



# **Interpretable Convolutional Neural Networks**

Quanshi Zhang, Ying Nian Wu, and Song-Chun Zhu  
University of California, Los Angeles

# Outline

- **Background**
- **Motivation & Goal**
- **Analysis**
- **Proposed Method**
- **Experiments**
- **Conclusion**

# Background



# Motivation & Goal

## **Motivation:**

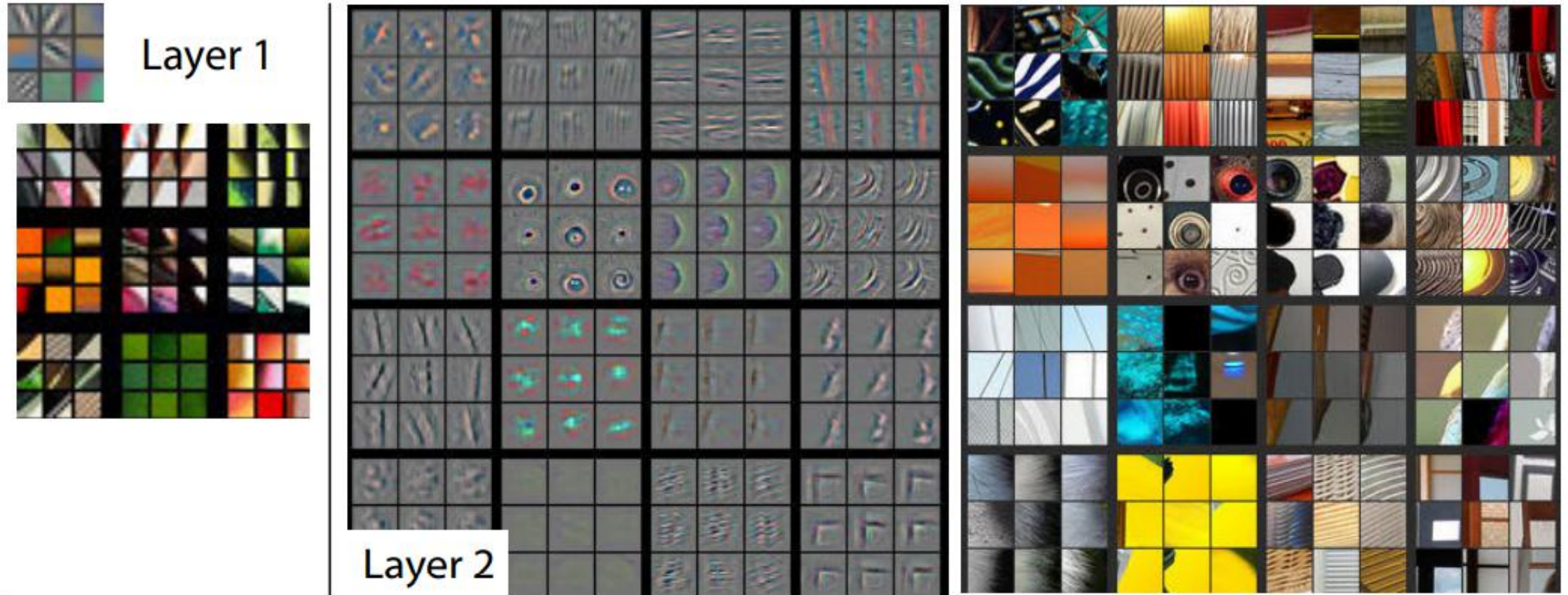
Without any additional human supervision, can we modify a CNN to make its conv-layers obtain interpretable knowledge representations?

## **Goal:**

1. We slightly revise a CNN to improve its interpretability, which can be broadly applied to CNNs with different structures.
2. We learn interpretable filters for a CNN without any additional annotations of object parts or textures for supervision. Training samples are the same as the original CNN.
3. We do not hope the network interpretability greatly affects the discrimination power.

