

# CAM

Learning Deep Features for Discriminative Localization

# Problem

The ability to localize objects from CNNs is lost in the FC layers

Need to avoid the use of FC layers

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input (224 × 224 RGB image)					
conv3-64	conv3-64 <b>LRN</b>	conv3-64 <b>conv3-64</b>	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 <b>conv3-128</b>	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 <b>conv1-256</b>	conv3-256 conv3-256 <b>conv3-256</b>	conv3-256 conv3-256 conv3-256 <b>conv3-256</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

# Solution

The ability to localize objects from CNNs is lost in the FC layers

Need to avoid the use of FC layers

Replace multiple FC layers with Global Average Pooling GAP

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input ( $224 \times 224$ RGB image)					
conv3-64	conv3-64 <b>LRN</b>	conv3-64 <b>conv3-64</b>	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 <b>conv3-128</b>	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 <b>conv1-256</b>	conv3-256 conv3-256 <b>conv3-256</b>	conv3-256 conv3-256 conv3-256 <b>conv3-256</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
<b>GAP</b>					
FC-1000					
soft-max					

# What the paper is about

The ability to localize objects from CNNs is lost in the FC layers

Need to avoid the use of FC layers

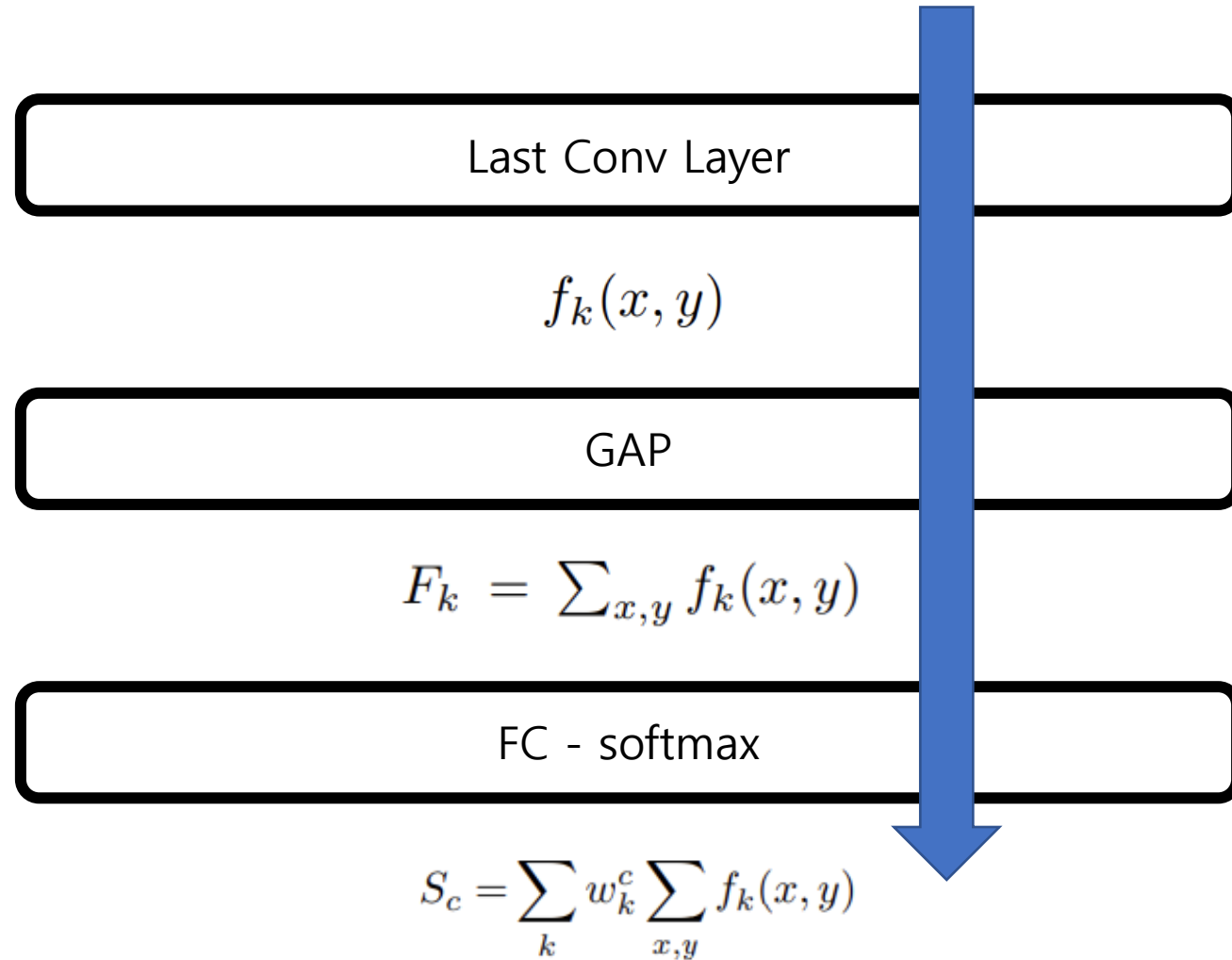
Replace multiple FC layers with Global Average Pooling GAP

