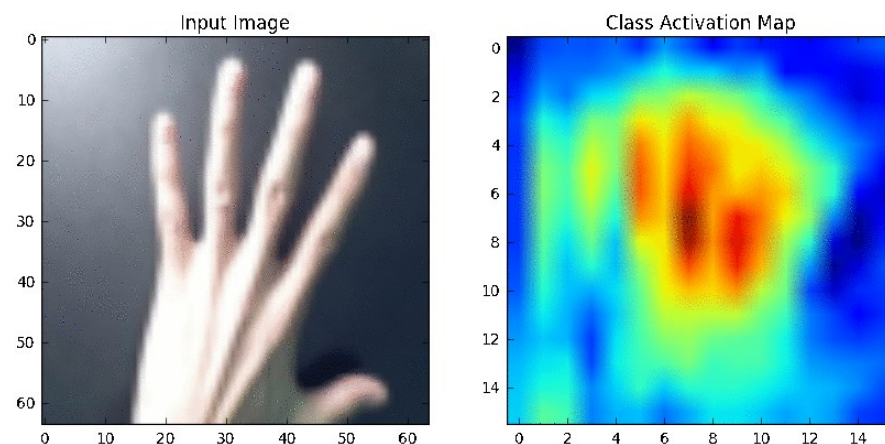


Visualization of CNN

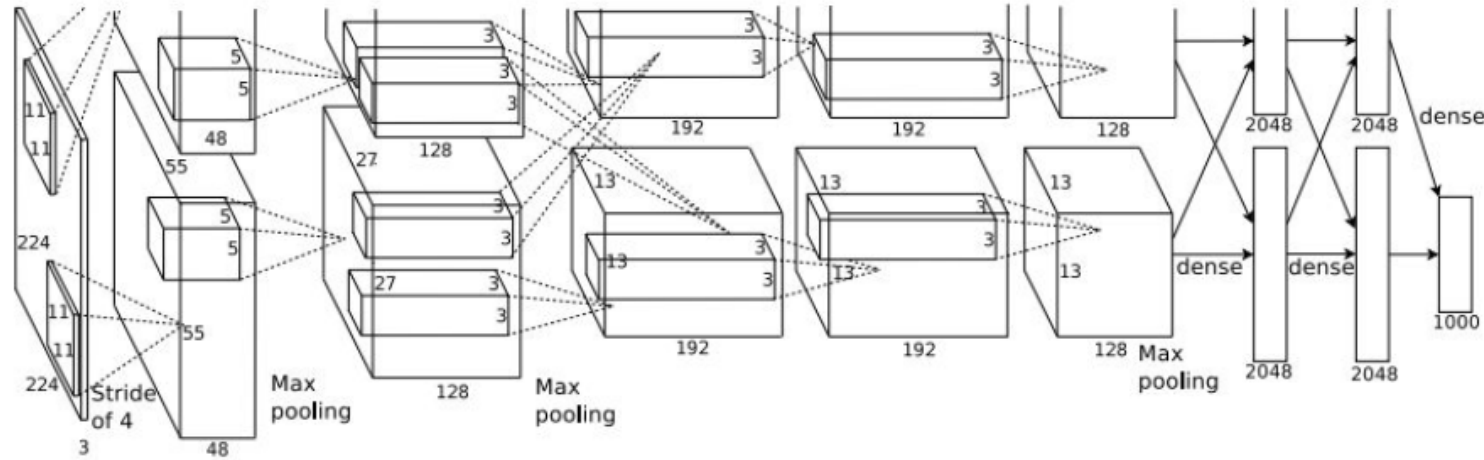


What's going on inside CNN?

This image is CC0 public domain



Input Image:
3 x 224 x 224



Class Scores:
1000 numbers

↑ ↑ ↑ ↑ ↑ ↑ ↑
What are the intermediate features looking for?