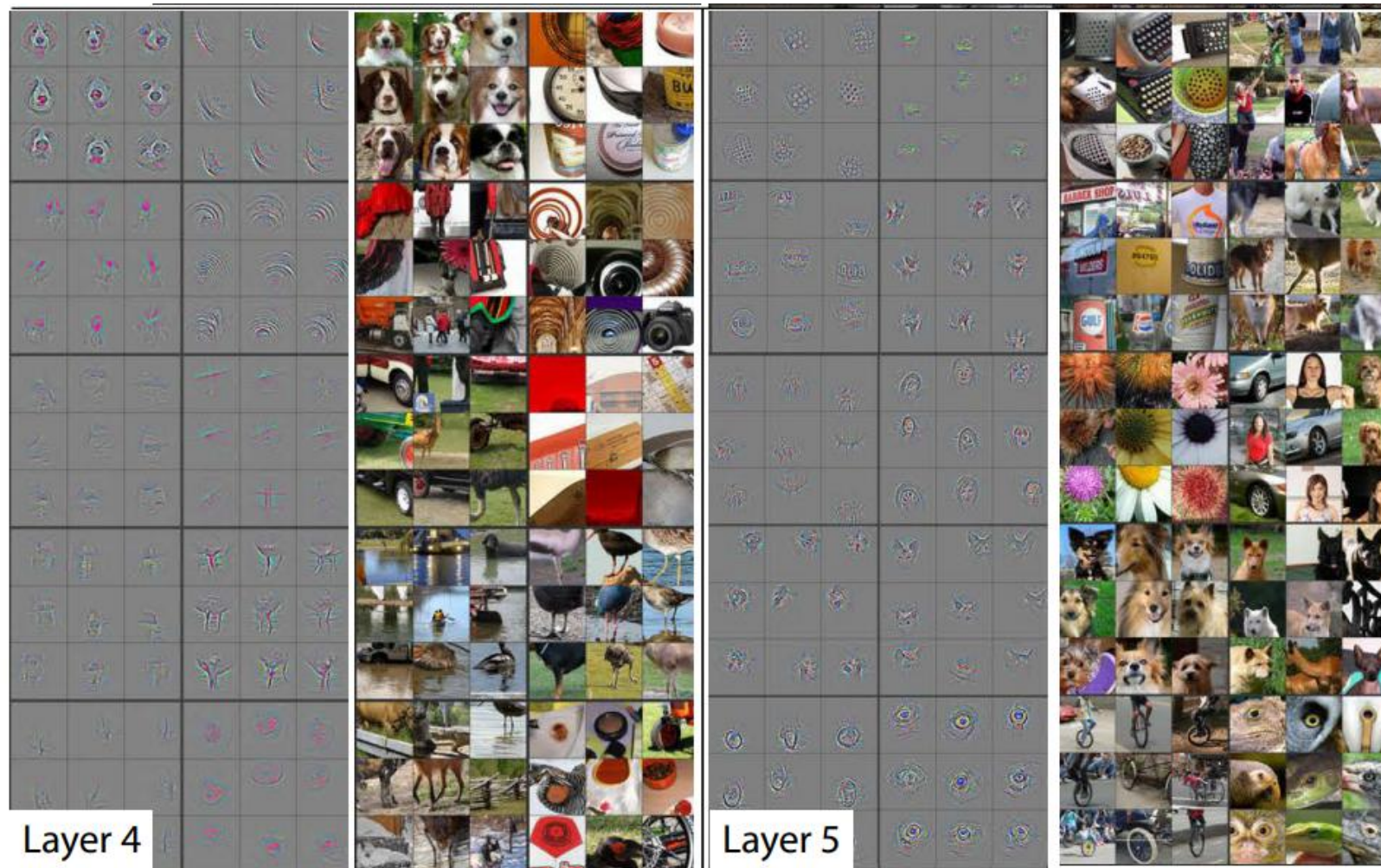


# Analysis



**Filters in low conv-layers usually describe simple textures.**

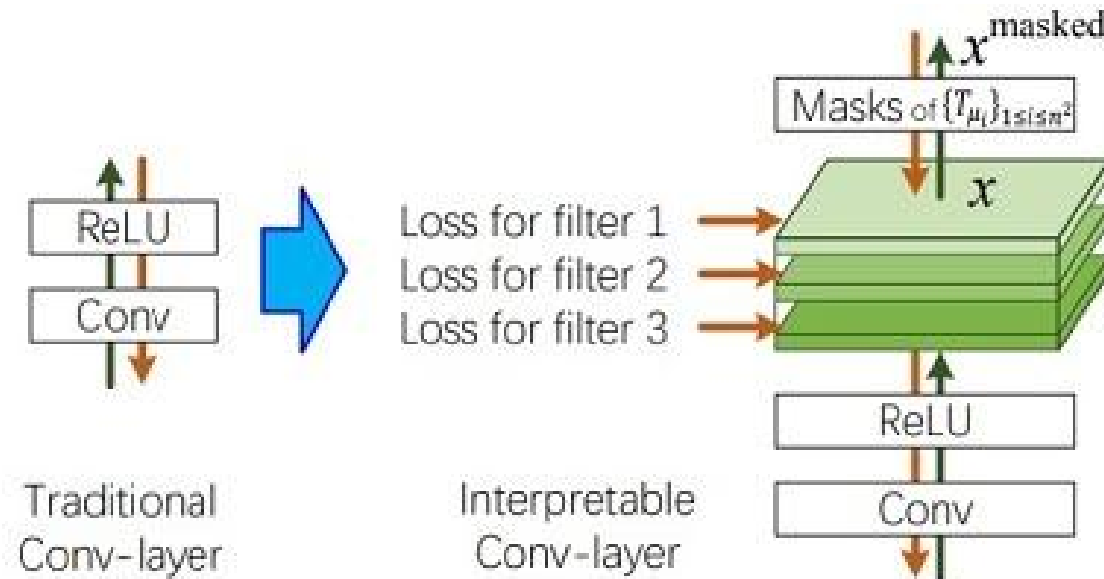
# Analysis



**Filters in high conv-layers are more likely to represent object parts.**

# Proposed Method

This paper proposes a simple yet effective loss to push a filter in a specific conv-layer of a CNN towards the representation of an object part.





