

Supplementary material for "RISE: Randomized Input Sampling for Explanation of Black-box Models"

Vitali Petsiuk
vpetsiuk@bu.edu

Boston University
Boston, USA

Abir Das
dasabir@bu.edu

Kate Saenko
saenko@bu.edu

Algorithm to compute deletion score.

Algorithm 1

```

1: procedure DELETION
2:   Input: black box  $f$ , image  $I$ , importance map  $S$ , number of pixels  $N$  removed per step
3:   Output: deletion score  $d$ 
4:    $n \leftarrow 0$ 
5:    $h_n \leftarrow f(I)$ 
6:   while  $I$  has non-zero pixels do
7:     According to  $S$ , set next  $N$  pixels in  $I$  to 0
8:      $n \leftarrow n + 1$ 
9:      $h_n \leftarrow f(I)$ 
10:   $d \leftarrow \text{AreaUnderCurve}(h_i \text{ vs. } i/n, \forall i = 0, \dots, n)$ 
11:  return  $d$ 

```

Algorithm to compute insertion score.

Algorithm 2

```

1: procedure INSERTION
2:   Input: black box  $f$ , image  $I$ , importance map  $S$ , number of pixels  $N$  removed per step
3:   Output: insertion score  $d$ 
4:    $n \leftarrow 0$ 
5:    $I' \leftarrow \text{Blur}(I)$ 
6:    $h_n \leftarrow f(I)$ 
7:   while  $I \neq I'$  do
8:     According to  $S$ , set next  $N$  pixels in  $I'$  to corresponding pixels in  $I$ 
9:      $n \leftarrow n + 1$ 
10:     $h_n \leftarrow f(I')$ 
11:   $d \leftarrow \text{AreaUnderCurve}(h_i \text{ vs. } i/n, \forall i = 0, \dots, n)$ 
12:  return  $d$ 