Visualize Patches that Maximally Activate Neurons

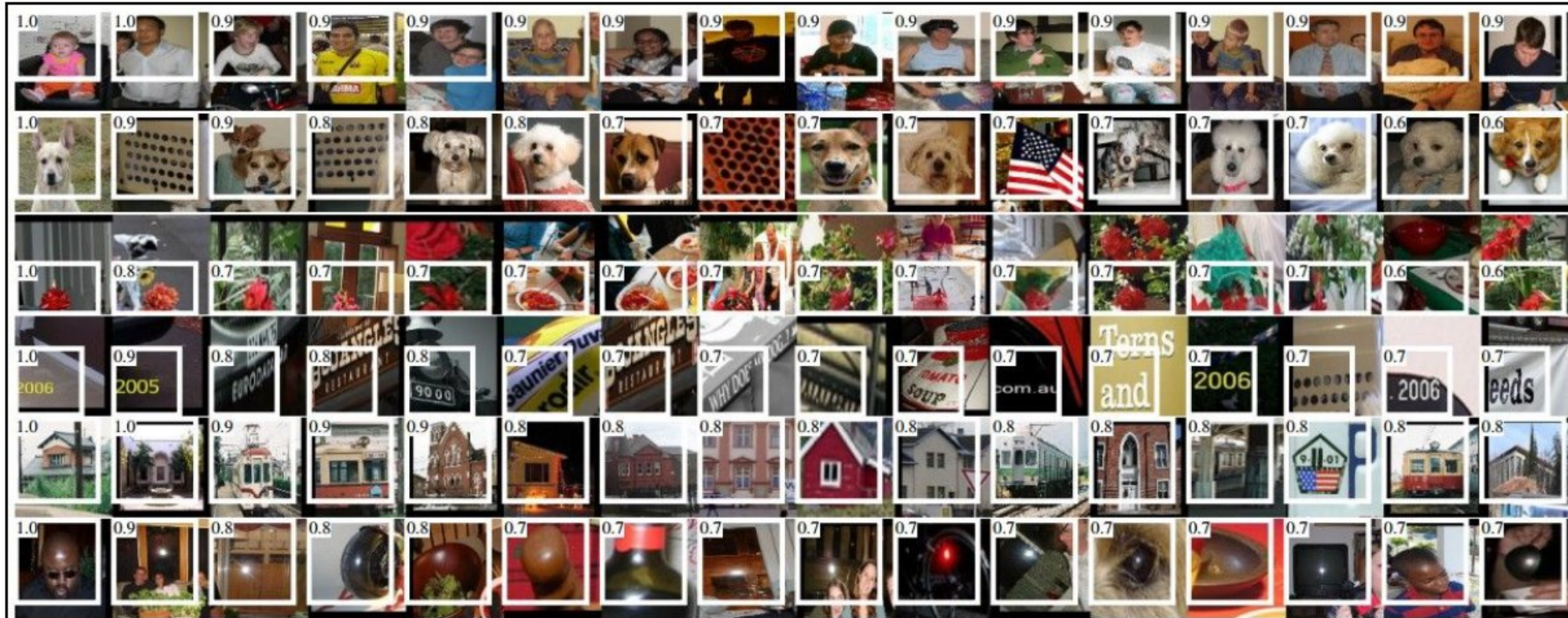
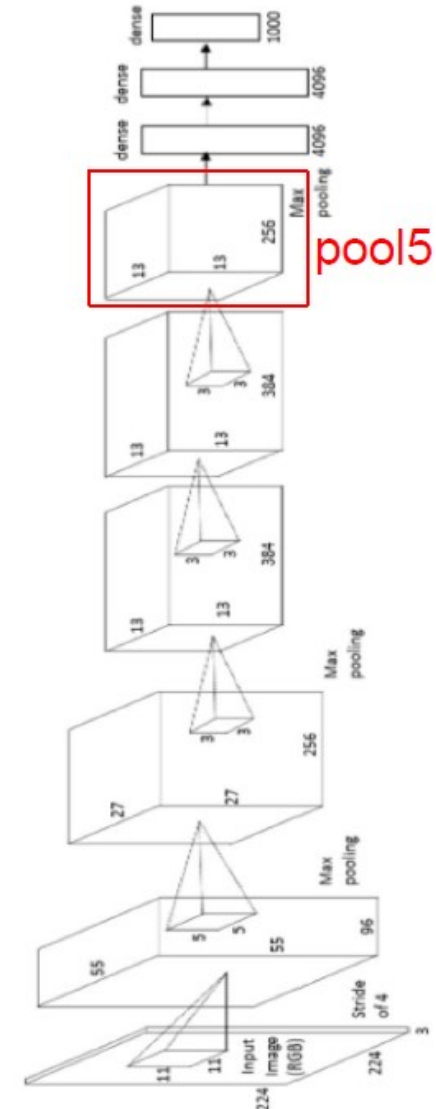


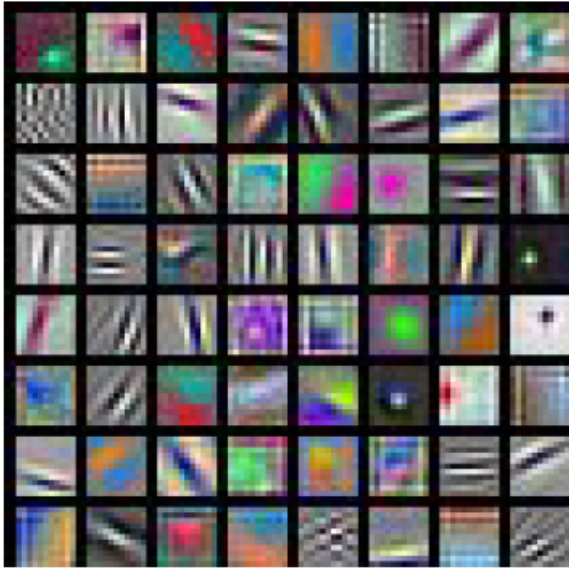
Figure 4: Top regions for six pool_5 units. Receptive fields and activation values are drawn in white. Some units are aligned to concepts, such as people (row 1) or text (4). Other units capture texture and material properties, such as dot arrays (2) and specular reflections (6).

Rich feature hierarchies for accurate object detection and semantic segmentation
[Girshick, Donahue, Darrell, Malik]

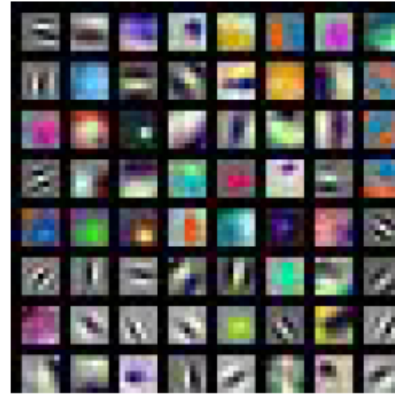


Visualize Filters

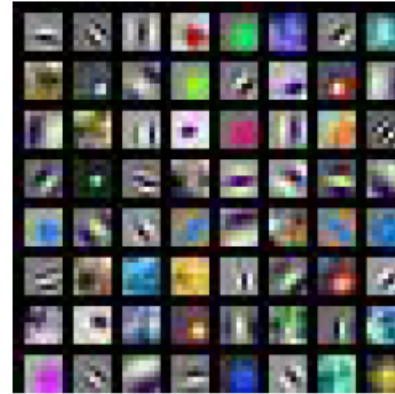
- Only interpretable on the first layer



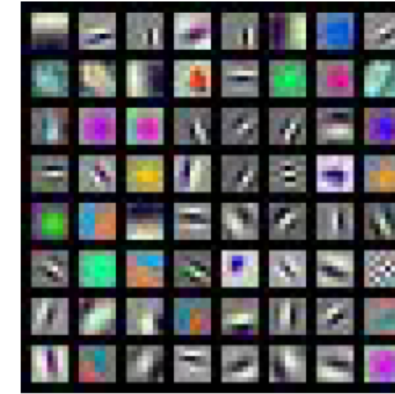
AlexNet:
 $64 \times 3 \times 11 \times 11$