Applying GAP for accurate discriminative localization

CAM: Class Activation Mapping

	LRN	conv3-64	conv3-64	conv3-64	conv3-64
maxpool					
conv3-128	conv3-128	conv3-128	conv3-128	conv3-128	conv3-128
		<b>conv3-128</b>	conv3-128	conv3-128	conv3-128
maxpool					
conv3-256	conv3-256	conv3-256	conv3-256	conv3-256	conv3-256
conv3-256	conv3-256	conv3-256	conv3-256	conv3-256	conv3-256
			<b>conv1-256</b>	<b>conv3-256</b>	conv3-256
					<b>conv3-256</b>
maxpool					
conv3-512	conv3-512	conv3-512	conv3-512	conv3-512	conv3-512
conv3-512	conv3-512	conv3-512	conv3-512	conv3-512	conv3-512
			<b>conv1-512</b>	<b>conv3-512</b>	conv3-512
					<b>conv3-512</b>
maxpool					
conv3-512	conv3-512	conv3-512	conv3-512	conv3-512	conv3-512
conv3-512	conv3-512	conv3-512	conv3-512	conv3-512	conv3-512
			<b>conv1-512</b>	<b>conv3-512</b>	conv3-512
					<b>conv3-512</b>
GAP					
FC-1000					
soft-max					

# CAM: Class Activation Mapping



# CAM: Class Activation Mapping

Last Conv Layer

$$f_k(x, y)$$

GAP

$$F_k = \sum_{x,y} f_k(x, y)$$

FC - softmax

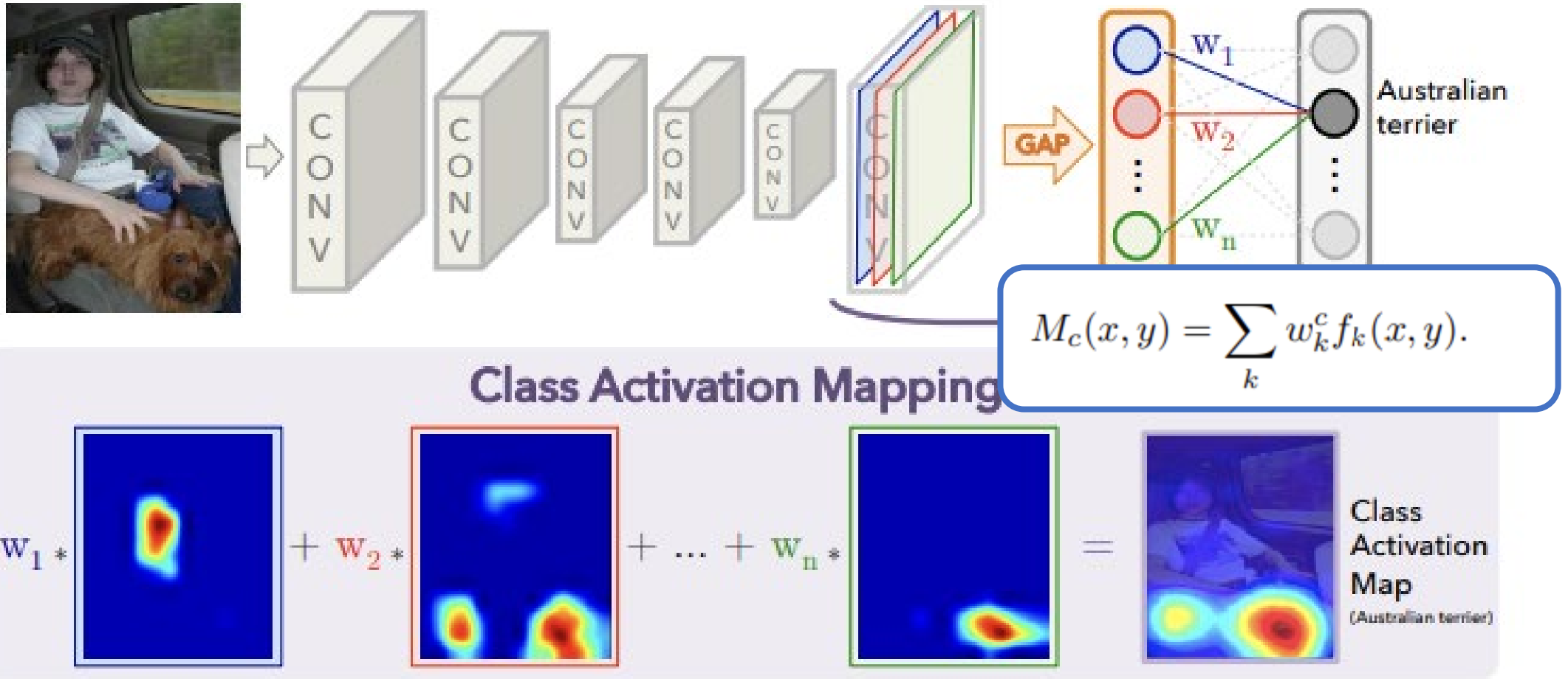
$$S_c = \sum_k w_k^c \sum_{x,y} f_k(x, y)$$

