

---

# Explainable Convolutional Neural Network

---

발표자 : 백인성

# 목차

---

1. Introduction
2. [CNN]Convolutional Neural Network
3. [CAM]Class Activation Map
4. Interpretable CNN(Convolutional Neural Network)
5. Conclusion

---

# 1. Introduction

---

# Introduction

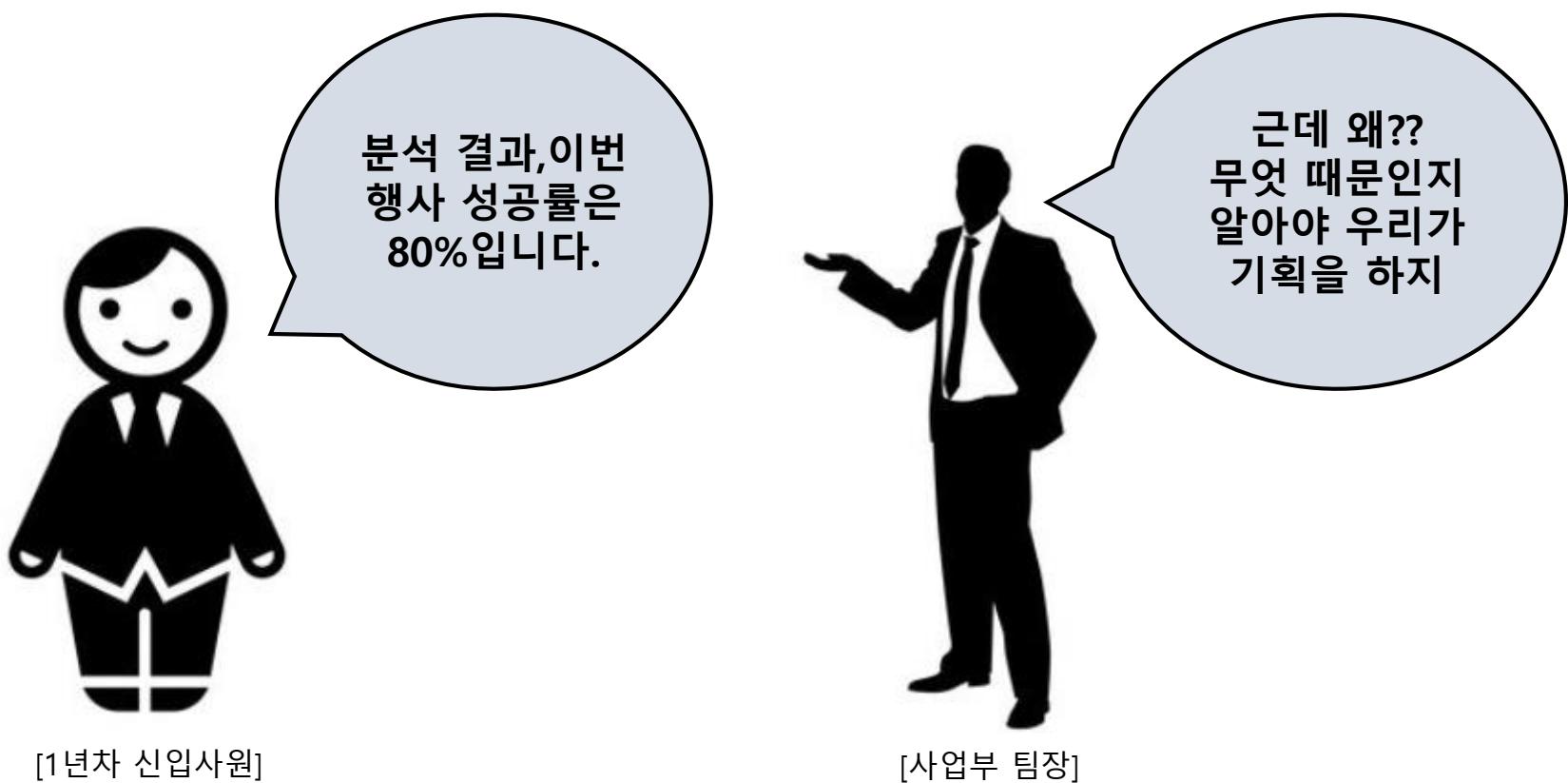
## ❖ 예측 모델에서 해석이 왜 중요한가?

- 모델링이 핵심인 전문가는 예측 성능 향상이 주 목표이지만
- 실제 적용되는 산업 분야에서는 해석은 예측 성능만큼 중요함



# Introduction

- ❖ 예측 모델에서 해석이 왜 중요한가?
  - 데이터 분석가는 모델의 결과를 중심으로 이야기 하지만
  - 실제로 사용하는 실무자들은 그 결과가 나온 이유를 더 궁금해 함



# Introduction

- ❖ 예측 모델에서 해석이 왜 필요한가?
  - 현업 관계자와 소통을 원활하게 하기 위해
  - 모델의 예측 성능에 대한 타당성을 직관적으로 보여주기 위해



[A사 팀장]

[B사 팀장]

# Introduction

---

- ❖ 예측 모델의 해석이 필요한 이유는 아래의 3가지 이유로 요약할 수 있음

**1. 예측 모델을 실제로 사용하는 사람들의 이해를 돋기 위해**

**2. 의사결정하는 사람들을 조금 더 쉽게 설득하기 위해**

**3. 구축한 예측 모델의 결과의 타당함을 직관적으로 보이기 위해**

# Introduction

## ❖ Interpretable Convolutional Neural Network(2018)

- CNN 모델의 결과를 해석하기 위한 방법론 중 하나
- 2018년도 CVPR에서 소개 되었고, 1년간 인용 회수는 65회

### Interpretable convolutional neural networks

[Q Zhang, Y Nian Wu, SC Zhu - Proceedings of the IEEE ...](#), 2018 - [openaccess.thecvf.com](#)

This paper proposes a method to modify a traditional convolutional neural network (CNN) into an interpretable CNN, in order to clarify knowledge representations in high conv-layers of the CNN. In an interpretable CNN, each filter in a high conv-layer represents a specific object part. Our interpretable CNNs use the same training data as ordinary CNNs without a need for any annotations of object parts or textures for supervision. The interpretable CNN automatically assigns each filter in a high conv-layer with an object part during the learning ...

☆ 99 65회 인용 관련 학술자료 전체 7개의 버전 >>

### Interpretable Convolutional Neural Networks

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

#### Abstract

*This paper proposes a method to modify a traditional convolutional neural network (CNN) into an interpretable CNN, in order to clarify knowledge representations in high conv-layers of the CNN. In an interpretable CNN, each filter in a high conv-layer represents a specific object part. Our interpretable CNNs use the same training data as ordinary CNNs without a need for any annotations of object parts or textures for supervision. The interpretable CNN automatically assigns each filter in a high conv-layer with an object part during the learning process. We can apply our method to different types of CNNs with various structures. The explicit knowledge representation in an interpretable CNN can help people understand the logic inside a CNN, i.e. what patterns are memorized by the CNN for prediction. Experiments have shown that filters in an interpretable CNN are more semantically meaningful than those in a traditional CNN. The code is available at <https://github.com/zqs1022/interpretableCNN>.*

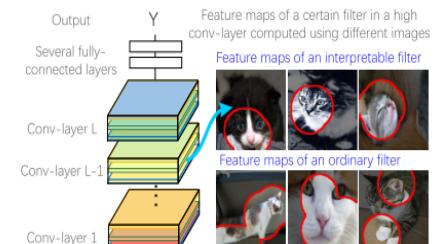


Figure 1. Comparison of a filter's feature maps in an interpretable CNN and those in a traditional CNN.

In fact, we can roughly consider the first two semantics as object-part patterns with specific shapes, and summarize the last four semantics as texture patterns without clear contours. Moreover, filters in low conv-layers usually describe simple textures, whereas filters in high conv-layers are more likely to represent object parts.

Therefore, in this study, we aim to train each filter in a

# Introduction

## ❖ 논문의 저자 (Quanshi Zhang)

- 상하이 대학교 부교수
- 딥러닝 모델의 해석력을 높이려는 연구를 중점적으로 진행



Quanshi Zhang

张 奎 石

Associate Professor

John Hopcroft Center for Computer Science  
School of electronic information and electrical  
engineering

Shanghai Jiao Tong University

Email: zqs1022 [AT] sjtu.edu.cn [\[知乎\]](#)

### 招生 Prospective Ph.D., Master, and undergraduate students:

I am looking for highly motivated students to work together on interpretability of neural networks, unsupervised and weakly-supervised learning, graph mining, and other frontier topics in machine learning and computer vision.