# Proposed Method

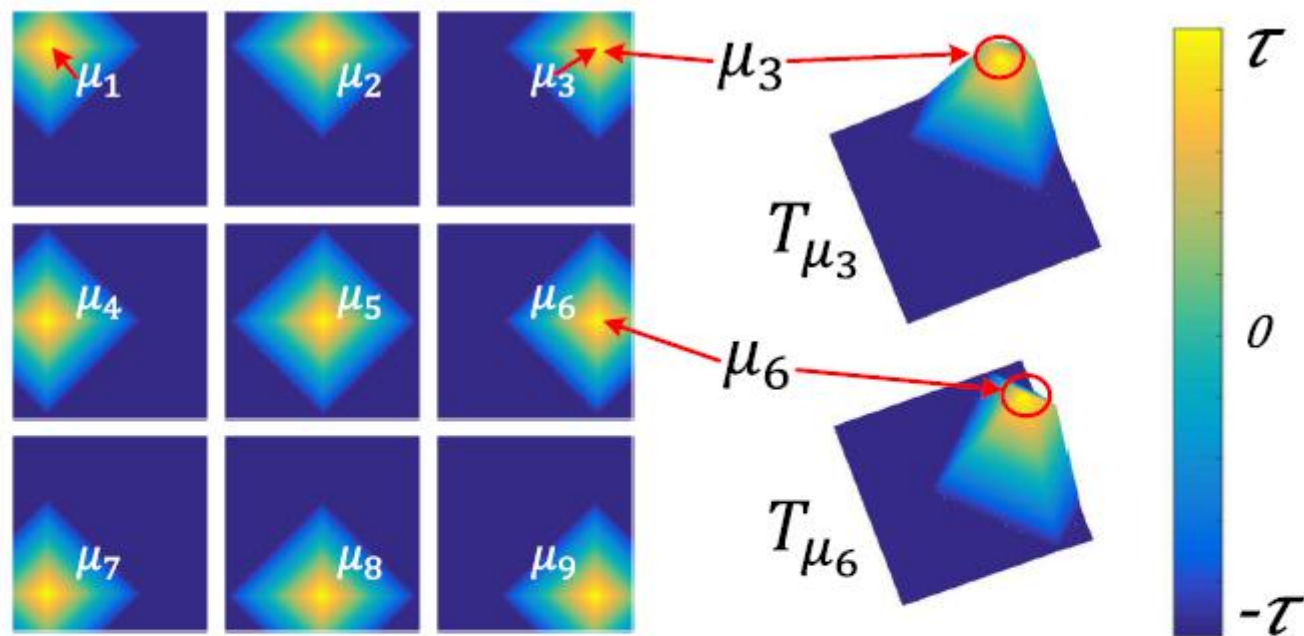


Figure 3. Templates of  $T_{\mu_i}$ . Each template  $T_{\mu_i}$  matches to a feature map  $x$  when the target part mainly triggers the  $i$ -th unit in  $x$ . In fact, the algorithm also supports a round template based on the L-2 norm distance. Here, we use the L-1 norm distance instead to speed up the computation.