```

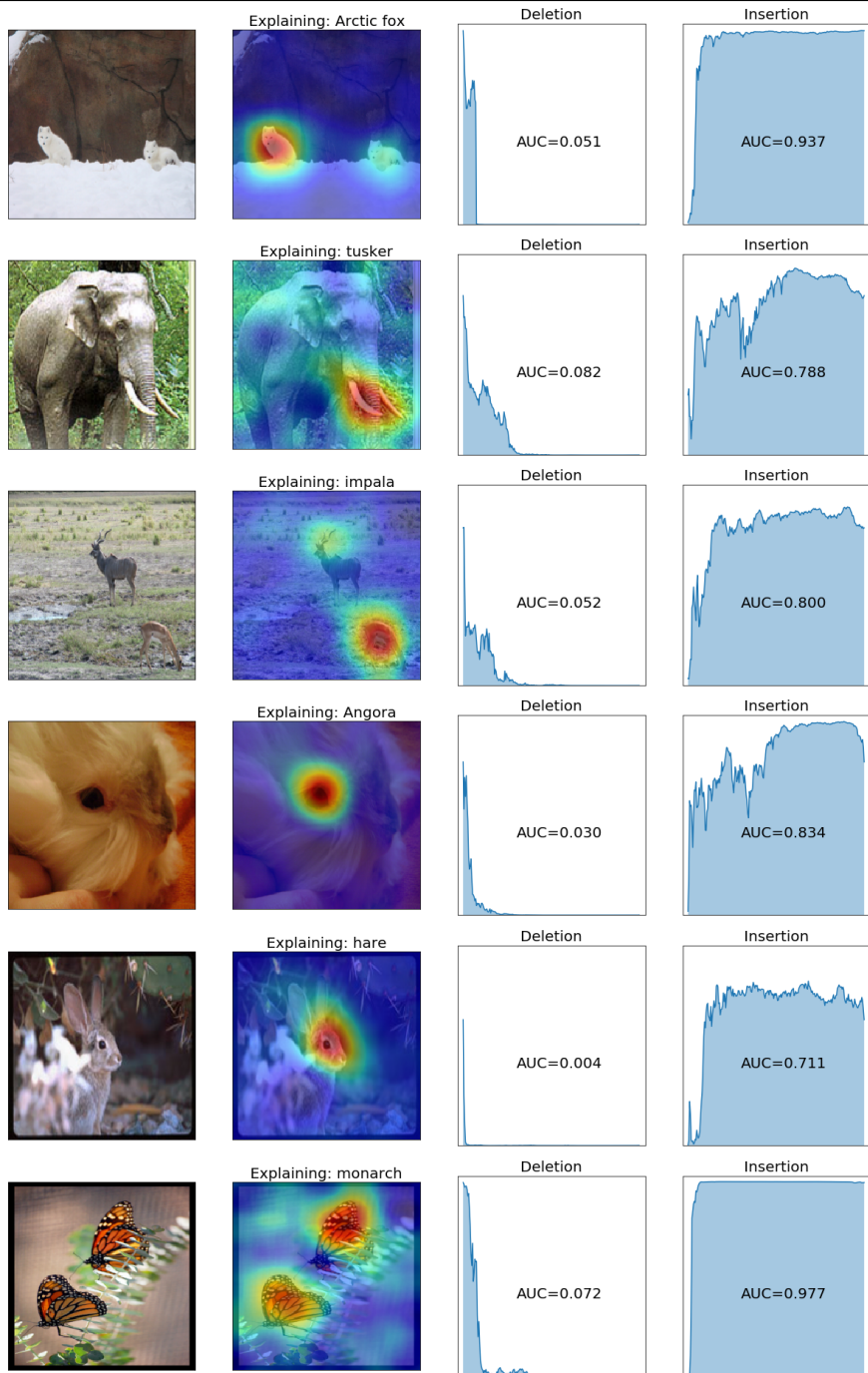


Figure 1: RISE generated importance maps (second column) for representative images (first column) with deletion (third column) and insertion curves (fourth column).

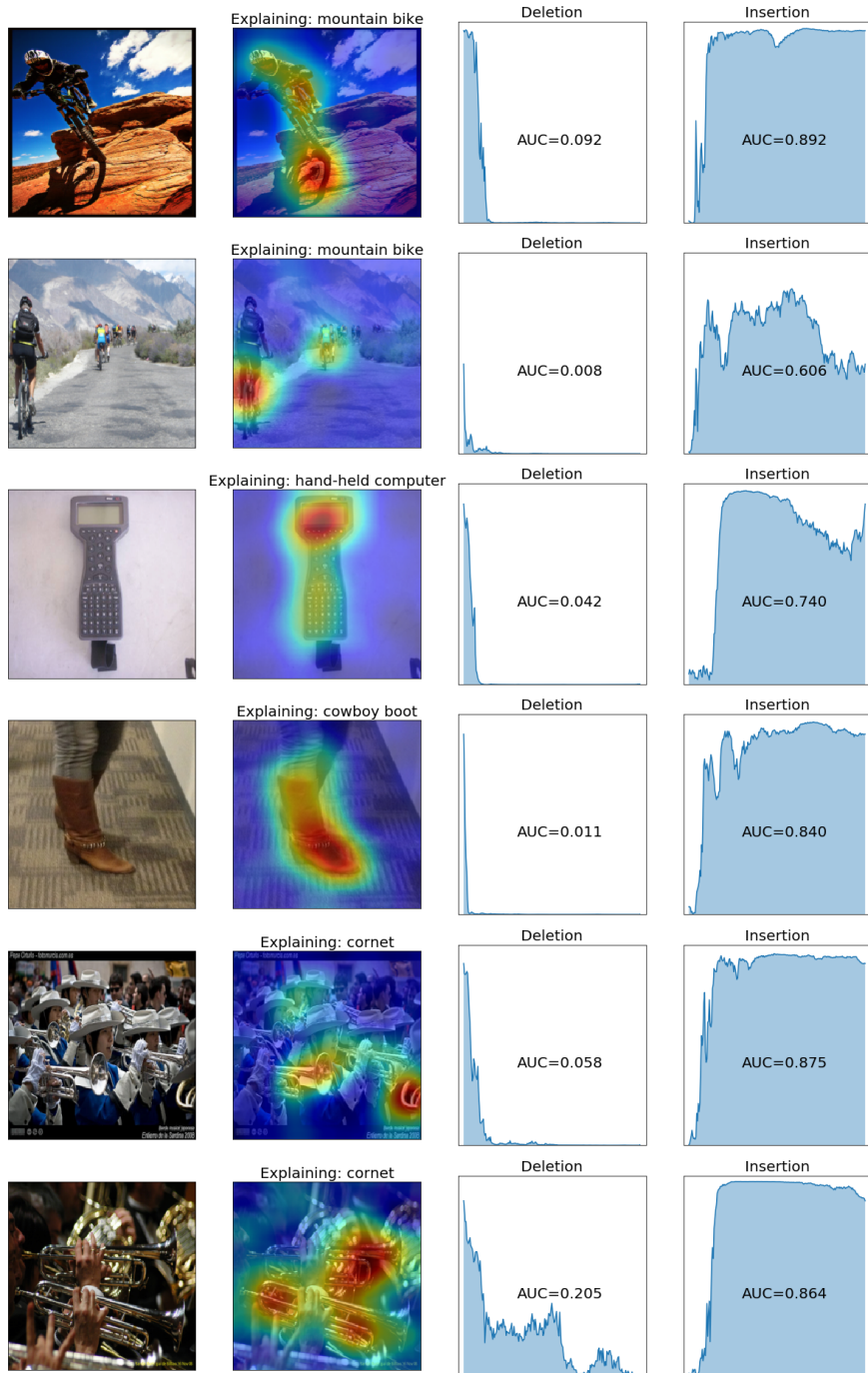


Figure 2: RISE generated importance maps (second column) for representative images (first column) with deletion (third column) and insertion curves (fourth column).