Please read "[写给学生](#)" and send me your CV and transcripts.

#### 제목

[Interpretable convolutional neural networks](#)

Q Zhang, Y Nian Wu, SC Zhu

Proceedings of the IEEE Conference on Computer Vision and Pattern ...

[Visual interpretability for deep learning: a survey](#)

Q Zhang, SC Zhu

Frontiers of Information Technology & Electronic Engineering 19 (1), 27-39

[Prediction of human emergency behavior and their mobility following large-scale disaster](#)

X Song, Q Zhang, Y Sekimoto, R Shibasaki

Proceedings of the 20th ACM SIGKDD international conference on Knowledge ...

[Modeling and probabilistic reasoning of population evacuation during large-scale disaster](#)

X Song, Q Zhang, Y Sekimoto, T Horanont, S Ueyama, R Shibasaki

Proceedings of the 19th ACM SIGKDD international conference on Knowledge ...

[Unsupervised skeleton extraction and motion capture from 3D deformable matching](#)

Q Zhang, X Song, X Shao, R Shibasaki, H Zhao

Neurocomputing 100, 170-182

[Interpreting cnn knowledge via an explanatory graph](#)

Q Zhang, R Cao, F Shi, YN Wu, SC Zhu

Thirty-Second AAAI Conference on Artificial Intelligence

---

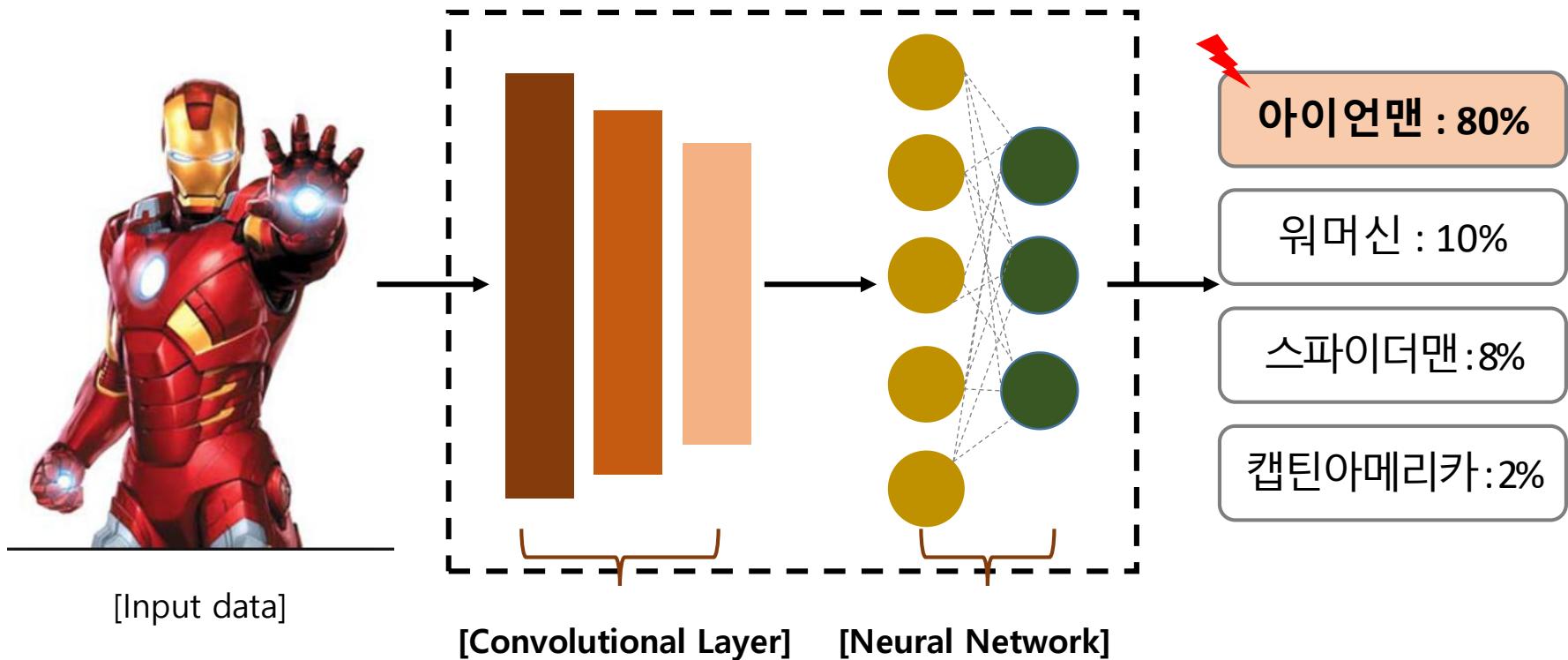
## 2. [CNN]Convolutional Neural Network

---

# [CNN]Convolutional Neural Network

- ❖ Convolutional Neural Network 기본 구조

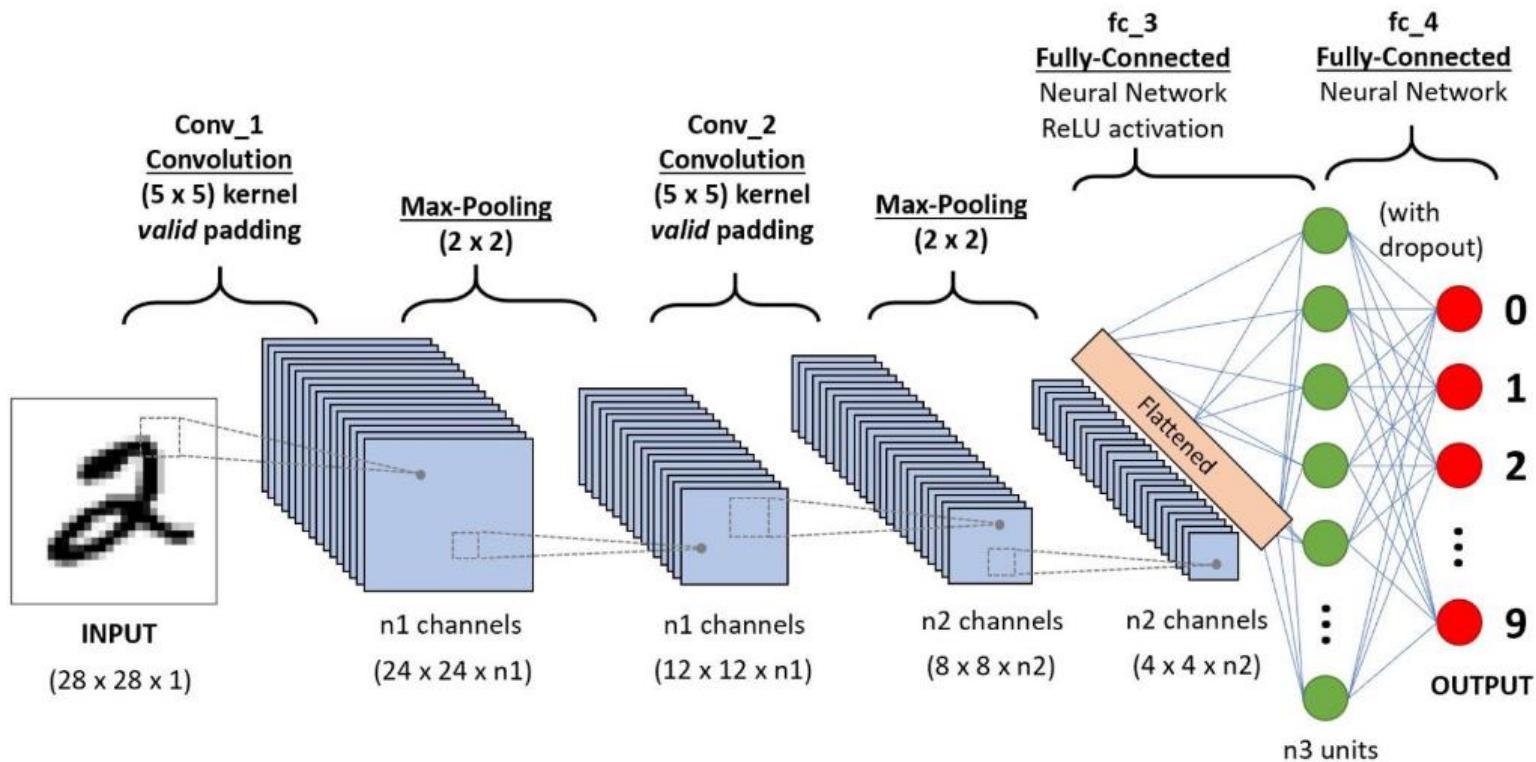
- Neural Network에 Convolution Layer를 사용한 방법론
- Object Detection, Classification 등 Visual Task에서 좋은 Performance 보임