# Proposed Method

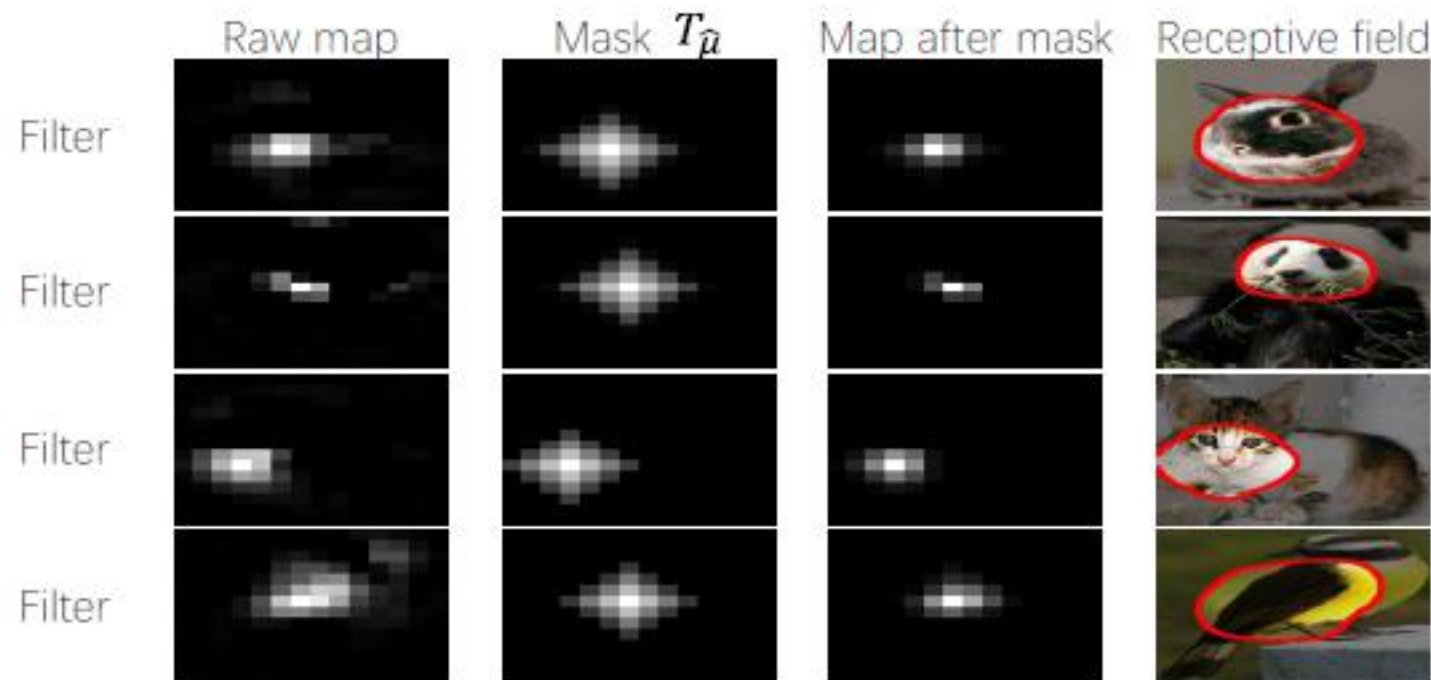


Figure 4. Given an input image  $I$ , from the left to the right, we consequently show the feature map of a filter after the ReLU layer  $x$ , the assigned mask  $T_{\hat{\mu}}$ , the masked feature map  $x^{\text{masked}}$ , and the image-resolution RF of activations in  $x^{\text{masked}}$  computed by [40].

# Proposed Method

## Filter Loss Function:

$$\text{Loss}_f = -MI(\mathbf{X}; \mathbf{T}) = - \sum_T p(T) \sum_x p(x|T) \log \frac{p(x|T)}{p(x)}$$

$MI(\cdot)$  denotes the mutual information

$$\forall T \in \mathbf{T}, \quad p(x|T) = \frac{1}{Z_T} \exp [\text{tr}(x \cdot T)]$$

$$Z_T = \sum_{x \in \mathbf{X}} \exp(\text{tr}(x \cdot T))$$

$$\text{tr}(x \cdot T) = \sum_{ij} x_{ij} t_{ij}. \quad p(x) = \sum_T p(T) p(x|T)$$

$$p(x) = \sum_T p(T) p(x|T)$$

注：最大化mutual information，因此loss function前面有负号

Mutual Information [https://www.wikiwand.com/en/Mutual\\_information](https://www.wikiwand.com/en/Mutual_information)

# Experiments

## Average part interpretability

|                        | bird         | cat          | cow          | dog          | horse        | sheep        | Avg.         |
|------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| AlexNet                | 0.332        | 0.363        | 0.340        | 0.374        | 0.308        | 0.373        | 0.348        |
| AlexNet, interpretable | <b>0.770</b> | <b>0.565</b> | <b>0.618</b> | <b>0.571</b> | <b>0.729</b> | <b>0.669</b> | <b>0.654</b> |
| VGG-16                 | 0.519        | 0.458        | 0.479        | 0.534        | 0.440        | 0.542        | 0.495        |
| VGG-16, interpretable  | <b>0.818</b> | <b>0.653</b> | <b>0.683</b> | <b>0.900</b> | <b>0.795</b> | <b>0.772</b> | <b>0.770</b> |
| VGG-M                  | 0.357        | 0.365        | 0.347        | 0.368        | 0.331        | 0.373        | 0.357        |
| VGG-M, interpretable   | <b>0.821</b> | <b>0.632</b> | <b>0.634</b> | <b>0.669</b> | <b>0.736</b> | <b>0.756</b> | <b>0.708</b> |
| VGG-S                  | 0.251        | 0.269        | 0.235        | 0.275        | 0.223        | <b>0.287</b> | 0.257        |
| VGG-S, interpretable   | <b>0.526</b> | <b>0.366</b> | <b>0.291</b> | <b>0.432</b> | <b>0.478</b> | 0.251        | <b>0.390</b> |

Table 1. Part interpretability of filters in CNNs for single-category classification based on the Pascal VOC Part dataset [3].