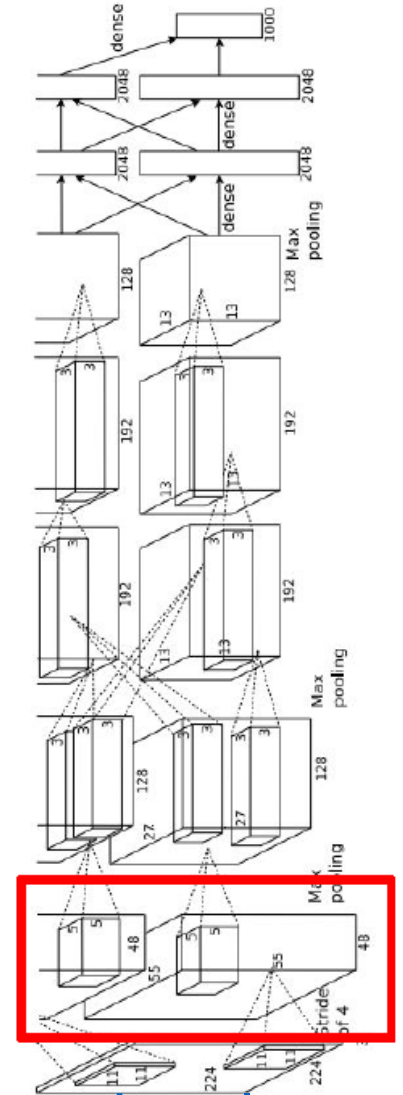
ResNet-18:
 $64 \times 3 \times 7 \times 7$



ResNet-101:
 $64 \times 3 \times 7 \times 7$



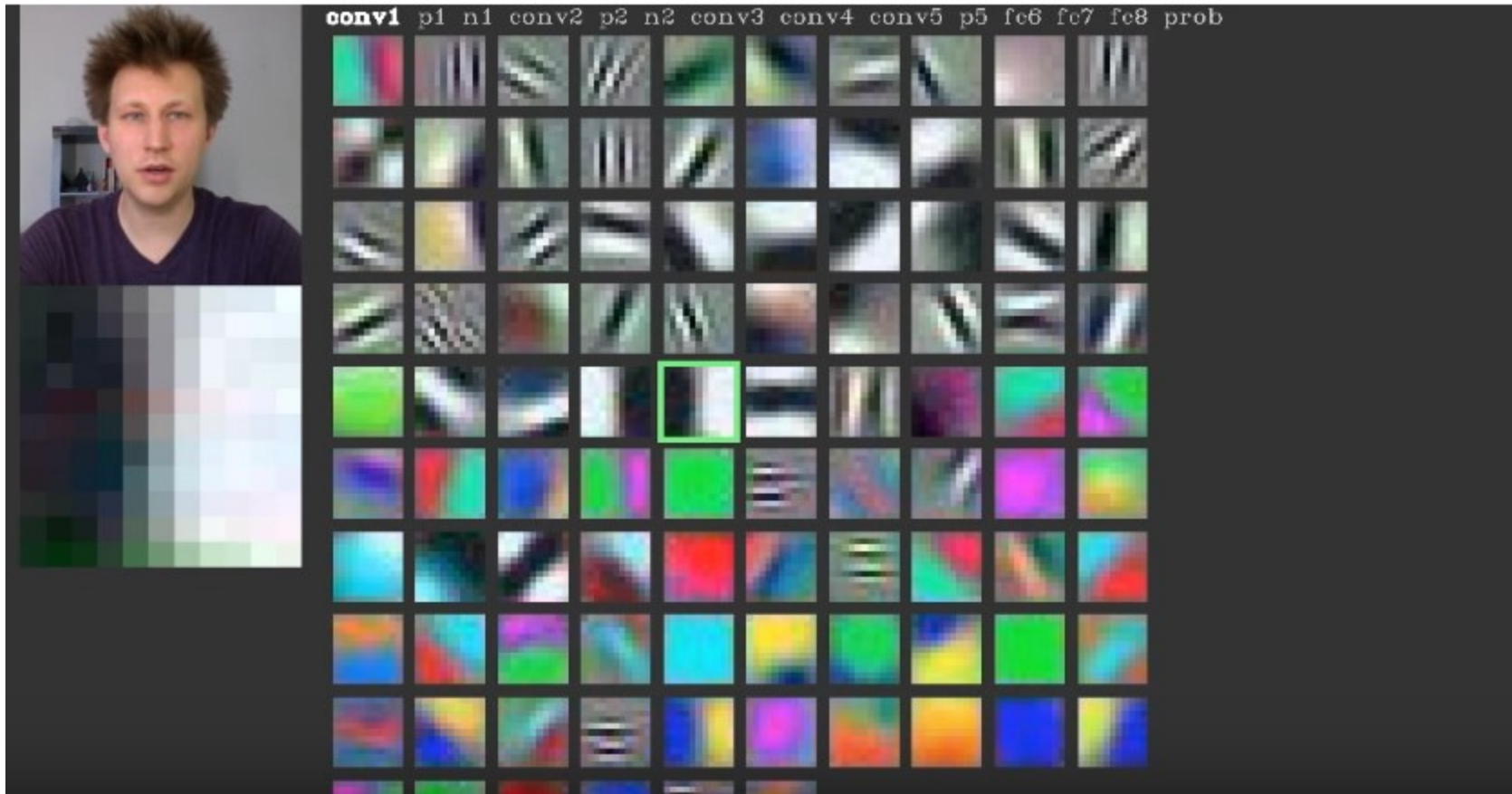
DenseNet-121:
 $64 \times 3 \times 7 \times 7$