# [CNN]Convolutional Neural Network

## ❖ Convolutional Neural Network 세부 구조

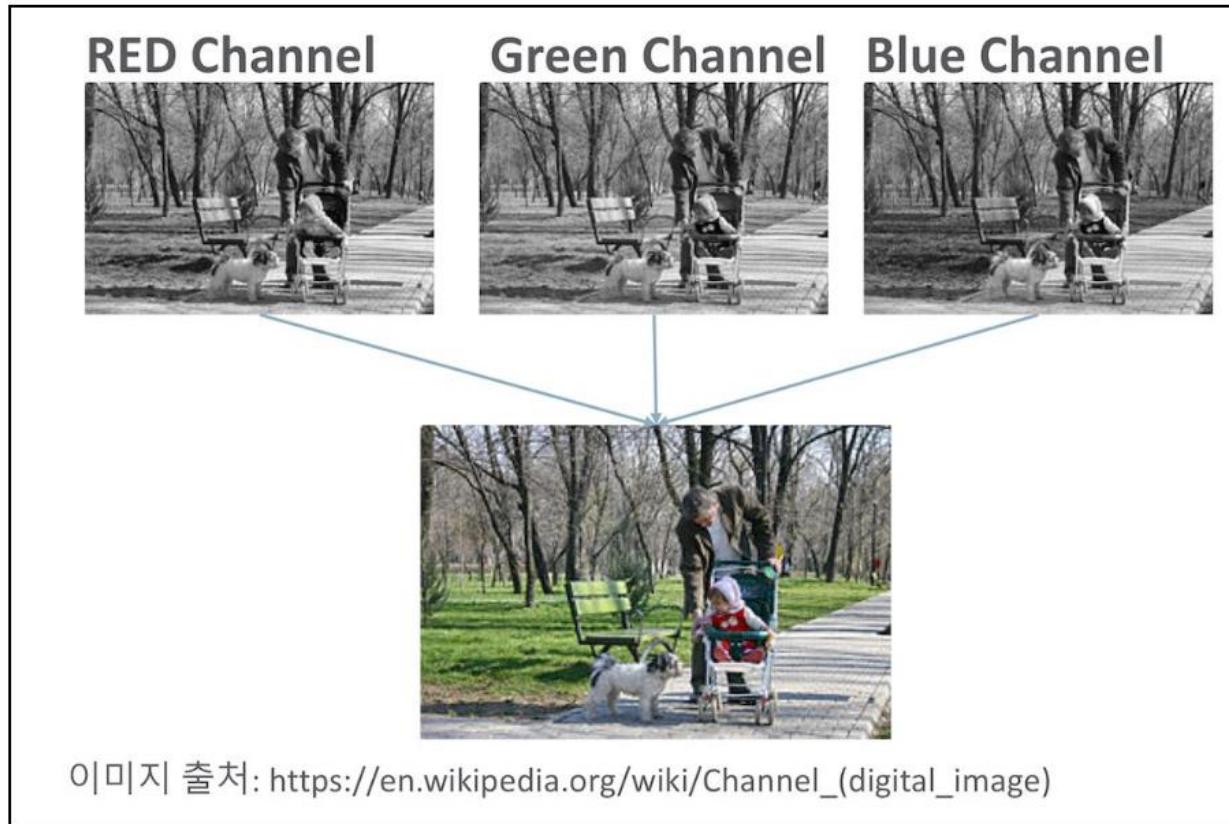
- Input data에 Convolution & Pooling을 활용하여 특징 추출
- 최종적으로 해당 Input이 어떤 Class에 속하는지 결정



# [CNN]Convolutional Neural Network

## ❖ Channel & Pooling

- Channel : Input이 이미지인 경우 그 이미지를 표현하는 R/G/B 각각을 의미함



# [CNN]Convolutional Neural Network

## ❖ Filter & Convolution

- Filter : Input의 특징을 찾아내기 위한 공용 parameter
- Convolution(합성곱) : Filter를 이동 시키면서 곱한 결과를 합산한 것

0	0	0	0	0	0	...
0	156	155	156	158	158	...
0	153	154	157	159	159	...
0	149	151	155	158	159	...
0	146	146	149	153	158	...
0	145	143	143	148	158	...
...	...	...	...	...	...	...

Input Channel #1 (Red)

0	0	0	0	0	0	...
0	167	166	167	169	169	...
0	164	165	168	170	170	...
0	160	162	166	169	170	...
0	156	156	159	163	168	...
0	155	153	153	158	168	...
...	...	...	...	...	...	...

Input Channel #2 (Green)

0	0	0	0	0	0	...
0	163	162	163	165	165	...
0	160	161	164	166	166	...
0	156	158	162	165	166	...
0	155	155	158	162	167	...
0	154	152	152	157	167	...
...	...	...	...	...	...	...

Input Channel #3 (Blue)

-1	-1	1
0	1	-1
0	1	1

Kernel Channel #1



308

+

1	0	0
1	-1	-1
1	0	-1

Kernel Channel #2



-498

0	1	1
0	1	0
1	-1	1

Kernel Channel #3

+

164

$$+ 1 = -25$$

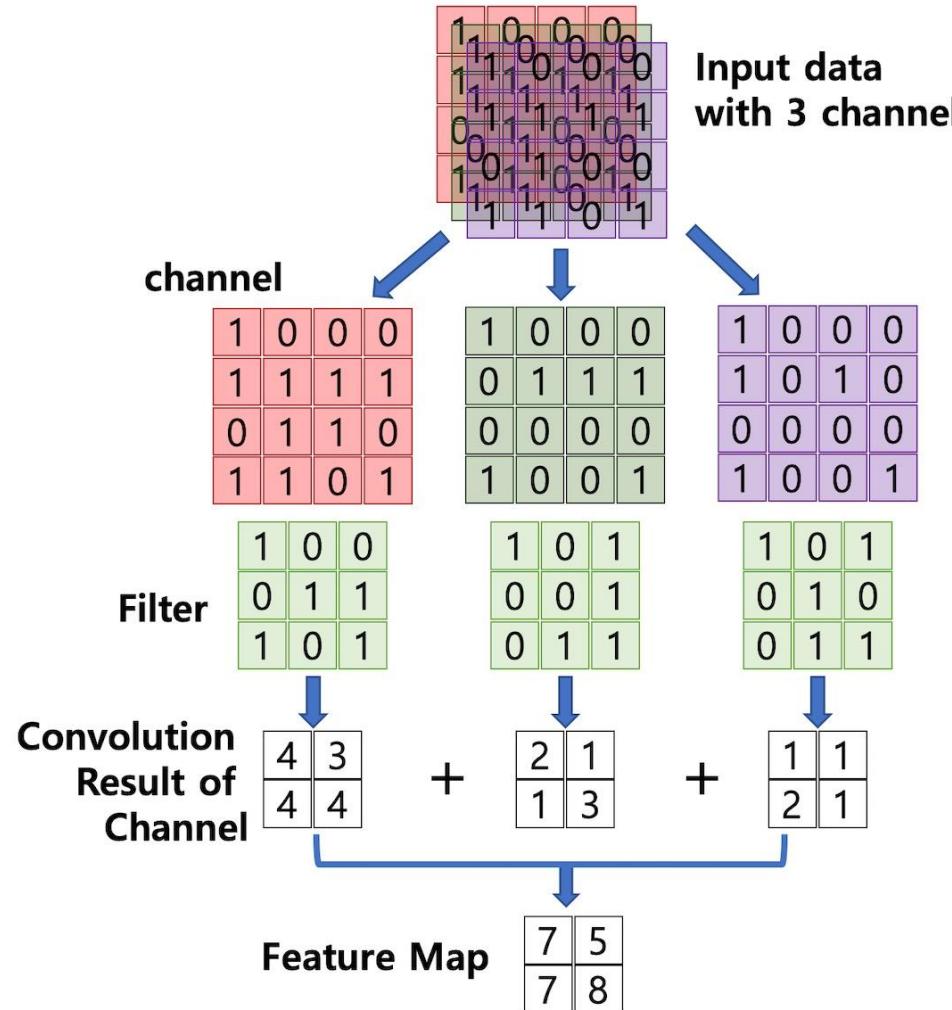


Bias = 1

-25					...
					...
					...
					...
...	...	...	...	...	...

