sum_j d(i, j) P(j|x)$$

- $d(i, j)$  predefined distance between two classes  $i, j$ ;  $d(i, i) = 0$
- $P(j|x)$  the probability that the classifier assigns

# Semantic Consistency

$$F(x, y) \leq 0 \quad \forall x \in \mathcal{S}_{\text{in}}, y = f(x).$$

$$\mathbb{E} [d(i, j)] = \sum_j d(i, j) P(j|x)$$

- $d(i, j)$  predefined distance between two classes  $i, j$ ;  $d(i, i) = 0$
- $P(j|x)$  the probability that the classifier assigns

$$\mathbb{E} [d(i, j)] \leq \epsilon.$$

# Summation Error

$$F(x, y) \leq 0 \quad \forall x \in \mathcal{S}_{\text{in}}, y = f(x).$$

$$\mathbb{E}_j \left[ \left| \sum_{n=1}^N (j_n - i_n) \right| \right] = \sum_{j \in J^N} \left| \sum_{n=1}^N (j_n - i_n) \right| \prod_{n=1}^N P(j_n | x_n),$$

- $\{x_n\}_{n=1}^N$  handwritten images and corresponding true transaction values  $\{i_n\}_{n=1}^N$

# Summation Error

$$F(x, y) \leq 0 \quad \forall x \in \mathcal{S}_{\text{in}}, y = f(x).$$

$$\mathbb{E}_j \left[ \left| \sum_{n=1}^N (j_n - i_n) \right| \right] = \sum_{j \in J^N} \left| \sum_{n=1}^N (j_n - i_n) \right| \prod_{n=1}^N P(j_n | x_n),$$

- $\{x_n\}_{n=1}^N$  handwritten images and corresponding true transaction values  $\{i_n\}_{n=1}^N$

$$\mathbb{E}_j \left[ \left| \sum_{n=1}^N (j_n - i_n) \right| \right] \leq \epsilon.$$

# Convex-relaxable Specification

**Assumption 1.** We assume that  $\mathcal{S}_{\text{in}} \subseteq \mathbb{R}^n$ ,  $\mathcal{S}_{\text{out}} \subseteq \mathbb{R}^m$  are compact sets and that we have access to an efficiently computable<sup>1</sup> procedure that takes in  $F, \mathcal{S}_{\text{in}}, \mathcal{S}_{\text{out}}$  and produces a compact convex set  $\mathcal{C}(F, \mathcal{S}_{\text{in}}, \mathcal{S}_{\text{out}})$  such that

$$\mathcal{T}(F, \mathcal{S}_{\text{in}}, \mathcal{S}_{\text{out}}) := \{(x, y, z) : F(x, y) = z, x \in \mathcal{S}_{\text{in}}, y \in \mathcal{S}_{\text{out}}\} \subseteq \mathcal{C}(F, \mathcal{S}_{\text{in}}, \mathcal{S}_{\text{out}}). \quad (11)$$

When the above assumption holds we shall say that the specification

$$F(x, y) \leq 0 \quad \forall x \in \mathcal{S}_{\text{in}}, y = f(x) \quad (12)$$

is convex-relaxable.

# Bounds propagation

- Given bounds of inputs, forward these bounds layer by layer to get the bounds of the last layer.



# Convex Optimization

maximize  $z$

subject to  $(x_0, x_K, z) \in \mathcal{C}(F, \mathcal{S}_{\text{in}}, \mathcal{S}_{\text{out}})$

$$x_{k+1} \in \text{Relax}(g_k)(W_k x_k + b_k, l_k, u_k), k = 0, \dots, K-1$$

$$l_k \leq x_k \leq u_k, \quad k = 0, \dots, K$$

- Here,  $z$  is the value of  $F$
- Main idea is to relax activation function and  $F$  to get boundaries.

# Experiments

—  $\epsilon = 0.5$ , True specification

—  $\epsilon = 0.3$ , Lower bound

Incomplete verification:

May not find a proof even if the specification is true.

However, if they do find a proof, the specification is guaranteed to be true.

How tight is it?

# Experiments

**Verification bound:** This is the fraction of test examples  $x^{\text{nom}}$  for which the specification is provably satisfied over the set  $\mathcal{S}_{\text{in}}(x^{\text{nom}}, \delta)$  using our verification algorithm.

**Adversarial bound:** This is the fraction of test examples  $x^{\text{nom}}$  for which the falsification algorithm based on (8) *was not able to find a counter-example in the set*  $\mathcal{S}_{\text{in}}(x^{\text{nom}}, \delta)$ .

\_\_\_\_\_  $\epsilon = 0.7$ , Adversarial bound

\_\_\_\_\_  $\epsilon = 0.5$ , True specification  $\beta$

\_\_\_\_\_  $\epsilon = 0.3$ , Verification bound

Verification bound  $\leq \beta \leq$  Adversarial bound.