- <http://cs.stanford.edu/people/karpathy/convnetjs/demo/cifar10.html>

Visualizing Activations

- <https://www.youtube.com/watch?v=AgkflQ4lGaM>

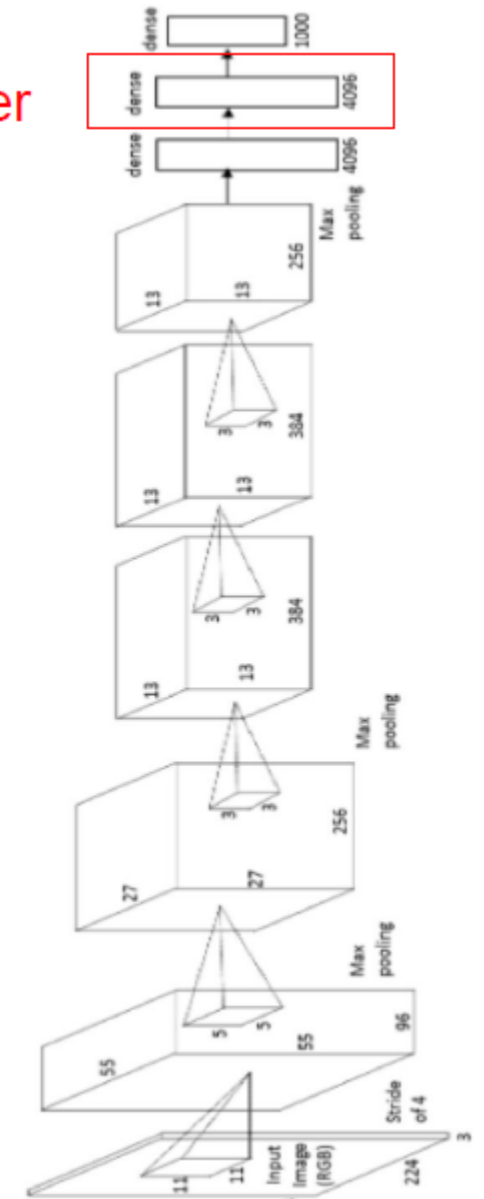


Visualizing the Representation

fc7 layer

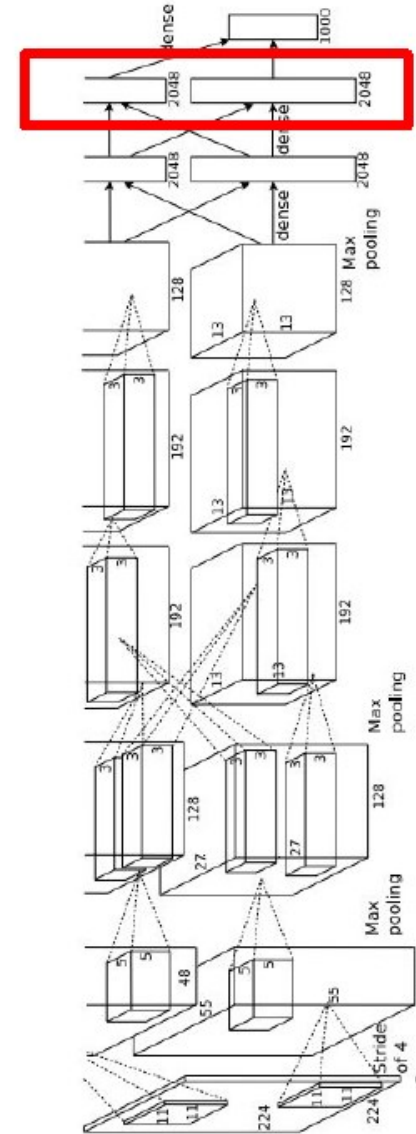
4096-dimensional “code” for an image
(layer immediately before the classifier)

