# Exploring Concept Contribution Spatially:
# Hidden Layer Interpretation with *Spatial Activation Concept Vector*

Andong Wang, Wei-Ning Lee

## I. Introduction

**Problem**

*TCAV* calculates concept contribution to a target class based on *a whole hidden layer*

→ Evaluation may be interfered with by *redundant background features*

$$\nabla h_{l,c}(f_l(\boldsymbol{x}))v_{TCAV}^l$$

**Spatial Activation Concept Vector (SACV)**

Which identifies the *relevant spatial locations* to the *query concept* while evaluating their contributions to the model prediction of the *target class*
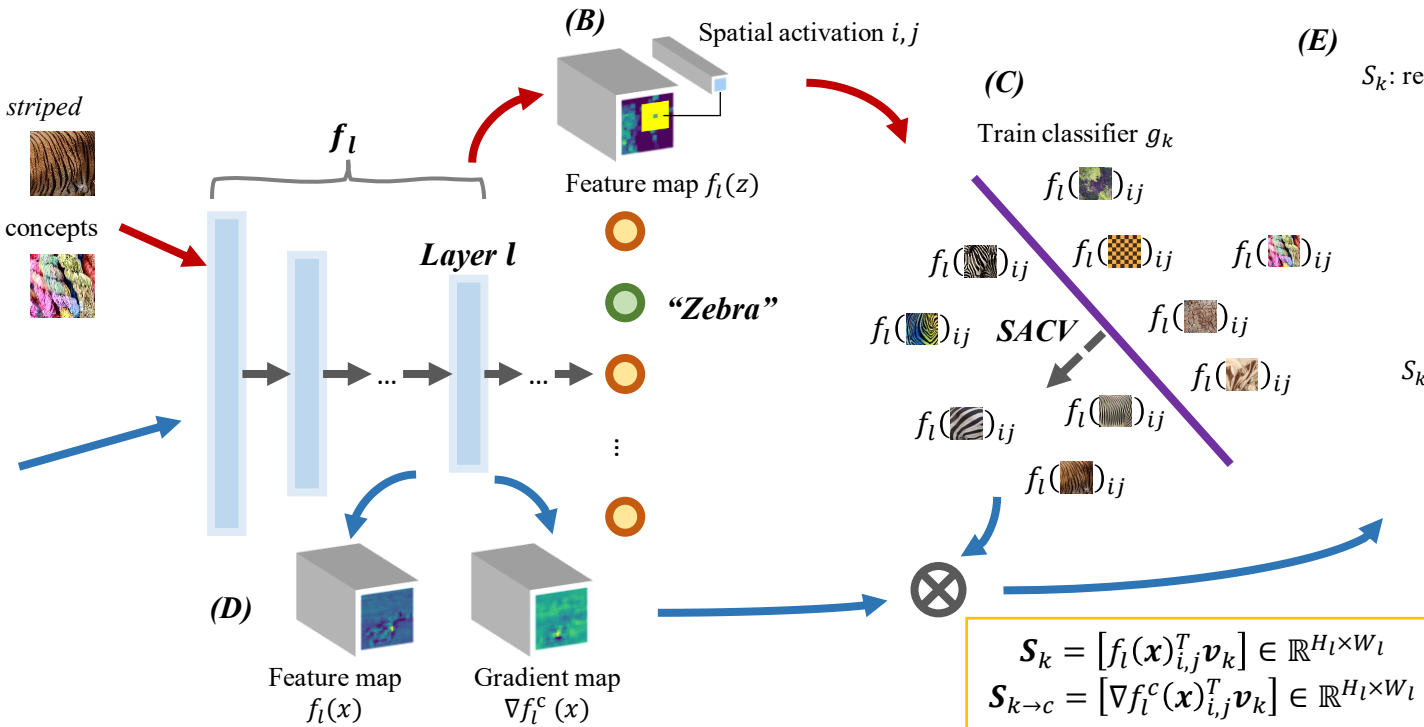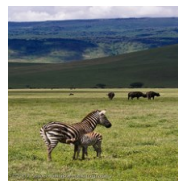
https://github.com/AntonotnaWang/Spatial-Activation-Concept-Vector
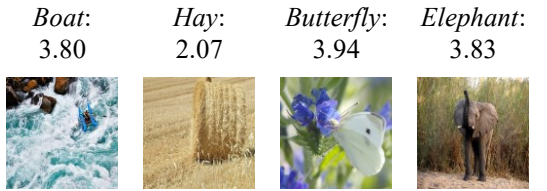
Large percentage of redundant background info

Input Image $x$

$f_l$

**Layer $l$**

Target object

## II. Method



*(A)* Sample images $\{z_k\}$ of concept *striped*

Sample images $\{z_R\}$ of random concepts

Input Image $x$

$f_l$

**Layer $l$**

*"Zebra"*

*(B)* Spatial activation $i, j$

Feature map $f_l(z)$

*(C)* Train classifier $g_k$

$f_l(\ )_{ij}$ ... *SACV*

*(E)* $S_k$: relevant spatial locations to *striped*

$S_{k \to c}$: contribution of *striped* to *zebra*

*(D)* Feature map $f_l(x)$   Gradient map $\nabla f_l^c(x)$

$$S_k = \left[f_l(\boldsymbol{x})_{i,j}^T \boldsymbol{v}_k\right] \in \mathbb{R}^{H_l \times W_l}$$
$$S_{k \to c} = \left[\nabla f_l^c(\boldsymbol{x})_{i,j}^T \boldsymbol{v}_k\right] \in \mathbb{R}^{H_l \times W_l}$$

## III. Results

*(A)* $S_k$: **relevant spatial locations to concept** *striped*

*Zebra* images

| 8.11 | 9.10 | 5.64 | 6.61 |

*Other* images

*Boat*: 3.80   *Hay*: 2.07   *Butterfly*: 3.94   *Elephant*: 3.83

(The numbers indicate $\max\{S_k\}$)

Raw image   features.2   features.5   features.10   features.25   features.30

$S_k$ on different layers

*(B)* $S_{k \to c}$: **contribution of concept** *striped* **to class** *zebra*
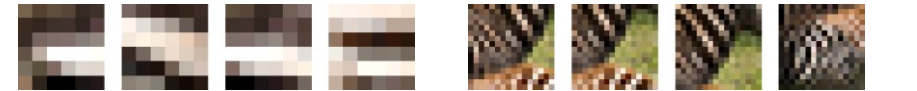
Receptive fields   features.5   features.10

Highest contribution

Lowest contribution