| Network               | Logistic log loss <sup>4</sup> | Softmax log loss |
|-----------------------|--------------------------------|------------------|
| VGG-16                | 0.710                          | 0.723            |
| VGG-16, interpretable | <b>0.938</b>                   | <b>0.897</b>     |
| VGG-M                 | 0.478                          | 0.502            |
| VGG-M, interpretable  | <b>0.770</b>                   | <b>0.734</b>     |
| VGG-S                 | 0.479                          | 0.435            |
| VGG-S, interpretable  | <b>0.572</b>                   | <b>0.601</b>     |

Table 2. Part interpretability of filters in CNNs that are trained for multi-category classification based on the VOC Part dataset [3]. Filters in our interpretable CNNs exhibited significantly better part interpretability than ordinary CNNs in all comparisons.



# Experiments

## Location instability

|                        | bird         | cat          | cow          | dog          | horse        | sheep        | Avg.         |
|------------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| AlexNet                | 0.153        | 0.131        | 0.141        | 0.128        | 0.145        | 0.140        | 0.140        |
| AlexNet, interpretable | <b>0.090</b> | <b>0.089</b> | <b>0.090</b> | <b>0.088</b> | <b>0.087</b> | <b>0.088</b> | <b>0.088</b> |
| VGG-16                 | 0.145        | 0.133        | 0.146        | 0.127        | 0.143        | 0.143        | 0.139        |
| VGG-16, interpretable  | <b>0.101</b> | <b>0.098</b> | <b>0.105</b> | <b>0.074</b> | <b>0.097</b> | <b>0.100</b> | <b>0.096</b> |
| VGG-M                  | 0.152        | 0.132        | 0.143        | 0.130        | 0.145        | 0.141        | 0.141        |
| VGG-M, interpretable   | <b>0.086</b> | <b>0.094</b> | <b>0.090</b> | <b>0.087</b> | <b>0.084</b> | <b>0.084</b> | <b>0.088</b> |
| VGG-S                  | 0.152        | 0.131        | 0.141        | 0.128        | 0.144        | 0.141        | 0.139        |
| VGG-S, interpretable   | <b>0.089</b> | <b>0.092</b> | <b>0.092</b> | <b>0.087</b> | <b>0.086</b> | <b>0.088</b> | <b>0.089</b> |

Table 4. Location instability of filters ( $\mathbf{E}_{f,k}[D_{f,k}]$ ) in CNNs that are trained for single-category classification using the Pascal VOC Part dataset [3]. Filters in our interpretable CNNs exhibited significantly lower localization instability than ordinary CNNs.

| Network                | Avg. location instability |
|------------------------|---------------------------|
| AlexNet                | 0.150                     |
| AlexNet, interpretable | <b>0.070</b>              |
| VGG-16                 | 0.137                     |
| VGG-16, interpretable  | <b>0.076</b>              |
| VGG-M                  | 0.148                     |
| VGG-M, interpretable   | <b>0.065</b>              |
| VGG-S                  | 0.148                     |
| VGG-S, interpretable   | <b>0.073</b>              |

Table 5. Location instability of filters ( $\mathbf{E}_{f,k}[D_{f,k}]$ ) in CNNs for single-category classification using the CUB200-2011 dataset.