can collect the code for many images



Last Layer : Nearest Neighbors

4096-dim vector



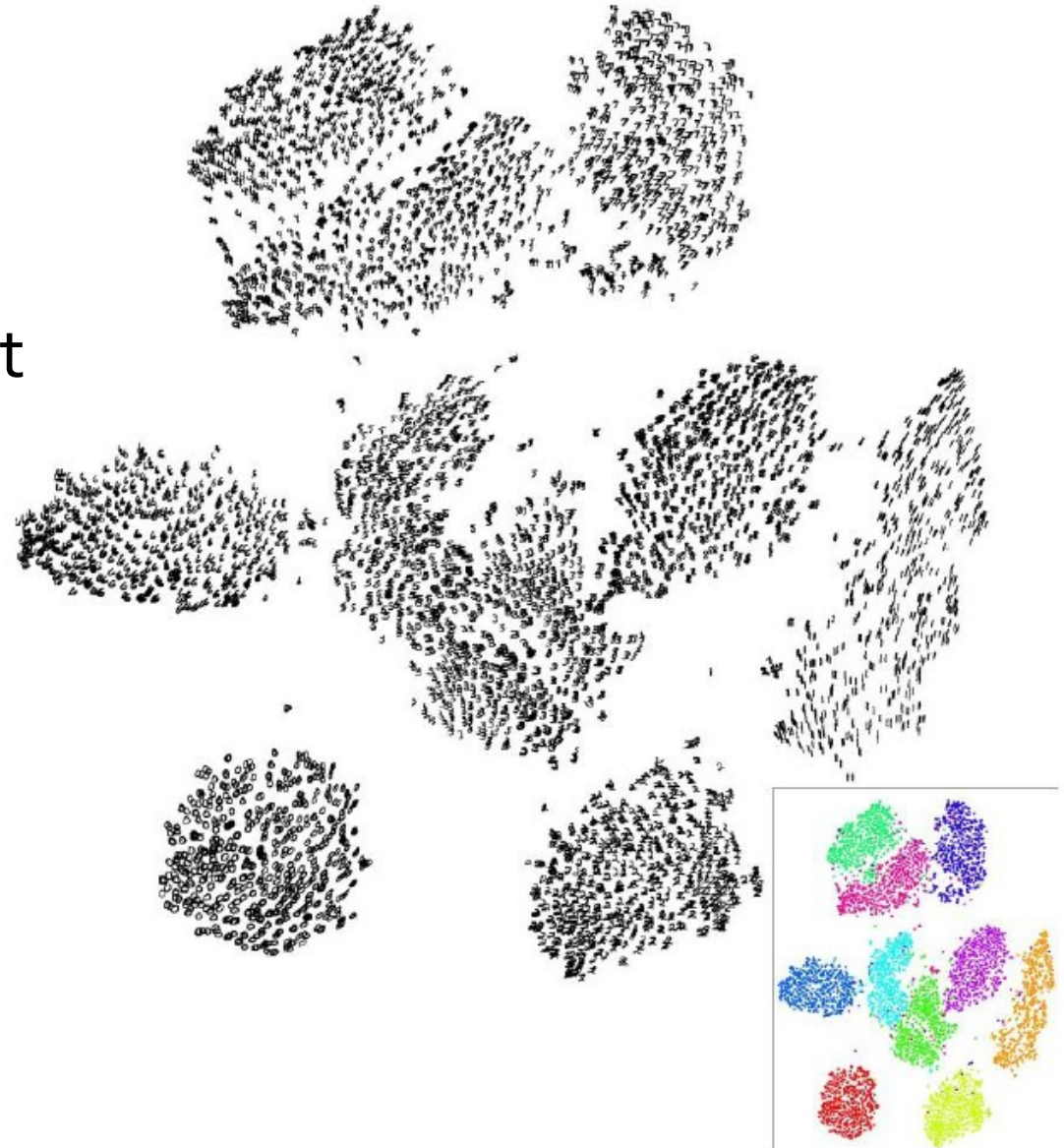
Test image L2 Nearest neighbors in feature space

Recall: Nearest neighbors in pixel space

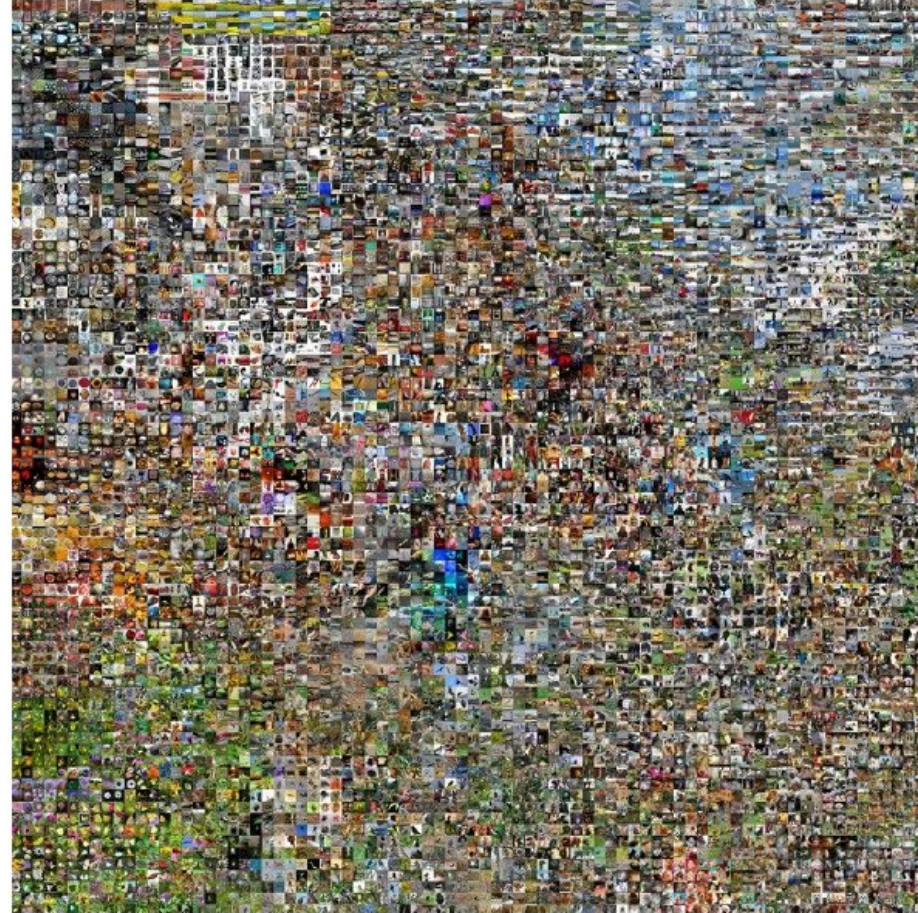


Last Layer: Dimensionality Reduction

- Visualize the “space” of FC7 feature vectors by reducing dimensionality of vectors from 4096 to 2 dimensions
- Simple algorithm: Principle Component Analysis(PCA)
- More complex: t-SNE