# [CNN]Convolutional Neural Network

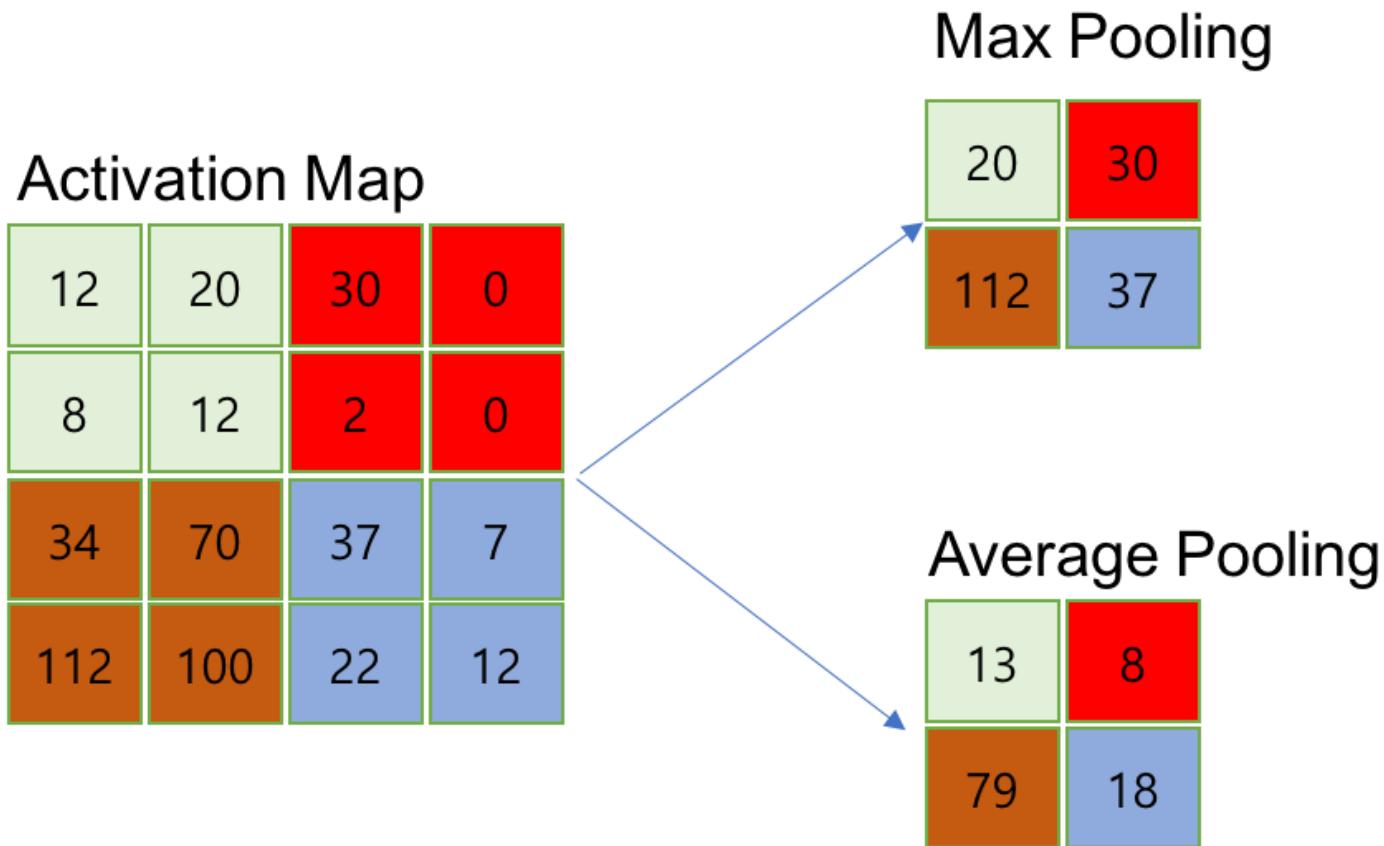
- ❖ Channel & Convolution & Filter 정리



# [CNN]Convolutional Neural Network

## ❖ Channel & Pooling

- Pooling : Convolution Layer의 출력 값을 입력으로 받아 크기를 줄이거나 특정 데이터를 강조하는 용도로 사용됨, Filter와 달리 Parameter가 아님