# Experiments

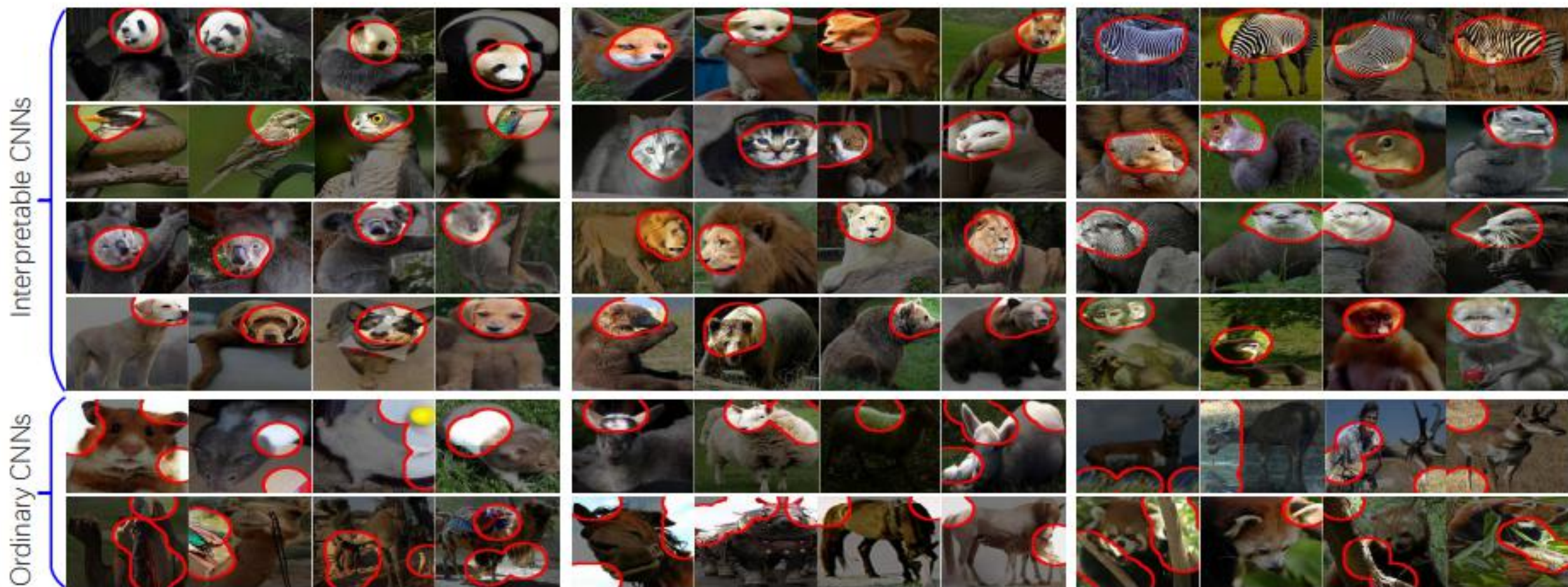
## Classification Accuracy

|               | multi-category        |                       |              | single-category |              |              |
|---------------|-----------------------|-----------------------|--------------|-----------------|--------------|--------------|
|               | ILSVRC Part           | VOC Part              |              | ILSVRC Part     | VOC Part     | CUB200       |
|               | logistic <sup>4</sup> | logistic <sup>4</sup> | softmax      |                 |              |              |
| AlexNet       | —                     | —                     | —            | <b>96.28</b>    | <b>95.40</b> | <b>95.59</b> |
| interpretable | —                     | —                     | —            | 95.38           | 93.93        | 95.35        |
| VGG-M         | 96.73                 | 93.88                 | 81.93        | <b>97.34</b>    | <b>96.82</b> | <b>97.34</b> |
| interpretable | <b>97.99</b>          | <b>96.19</b>          | <b>88.03</b> | 95.77           | 94.17        | 96.03        |
| VGG-S         | 96.98                 | 94.05                 | 78.15        | <b>97.62</b>    | <b>97.74</b> | <b>97.24</b> |
| interpretable | <b>98.72</b>          | <b>96.78</b>          | <b>86.13</b> | 95.64           | 95.47        | 95.82        |
| VGG-16        | —                     | 97.97                 | 89.71        | <b>98.58</b>    | <b>98.66</b> | <b>98.91</b> |
| interpretable | —                     | <b>98.50</b>          | <b>91.60</b> | 96.67           | 95.39        | 96.51        |

Table 7. Classification accuracy based on different datasets. In single-category classification, ordinary CNNs performed better, while in multi-category classification, interpretable CNNs exhibited superior performance.

# Experiments

## Performances



# Conclusion

1. This paper proposes a general method to enhance feature interpretability of CNNs.
2. We design a loss to push a filter in high conv-layers towards the representation of an object part during the learning process without any part annotations.