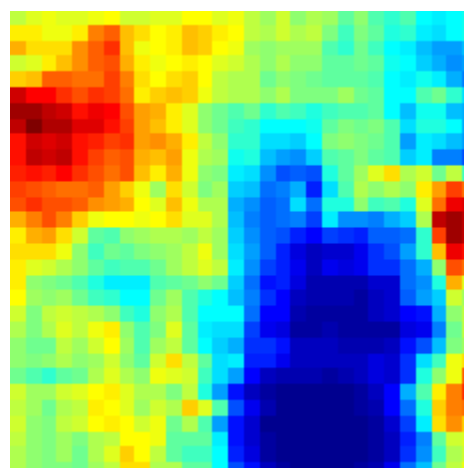
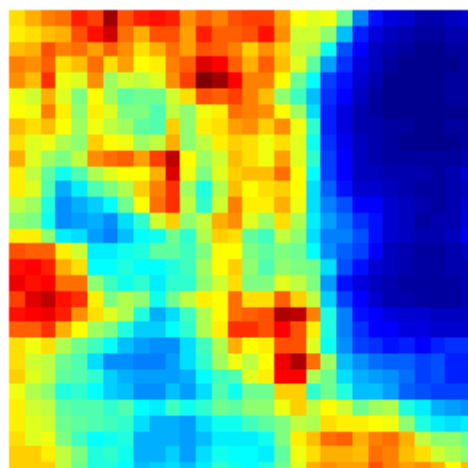
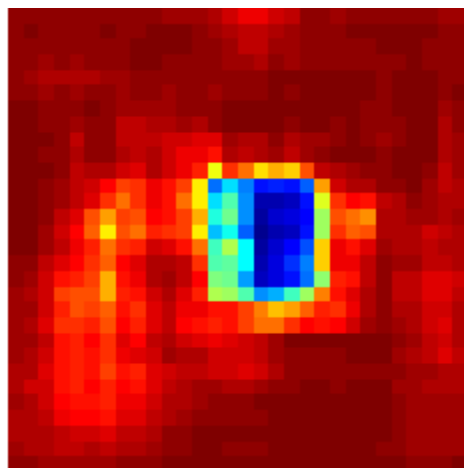
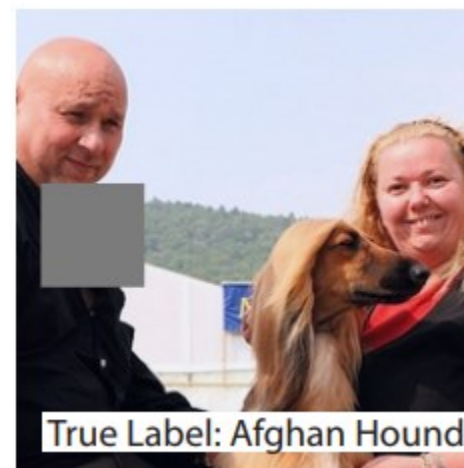


Last Layer: Dimensionality Reduction

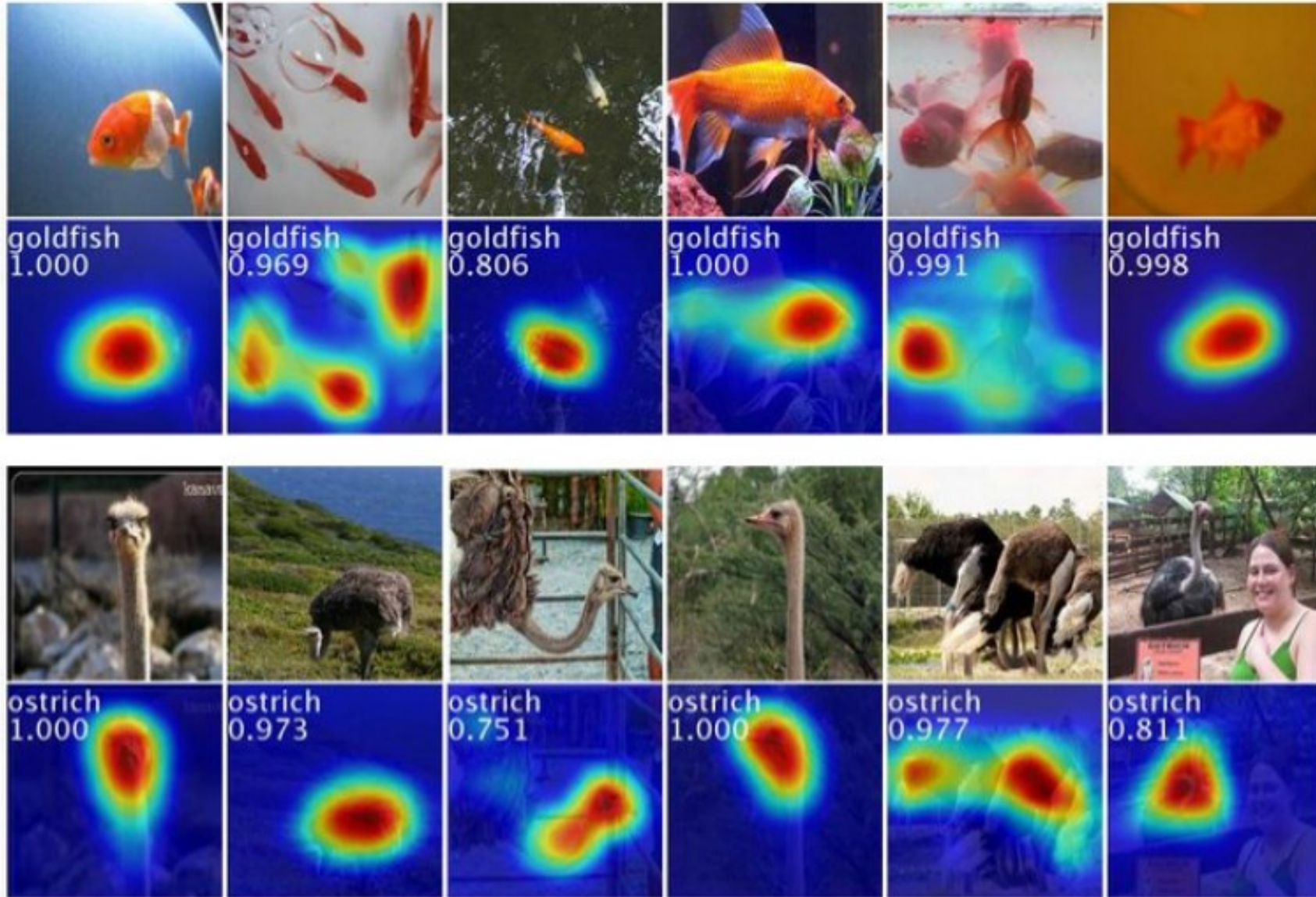


- <http://cs.stanford.edu/people/karpathy/cnnembed/>

Occlusion Experiments



Weakly Supervised Learning



Class activation map (CAM)

- **Identify important image regions** by projecting back the weights of output layer to convolutional feature maps.
- CAMs can be generated for each class in single image.
- Regions for each categories are different in given image.
 - palace, dome, church ...