$$S_c = \sum_k w_k^c \sum_{x,y} f_k(x, y)$$

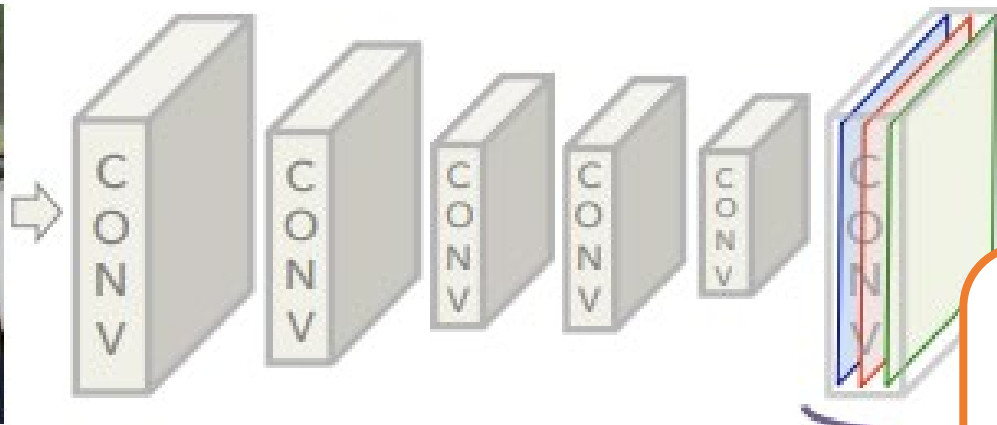
$$= \sum_{x,y} \sum_k w_k^c f_k(x, y).$$

$$M_c(x, y) = \sum_k w_k^c f_k(x, y).$$

# CAM: Class Activation Mapping



# CAM: Class Activation Mapping



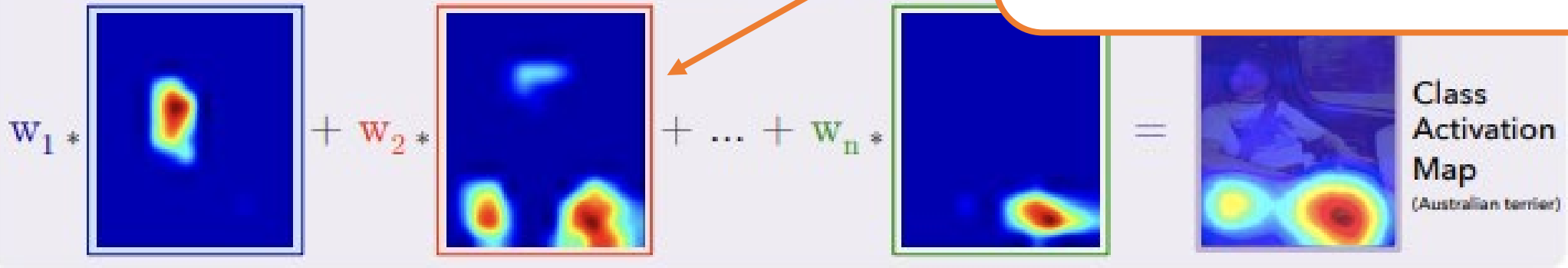
How is upsampling performed?