---

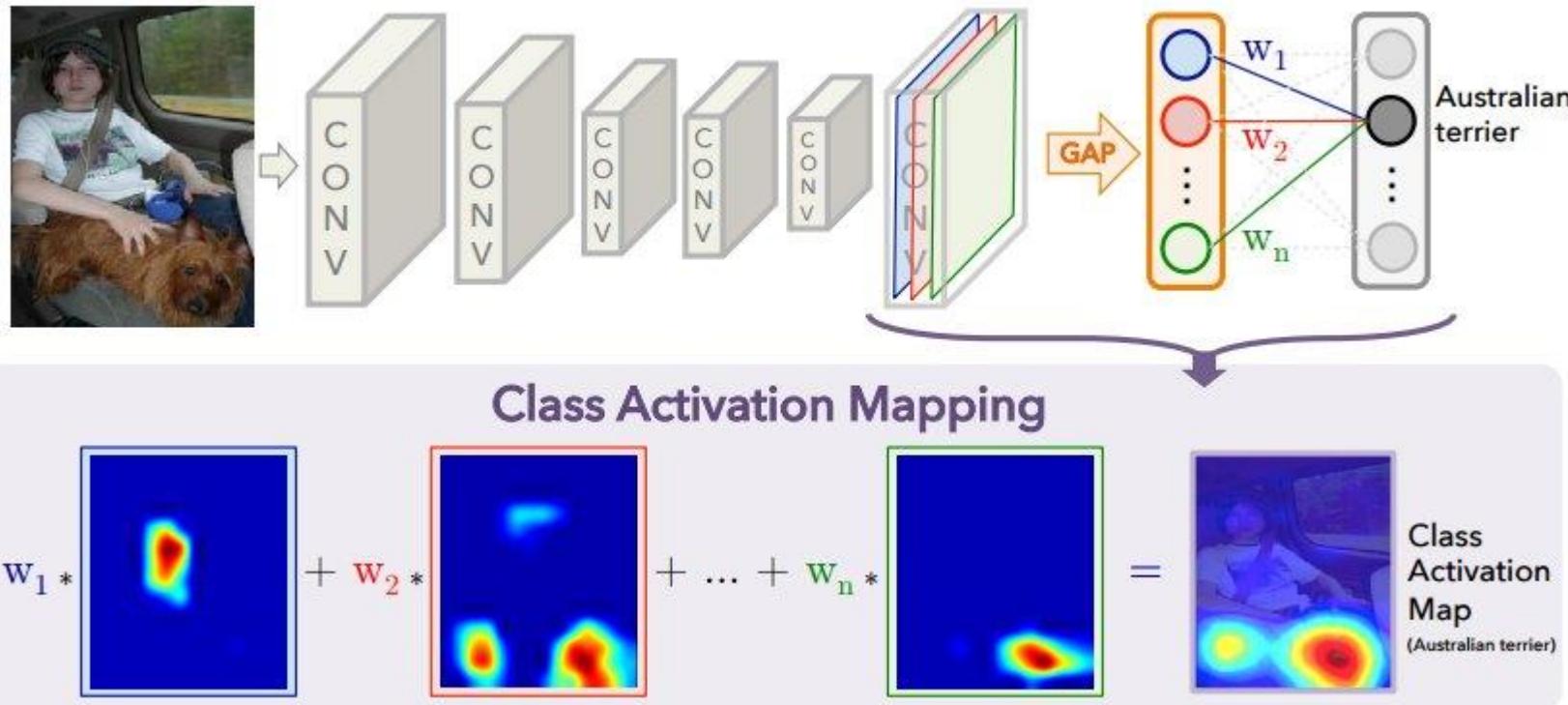
### 3. [CAM]Class Activation Map

---

# [CAM] Class Activation Map

## ❖ Class Activation Map 구조

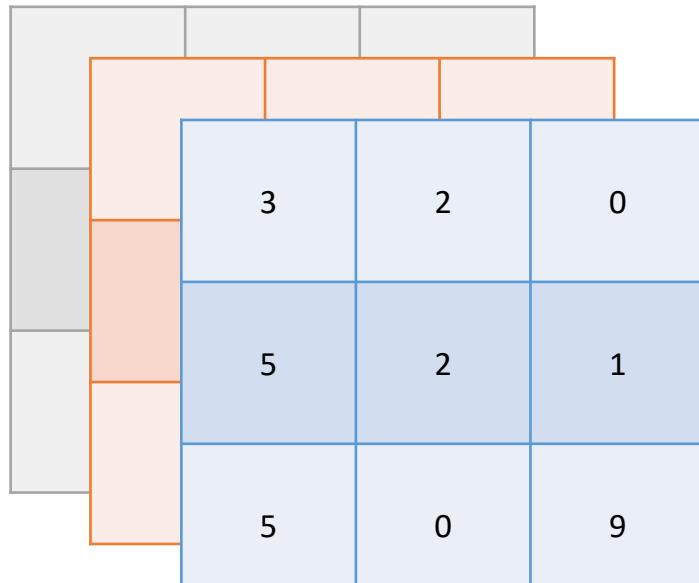
- CNN 구조에서 어떤 부분이 Class 결정에 큰 영향을 주었는지 확인 가능
- Flatten 대신 GAP을 적용, 마지막 판별 전까지 데이터 위치 정보 훼손되지 않음



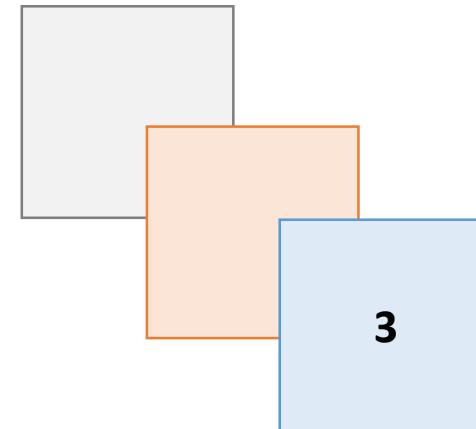
# [CAM]Class Activation Map

- ❖ GAP(Global Average Pooling)

- CNN 마지막에 Flatten 방법 대신 사용
- 결과는 Class개수와 동일, 채널별 평균값에 Softmax 함수 적용



[3x3x3(Channel)]

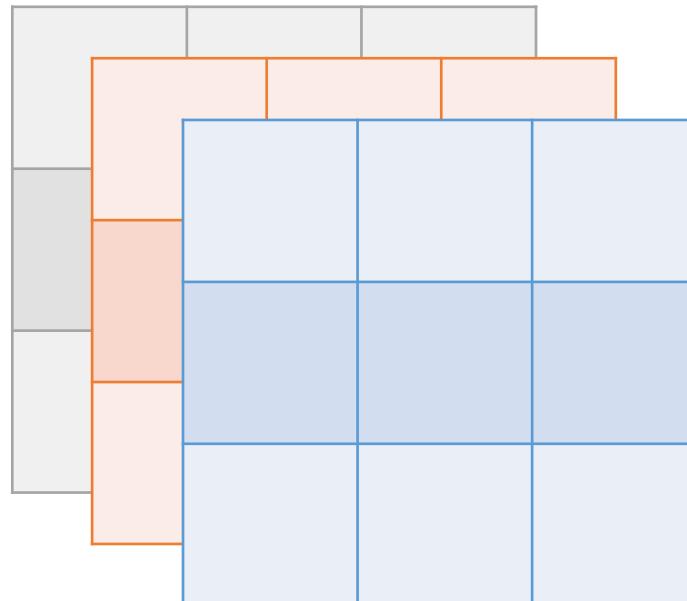


[1x1x3(channel)]

# [CAM] Class Activation Map

## ❖ Class Activation Map

- 해당 채널의 각 값에 Weight 값을 곱한 수치를 중요도로 해석
- 위처럼 계산된 모든 채널 값을 더하여  $N \times N$  행렬을 Heat-map으로 그림



1	0	9
3	2	7
4	3	2

$* w_1$



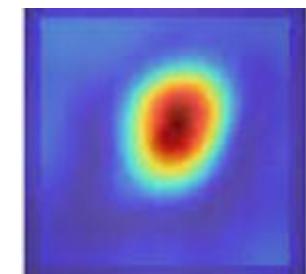
0	0	7
5	2	1
2	2	1

$* w_2$



3	3	0
5	2	1
5	0	9

$* w_3$



---

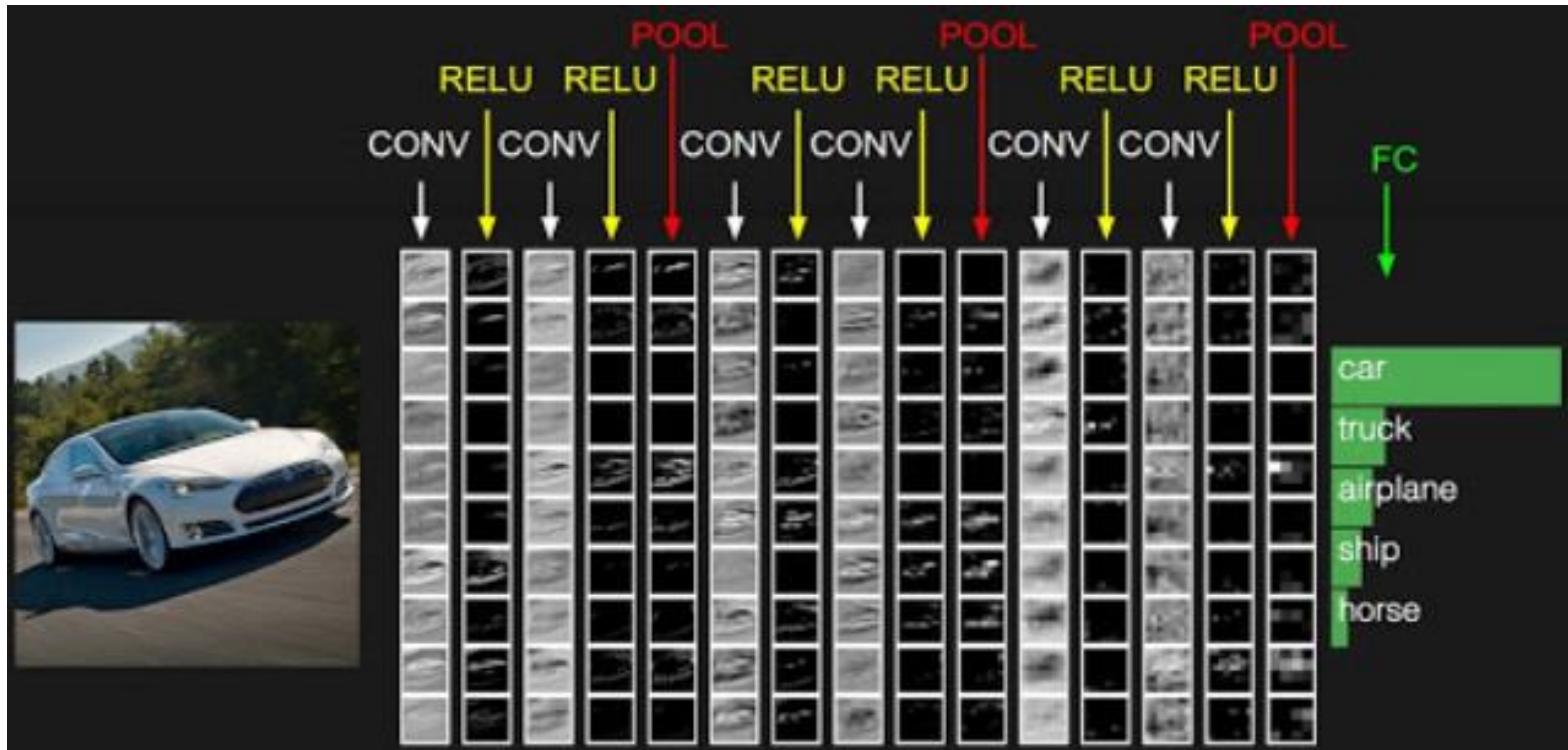
## 4. Interpretable Convolutional Neural Network