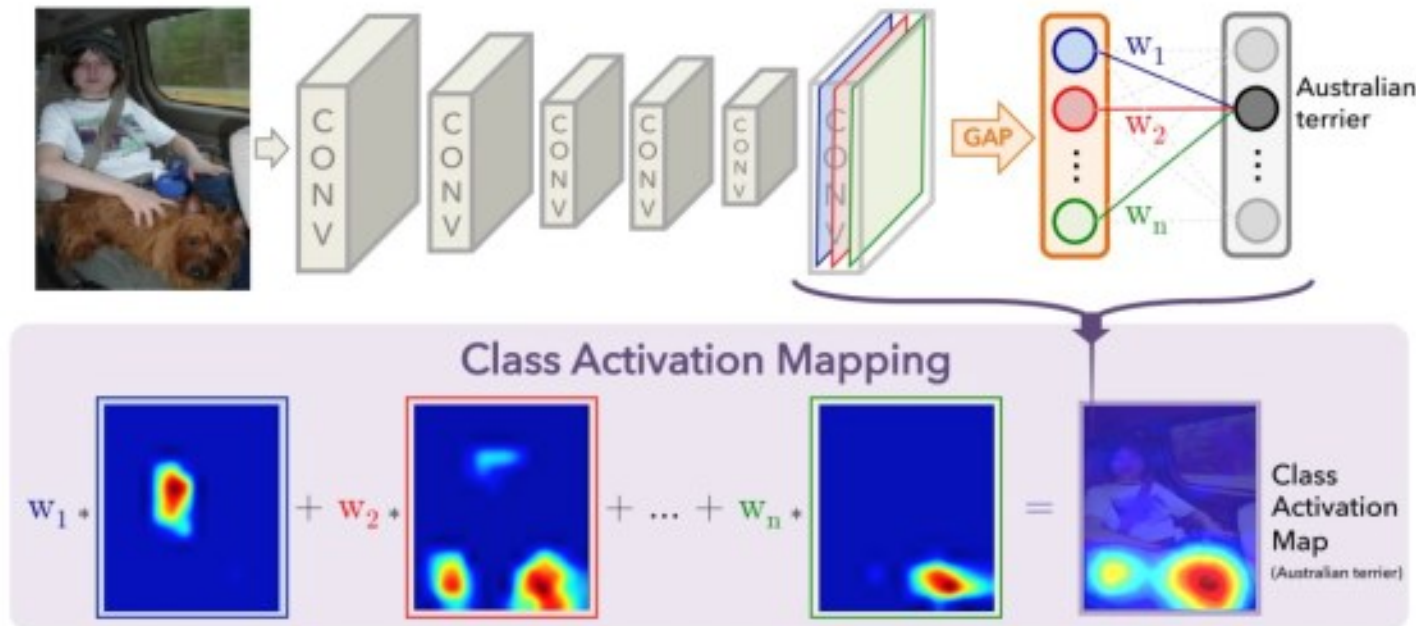


Figure 2. Class Activation Mapping: the predicted class score is mapped back to the previous convolutional layer to generate the class activation maps (CAMs). The CAM highlights the class-specific discriminative regions.

Results

- CAM on top 5 predictions on an image
- CAM for one object class in images

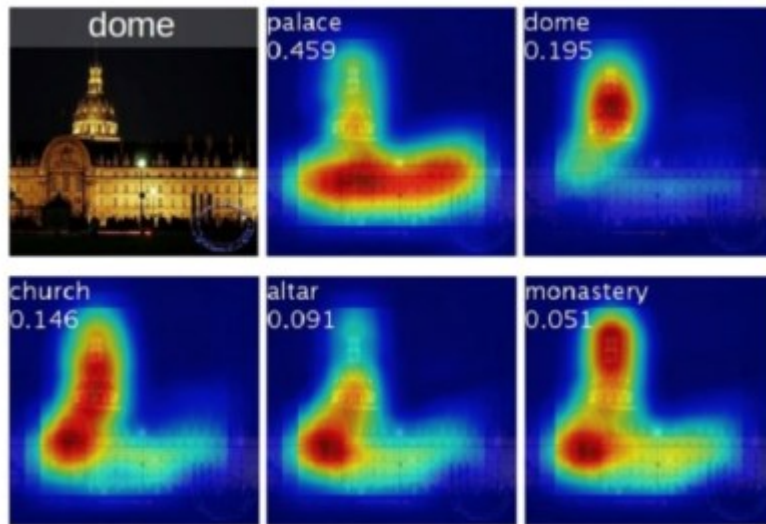


Figure 4. Examples of the CAMs generated from the top 5 predicted categories for the given image with ground-truth as dome. The predicted class and its score are shown above each class activation map. We observe that the highlighted regions vary across predicted classes e.g., *dome* activates the upper round part while *palace* activates the lower flat part of the compound.