...

Upsampling

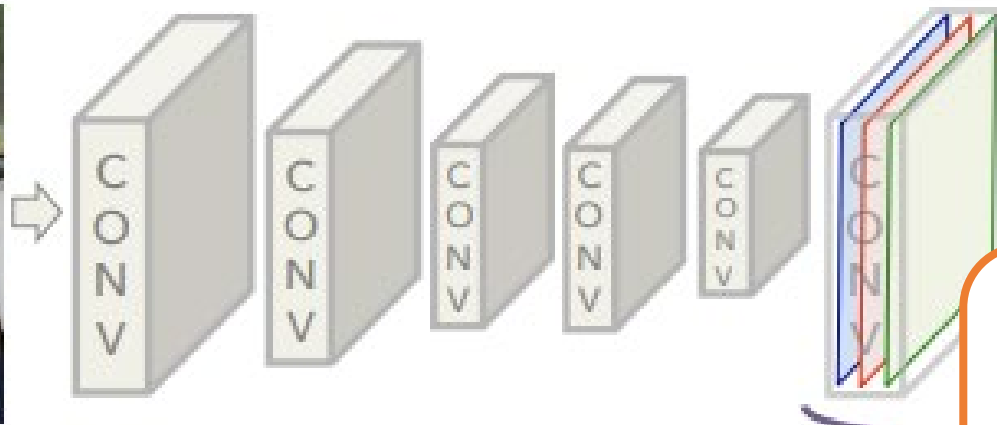
$$M_c(x, y) = \sum_k w_k^c f_k(x, y).$$

Class Activation Mapping





# CAM: Class Activation Mapping



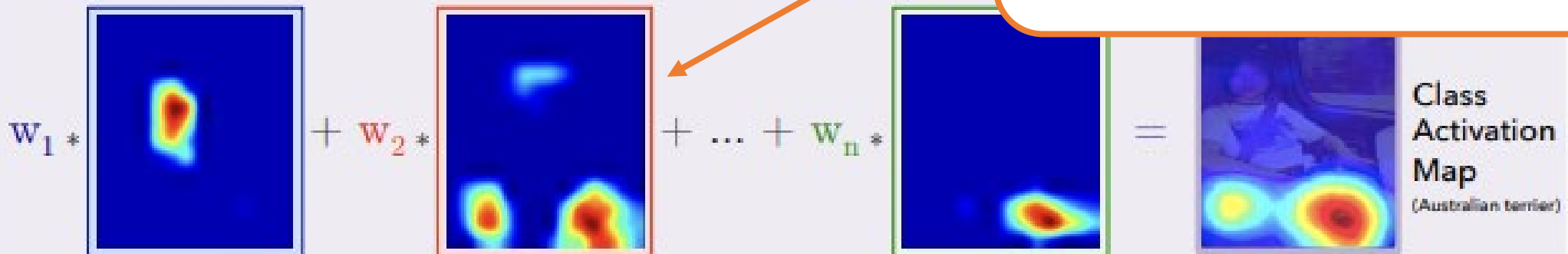
How is upsampling performed?

...