$|\text{Verification bound} - \beta| \leq \text{Adversarial bound} - \text{Verification bound}.$

# Experiments

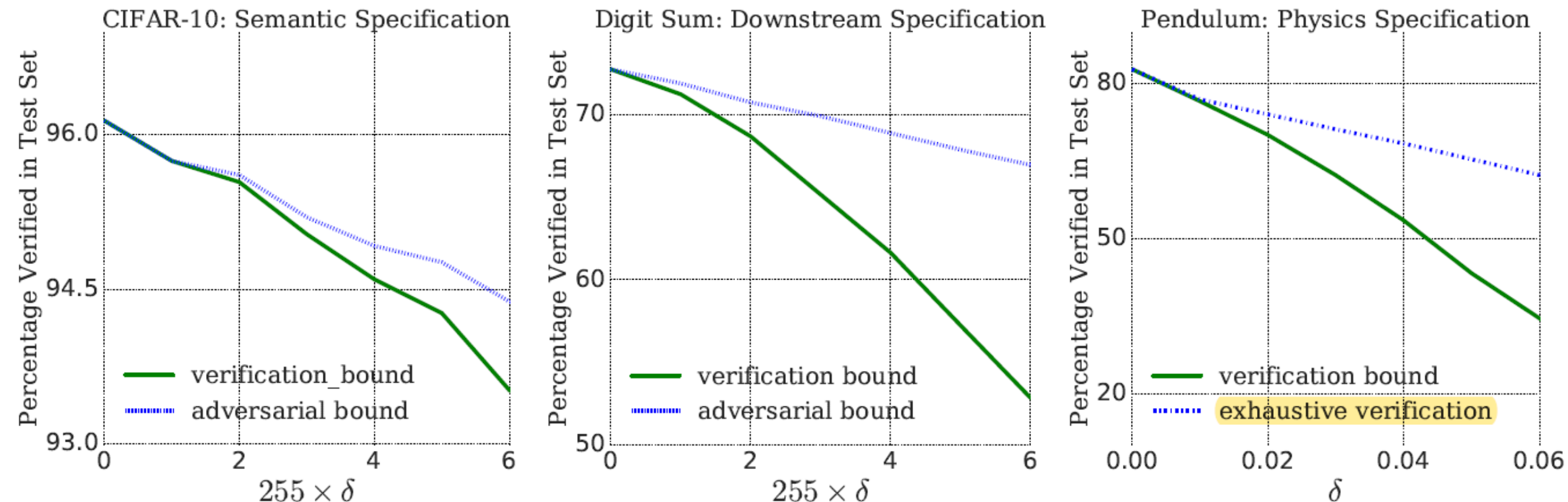
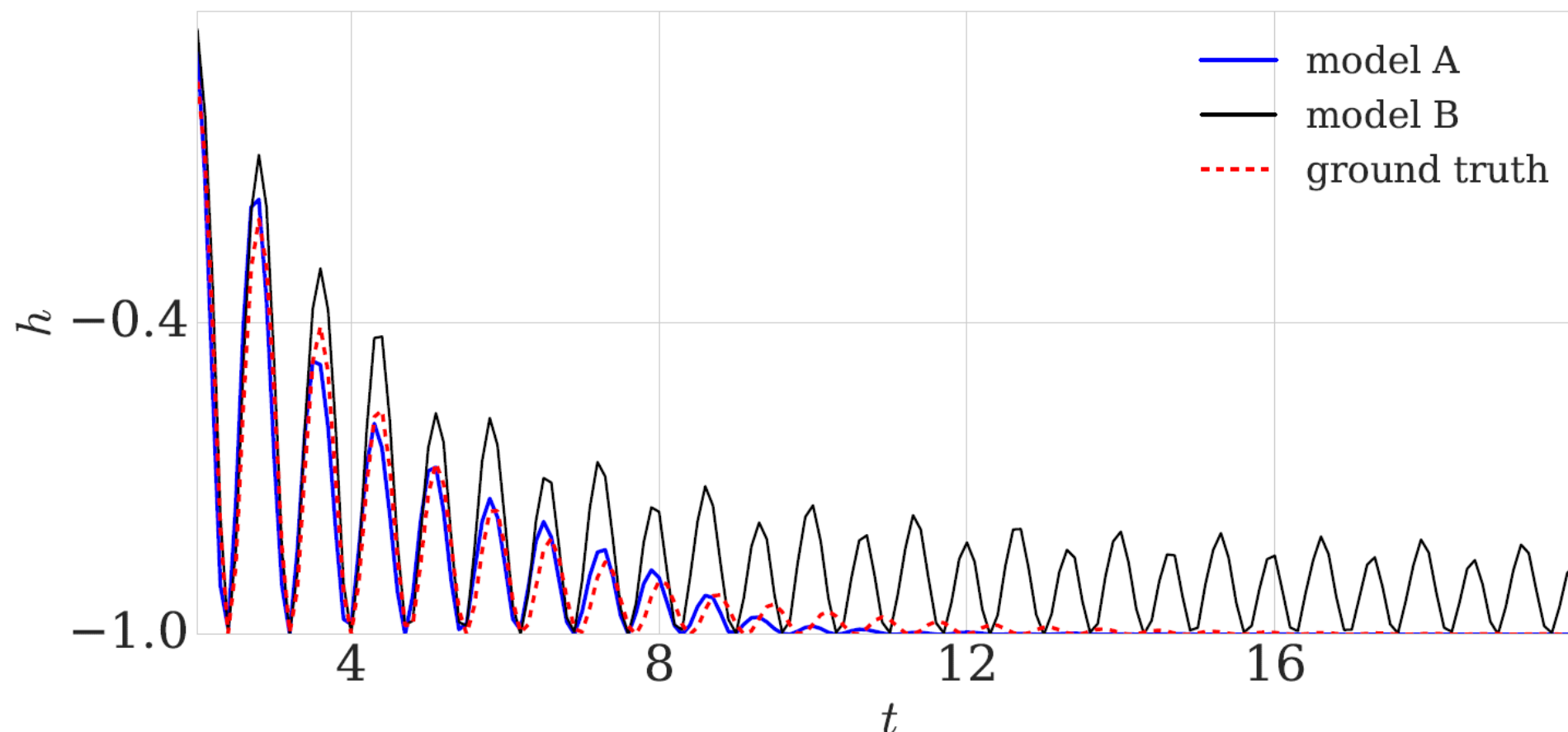