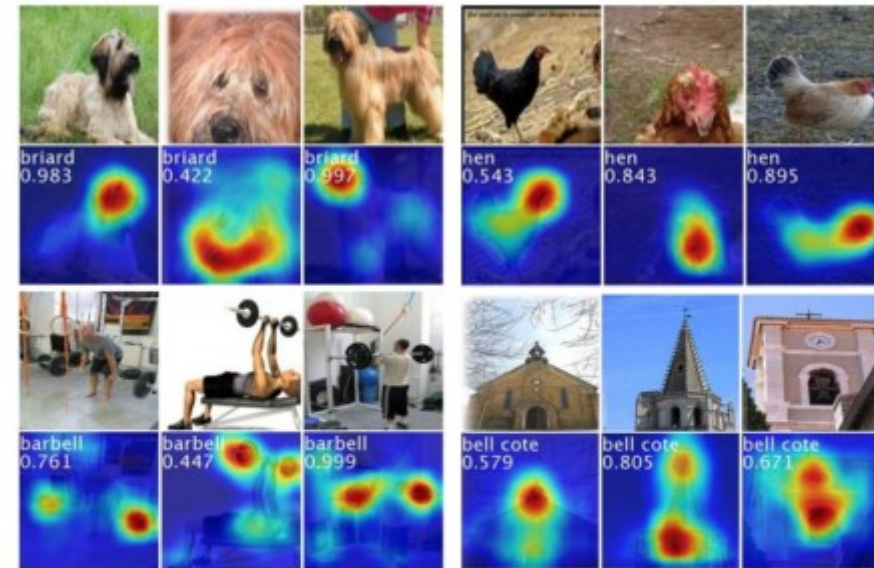
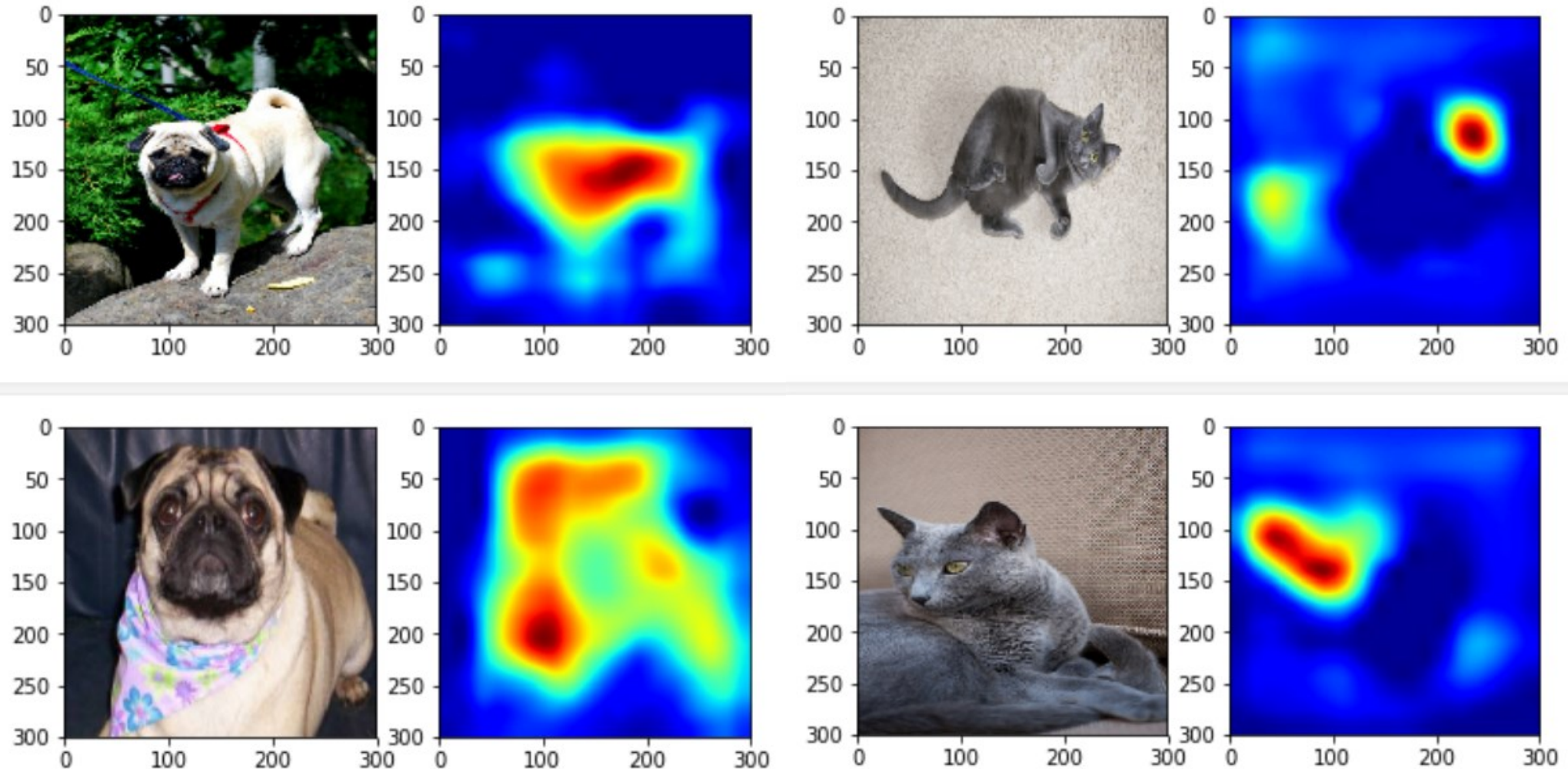


Figure 3. The CAMs of four classes from ILSVRC [20]. The maps highlight the discriminative image regions used for image classification e.g., the head of the animal for *briard* and *hen*, the plates in *barbell*, and the bell in *bell cote*.

Weakness of CAM (Weakly Supervised Localicztion)

- Focusing on discriminative features



Weakness of CAM (Weakly Supervised Localicztion)

- Focusing on discriminative features