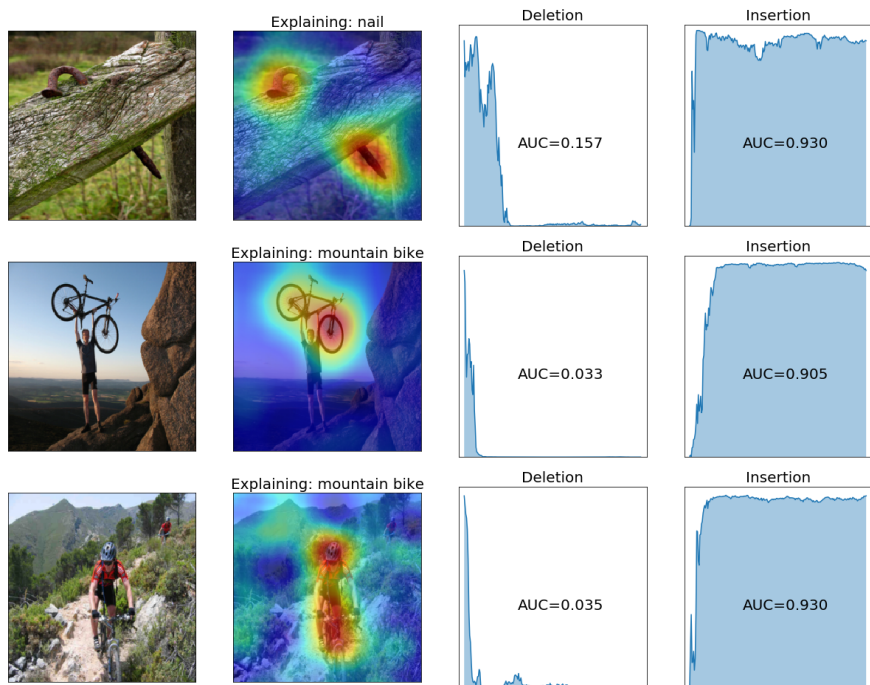


Figure 3: RISE generated importance maps (second column) for representative images (first column) with deletion (third column) and insertion curves (fourth column).

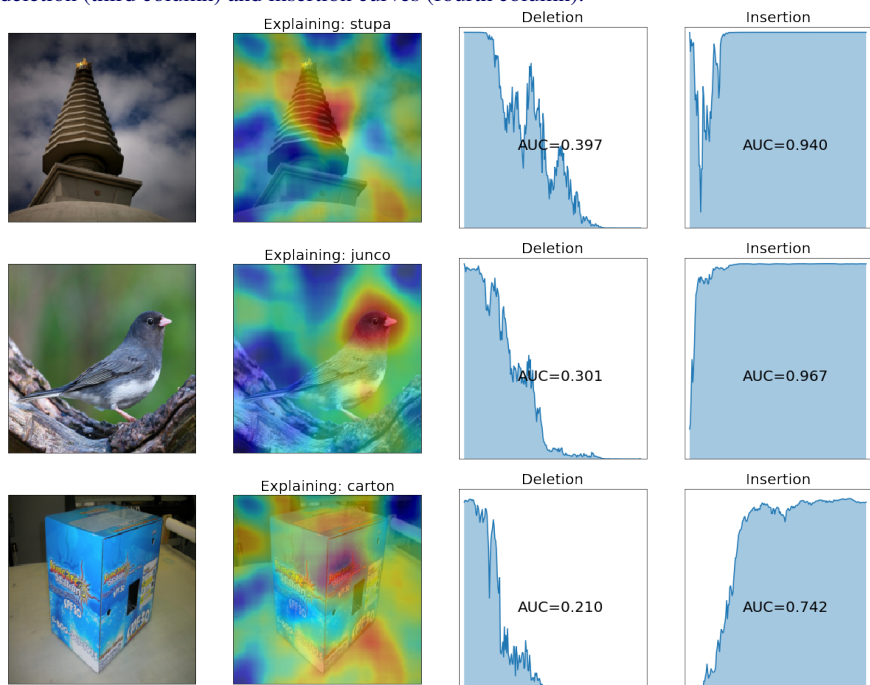


Figure 4: Failure cases. In some cases RISE does pick up more important features, but cannot get rid of the background noise (in part due to MC approximation with only a subset) like in rows 1 and 2.