Figure 1: For three specifications, we plot the verification bound and adversarial bound as a function of perturbation size on the input.

# Experiments

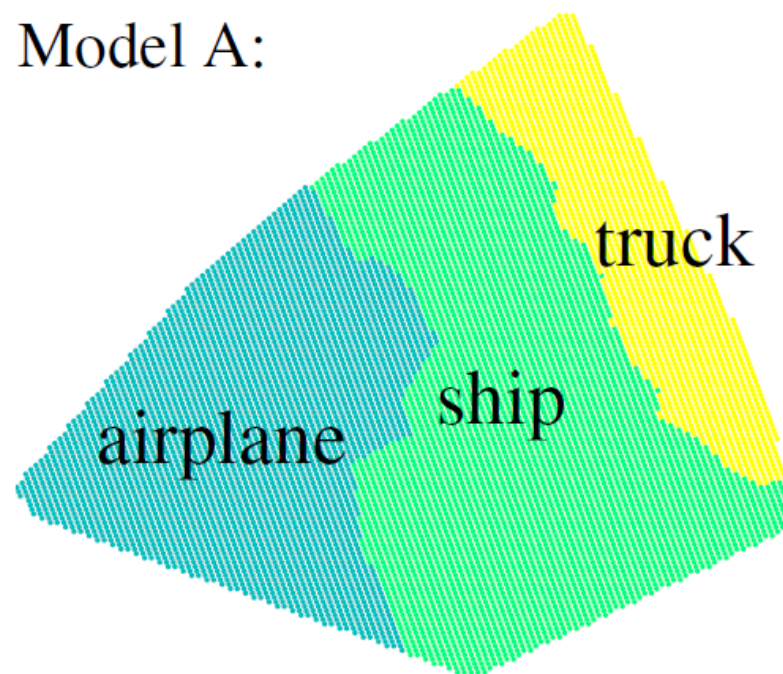
Are the specifications really satisfied in applications?  
Conservation of Energy



# Experiments

Are the specifications really satisfied in applications?  
Semantic Consistency

Model A:



Model B:

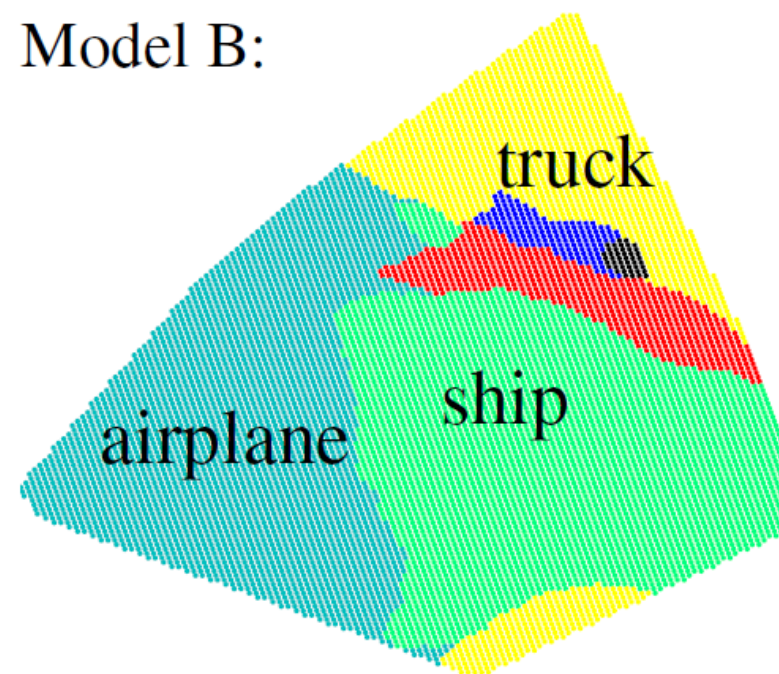


Figure 3: The projection of the decision boundaries onto a two dimensional surface formed by interpolating between three images belonging to the same semantic category (vehicles) - aeroplane (cyan), ship (green) and truck (yellow). The **red/blue/black** regions represent **bird/cat/frog** respectively).

# Thanks !