---

# Interpretable Convolutional Neural Network

## ❖ 핵심 컨셉

- CNN의 기본구조를 변형하지 않으면서 해석력을 높인 방법론
- Input의 정보를 최대로 압축해서 담고 있는 마지막 Conv-Layer만 사용



# Interpretable Convolutional Neural Network

- ❖ 일반적인 CNN filter와 Interpretable filter의 비교
  - 일반적인 filter를 사용한 경우에는 중요 feature를 정확히 찾지 못함
  - Interpretable Filter를 사용한 경우에는 얼굴 중심의 중요 feature를 찾음

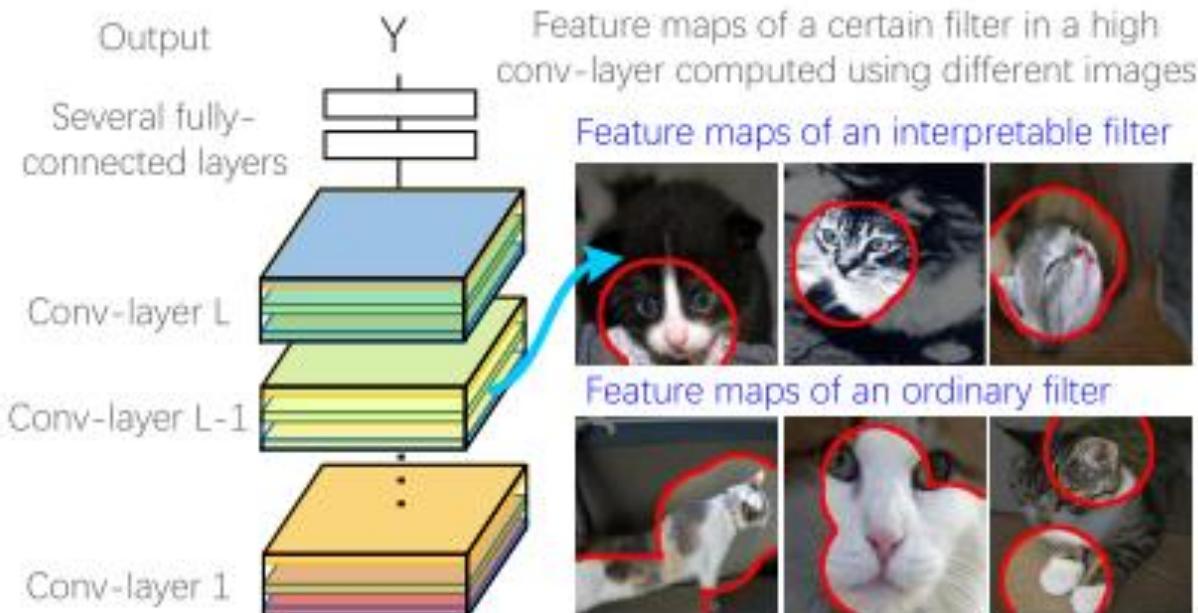
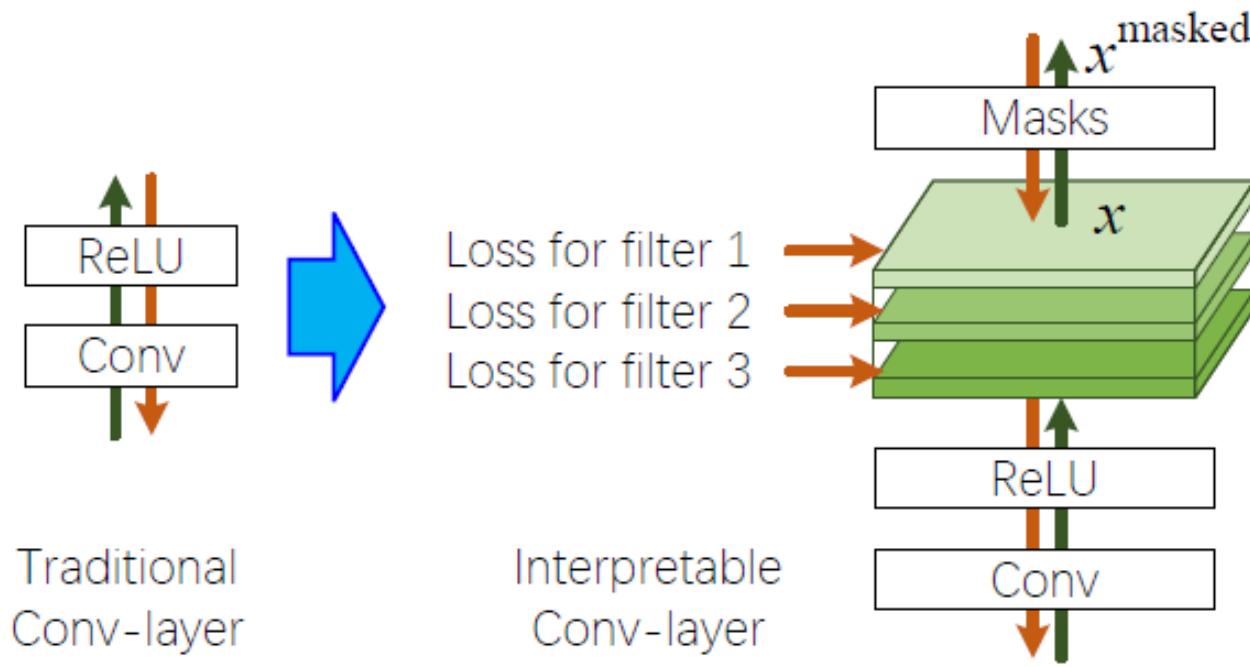


Figure 1. Comparison of a filter's feature maps in an interpretable CNN and those in a traditional CNN.

# Interpretable Convolutional Neural Network

- ❖ 전통적인 CNN 개선 방법
  - 마지막 Layer의 Filter별 각각 Loss를 더함
  - 각 Filter는 중요한 부분을 구별할 수 있는 Power를 지님



# Interpretable Convolutional Neural Network

## ❖ Filter 별로 계산하는 이유

- 반복된 Filter별 계산값 통해 더 객관적으로 중요 Feature를 선별할 수 있음
- Filter는 이미지의 특정 한 부분을 통해 활성화 함

2	1	0	1
3	7	19	11
0	5	14	5
0	1	2	1

[After ReLu Activation]

Step1.



1	0
0	5



37	96	55
28	77	46
5	15	19

Step2.



0	2
1	1



12	26	32
19	57	41
11	31	12

[Feature map]

# Interpretable Convolutional Neural Network

## ❖ Template 생성

- 트레이닝 데이터에서 같은 Class(ex 고양이)로 분류되는 데이터로 학습
- 공통적으로 가장 값이 큰 곳을  $\mu_i$ 로 지정한 Template을 구성

Step1.

1	0
0	5



37	96	55
28	77	46
5	15	19



$\mu_i$
[Template]

Step2.