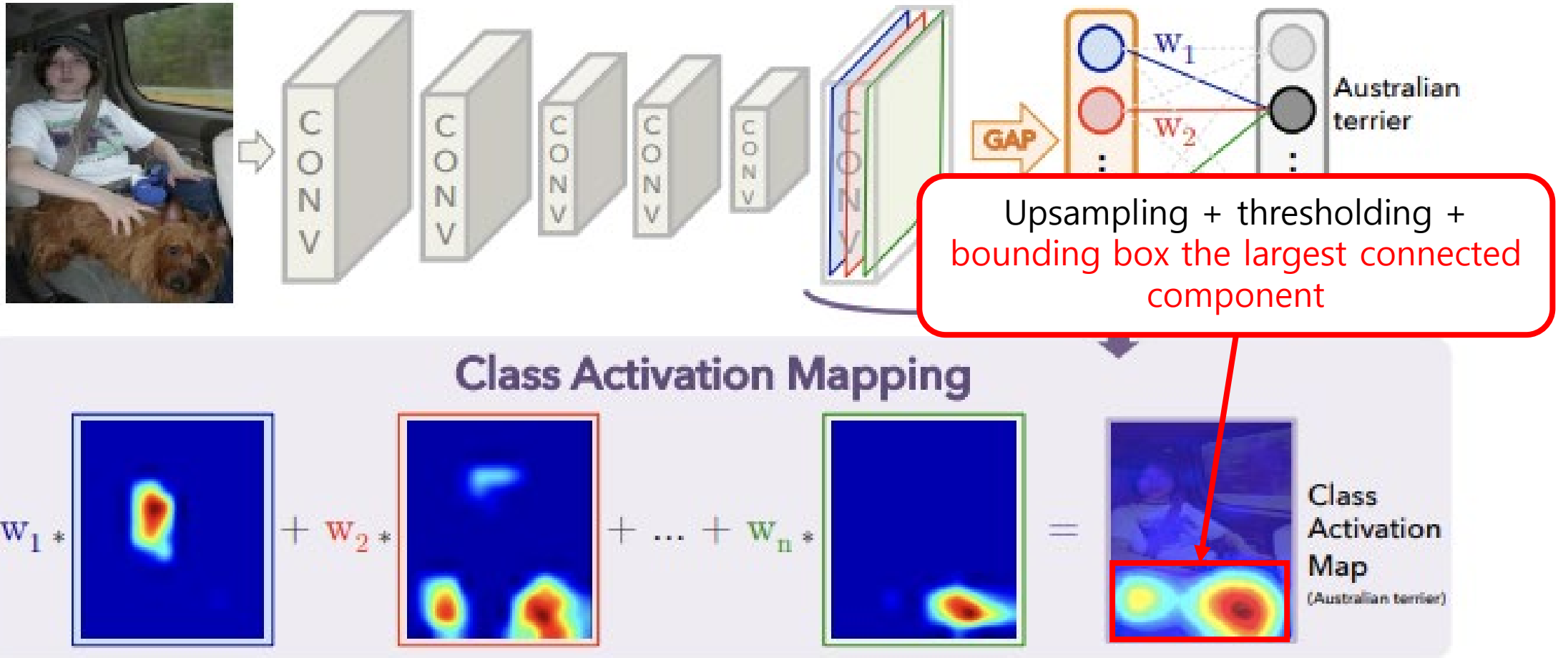
Upsampling + thresholding

$$M_c(x, y) = \sum_k w_k^c f_k(x, y).$$

Class Activation Mapping



# Localization



# GAP vs GMP

GAP

Encourages the network to identify the **full extent** of the image

The average of a map is maximized by finding all discriminative part of the object