0	2
1	1



12	26	32
19	57	41
11	31	12



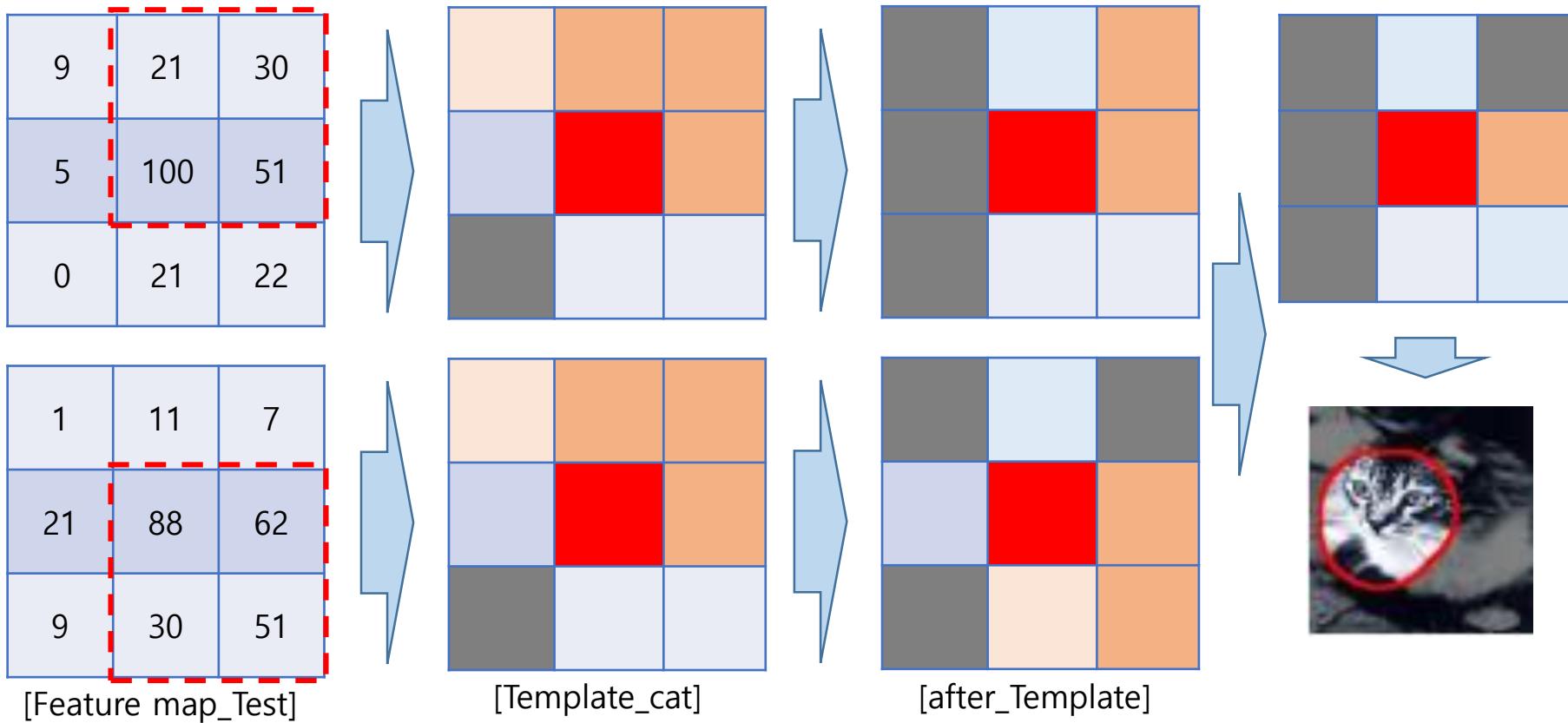
[Filter]

[Feature map\_Train]

# Interpretable Convolutional Neural Network

## ❖ Template 생성

- Test Data의 Top conv-layer의 값을 Template 값과 비교
- 학습을 진행하면서 비활성화 되는 부분을 중첩하면 중요 Feature 확인 가능



# Interpretable Convolutional Neural Network

---

- ❖ Loss Function
  - Filter별 Loss Function Term이 추가 됨

$$\begin{aligned}\text{Loss}_f &= -MI(\mathbf{X}; \mathbf{T}) \quad \text{for filter } f \\ &= - \sum_T p(T) \sum_x p(x|T) \log \frac{p(x|T)}{p(x)}\end{aligned}\quad (1)$$

$$\forall T \in \mathbf{T}, \quad p(x|T) = \frac{1}{Z_T} \exp [tr(x \cdot T)] \quad (2)$$

# Interpretable Convolutional Neural Network

## ❖ Loss Function

- Filter별 Loss Function Term이 추가 됨

$$\frac{\partial \text{Loss}}{\partial x_{ij}} = \lambda \frac{\partial \text{Loss}_f}{\partial x_{ij}} + \frac{1}{N} \sum_{i=k}^N \frac{\partial \mathbf{L}(\hat{y}_k, y_k^*)}{\partial x_{ij}} \quad (3)$$

$\lambda$  is a weight



$$\begin{aligned} \frac{\partial \text{Loss}_f}{\partial x_{ij}} &= \frac{1}{Z_T} \sum_T p(T) t_{ij} e^{tr(x \cdot T)} \left\{ tr(x \cdot T) - \log [Z_T p(x)] \right\} \\ &\approx \frac{p(\hat{T}) \hat{t}_{ij}}{Z_{\hat{T}}} e^{tr(x \cdot \hat{T})} \left\{ tr(x \cdot \hat{T}) - \log Z_{\hat{T}} - \log p(x) \right\} \quad (4) \end{aligned}$$

# Interpretable Convolutional Neural Network

---

- ❖ Loss Function
  - Filter별 Loss Function Term이 추가 됨

$$\begin{aligned}\text{Loss}_f = & -H(\mathbf{T}) + H(\mathbf{T}' = \{T^-, \mathbf{T}^+\} | \mathbf{X}) \\ & + \sum_x p(\mathbf{T}^+, x) H(\mathbf{T}^+ | X = x)\end{aligned}\quad (5)$$

$$H(\mathbf{T}' = \{T^-, \mathbf{T}^+\} | \mathbf{X}) = -\sum_x p(x) \sum_{T \in \{T^-, \mathbf{T}^+\}} p(T|x) \log p(T|x)\quad (6)$$

Equation 6: Low inter-category entropy

$$H(\mathbf{T}^+ | X = x) = \sum_{\mu} \tilde{p}(T_{\mu} | x) \log \tilde{p}(T_{\mu} | x)\quad (7)$$

Equation 7: Low spatial entropy

# Interpretable Convolutional Neural Network

- ❖ Masking Data
  - 기존 Data에 Mask(=Template)를 쓰운 이미지 결과

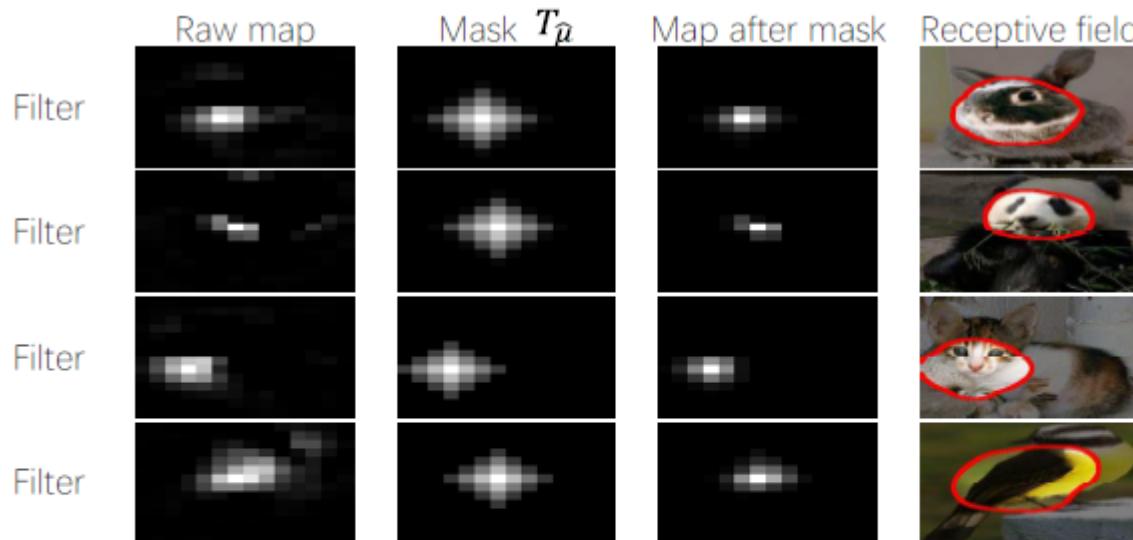
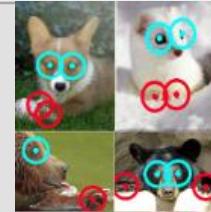


Figure 4. Given an input image  $I$ , from the left to the right, we consequently show the feature map of a filter after the ReLU layer  $x$ , the assigned mask  $T_{\hat{\mu}}$ , the masked feature map  $x^{\text{masked}}$ , and the image-resolution RF of activations in  $x^{\text{masked}}$  computed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