GMP

Encourages the network to identify **just one** discriminative part of image

The max of a map is dependent only on the most discriminative part of the object

# Weakly vs Fully Supervised

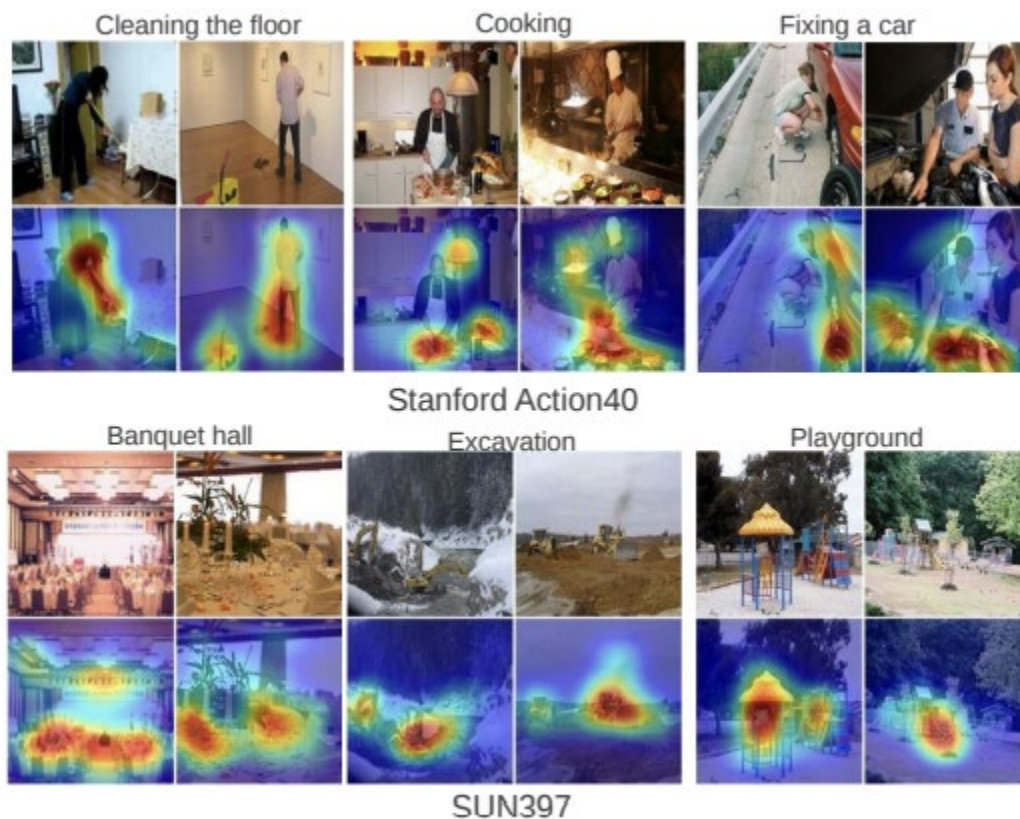
Exactly what is weakly vs fully supervised?

Table 3. Localization error on the ILSVRC test set for various weakly- and fully- supervised methods.

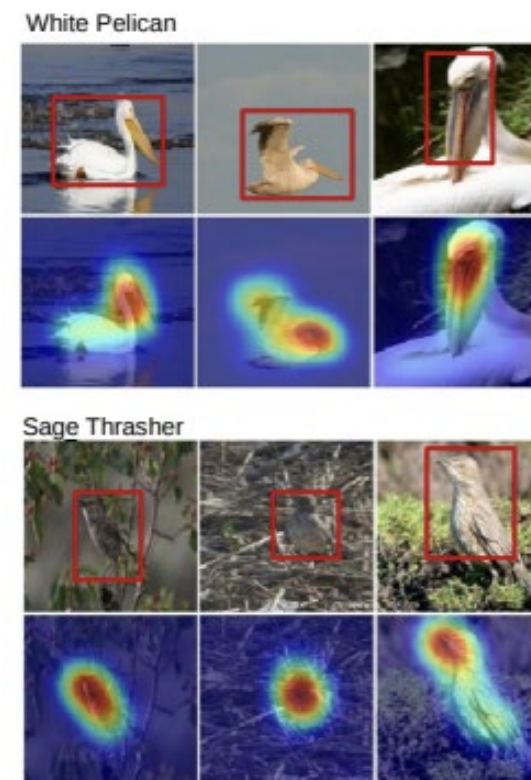
Method	supervision	top-5 test error
GoogLeNet-GAP (heuristics)	weakly	<b>37.1</b>
GoogLeNet-GAP	weakly	42.9
Backprop [22]	weakly	46.4
GoogLeNet [24]	full	26.7
OverFeat [21]	full	29.9
AlexNet [24]	full	34.2

# And more...

Generating localizable deep features for generic tasks



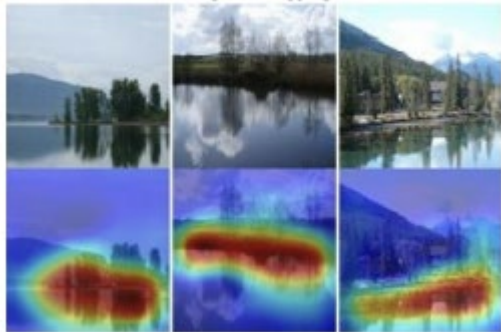
Fine grained recognition



# And more...

Concept localization

mirror in lake



Text detector



Visual question answering