# Interpretable Convolutional Neural Network

## ❖ 실험 세팅(Data Set)

- ILSVRC 2013 DET Animal-Part dataset
- CUB200-2011 dataset
- Pascal VOC Part dataset

Data	내용	이미지
Animal-Part	Heads and legs of 30 animals	
CUB200-2011	Bird image of 200 pieces	
Pascal VOC Part	107 object landmarks in six animal categories	 bird cat

# Interpretable Convolutional Neural Network

- ❖ 실험 결과(Animal-Part dataset - instability)

Normalized distance :

$$d_1(p_k, \hat{\mu}) = \frac{\|p_k - p(\hat{\mu})\|}{\sqrt{w^2 + h^2}}$$

Relative location deviation :

$$D_{f,k} = \sqrt{\text{var}_I[d_I(p_k, \hat{\mu})]}$$

	gold.	bird	frog	turt.	liza.	koala	lobs.	dog	fox	cat	lion	tiger	bear	rabb.	hams.	squi.
AlexNet	0.161	0.167	0.152	0.153	0.175	0.128	0.123	0.144	0.143	0.148	0.137	0.142	0.144	0.148	0.128	0.149
AlexNet, interpretable	<b>0.084</b>	<b>0.095</b>	<b>0.090</b>	<b>0.107</b>	<b>0.097</b>	<b>0.079</b>	<b>0.077</b>	<b>0.093</b>	<b>0.087</b>	<b>0.095</b>	<b>0.084</b>	<b>0.090</b>	<b>0.095</b>	<b>0.095</b>	<b>0.077</b>	<b>0.095</b>
VGG-16	0.153	0.156	0.144	0.150	0.170	0.127	0.126	0.143	0.137	0.148	0.139	0.144	0.143	0.146	0.125	0.150
VGG-16, interpretable	<b>0.076</b>	<b>0.099</b>	<b>0.086</b>	<b>0.115</b>	<b>0.113</b>	<b>0.070</b>	<b>0.084</b>	<b>0.077</b>	<b>0.069</b>	<b>0.086</b>	<b>0.067</b>	<b>0.097</b>	<b>0.081</b>	<b>0.079</b>	<b>0.066</b>	<b>0.065</b>
VGG-M	0.161	0.166	0.151	0.153	0.176	0.128	0.125	0.145	0.145	0.150	0.140	0.145	0.144	0.150	0.128	0.150
VGG-M, interpretable	<b>0.088</b>	<b>0.088</b>	<b>0.089</b>	<b>0.108</b>	<b>0.099</b>	<b>0.080</b>	<b>0.074</b>	<b>0.090</b>	<b>0.082</b>	<b>0.103</b>	<b>0.079</b>	<b>0.089</b>	<b>0.101</b>	<b>0.097</b>	<b>0.082</b>	<b>0.095</b>
VGG-S	0.158	0.166	0.149	0.151	0.173	0.127	0.124	0.143	0.142	0.148	0.138	0.142	0.143	0.148	0.128	0.146
VGG-S, interpretable	<b>0.087</b>	<b>0.101</b>	<b>0.093</b>	<b>0.107</b>	<b>0.096</b>	<b>0.084</b>	<b>0.078</b>	<b>0.091</b>	<b>0.082</b>	<b>0.101</b>	<b>0.082</b>	<b>0.089</b>	<b>0.097</b>	<b>0.091</b>	<b>0.076</b>	<b>0.098</b>
AlexNet	horse	zebra	swine	hippo	catt.	sheep	ante.	camel	otter	arma.	monk.	elep.	red pa.	gia.pa.	<b>Avg.</b>	
AlexNet, interpretable	0.152	0.154	0.141	0.141	0.144	0.155	0.147	0.153	0.159	0.160	0.139	0.125	0.140	0.125	0.146	
VGG-16	<b>0.098</b>	<b>0.084</b>	<b>0.091</b>	<b>0.089</b>	<b>0.097</b>	<b>0.101</b>	<b>0.085</b>	<b>0.102</b>	<b>0.104</b>	<b>0.095</b>	<b>0.090</b>	<b>0.085</b>	<b>0.084</b>	<b>0.073</b>	<b>0.091</b>	
VGG-16, interpretable	0.150	0.153	0.141	0.140	0.140	0.150	0.144	0.149	0.154	0.163	0.136	0.129	0.143	0.125	0.144	
VGG-M	<b>0.106</b>	<b>0.077</b>	<b>0.094</b>	<b>0.083</b>	<b>0.102</b>	<b>0.097</b>	<b>0.091</b>	<b>0.105</b>	<b>0.093</b>	<b>0.100</b>	<b>0.074</b>	<b>0.084</b>	<b>0.067</b>	<b>0.063</b>	<b>0.085</b>	
VGG-M, interpretable	0.151	0.158	0.140	0.140	0.143	0.155	0.146	0.154	0.160	0.161	0.140	0.126	0.142	0.127	0.147	
VGG-S	<b>0.095</b>	<b>0.080</b>	<b>0.095</b>	<b>0.084</b>	<b>0.092</b>	<b>0.094</b>	<b>0.077</b>	<b>0.104</b>	<b>0.102</b>	<b>0.093</b>	<b>0.086</b>	<b>0.087</b>	<b>0.089</b>	<b>0.068</b>	<b>0.090</b>	
VGG-S, interpretable	0.149	0.155	0.139	0.140	0.141	0.155	0.143	0.154	0.158	0.157	0.140	0.125	0.139	0.125	0.145	

Table 3. Location instability of filters ( $E_{f,k}[D_{f,k}]$ ) in CNNs that are trained for single-category classification using the ILSVRC 2013 DET Animal-Part dataset [36]. Filters in our interpretable CNNs exhibited significantly lower localization instability than ordinary CNNs.

# Interpretable Convolutional Neural Network

- ❖ 실험 결과(Pascal VOC Part Dataset - instability)

Normalized distance :

$$d_1(p_k, \hat{\mu}) = \frac{\|p_k - p(\hat{\mu})\|}{\sqrt{w^2 + h^2}}$$

Relative location deviation :

$$D_{f,k} = \sqrt{\text{var}_{I[d_I(p_k, \hat{\mu})]}}$$

	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 ( $E_{f,k}[D_{f,k}]$ ) 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.

# Interpretable Convolutional Neural Network

- ❖ 실험 결과(CUB200-2011 dataset- instability)

Normalized distance :

$$d_1(p_k, \hat{\mu}) = \frac{\|p_k - p(\hat{\mu})\|}{\sqrt{w^2 + h^2}}$$

Relative location deviation :

$$D_{f,k} = \sqrt{\text{var}_{I[d_I(p_k, \hat{\mu})]}}$$

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 ( $E_{f,k}[D_{f,k}]$ ) in CNNs for single-category classification using the CUB200-2011 dataset.

# Interpretable Convolutional Neural Network

- ❖ 실험 결과(Multi Category - instability)

Normalized distance :

$$d_1(p_k, \hat{\mu}) = \frac{\|p_k - p(\hat{\mu})\|}{\sqrt{w^2 + h^2}}$$

Relative location deviation :

$$D_{f,k} = \sqrt{\text{var}_{I[d_I(p_k, \hat{\mu})]}}$$

Dataset Network	ILSVRC Part [36]		Pascal VOC Part [3]	
	Logistic log loss <sup>4</sup>		Logistic log loss <sup>4</sup>	Softmax log loss
VGG-16 interpretable	–	–	0.128	0.142
VGG-M interpretable	0.167	<b>0.096</b>	0.135	0.137
VGG-S interpretable	0.131	<b>0.083</b>	0.138	0.138

Table 6. Location instability of filters ( $E_{f,k}[D_{f,k}]$ ) in CNNs that are trained for multi-category classification. Filters in our interpretable CNNs exhibited significantly lower localization instability than ordinary CNNs in all comparisons.

# Interpretable Convolutional Neural Network

## ❖ 실험 결과(Multi Category - instability)

Normalized distance :

$$d_1(p_k, \hat{\mu}) = \frac{\|p_k - p(\hat{\mu})\|}{\sqrt{w^2 + h^2}}$$

Relative location deviation :

$$D_{f,k} = \sqrt{\text{var}_I[d_I(p_k, \hat{\mu})]}$$

	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 ( $E_{f,k}[D_{f,k}]$ ) 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.

# Interpretable Convolutional Neural Network

## ❖ 실험 결과(Single & Multi Category)

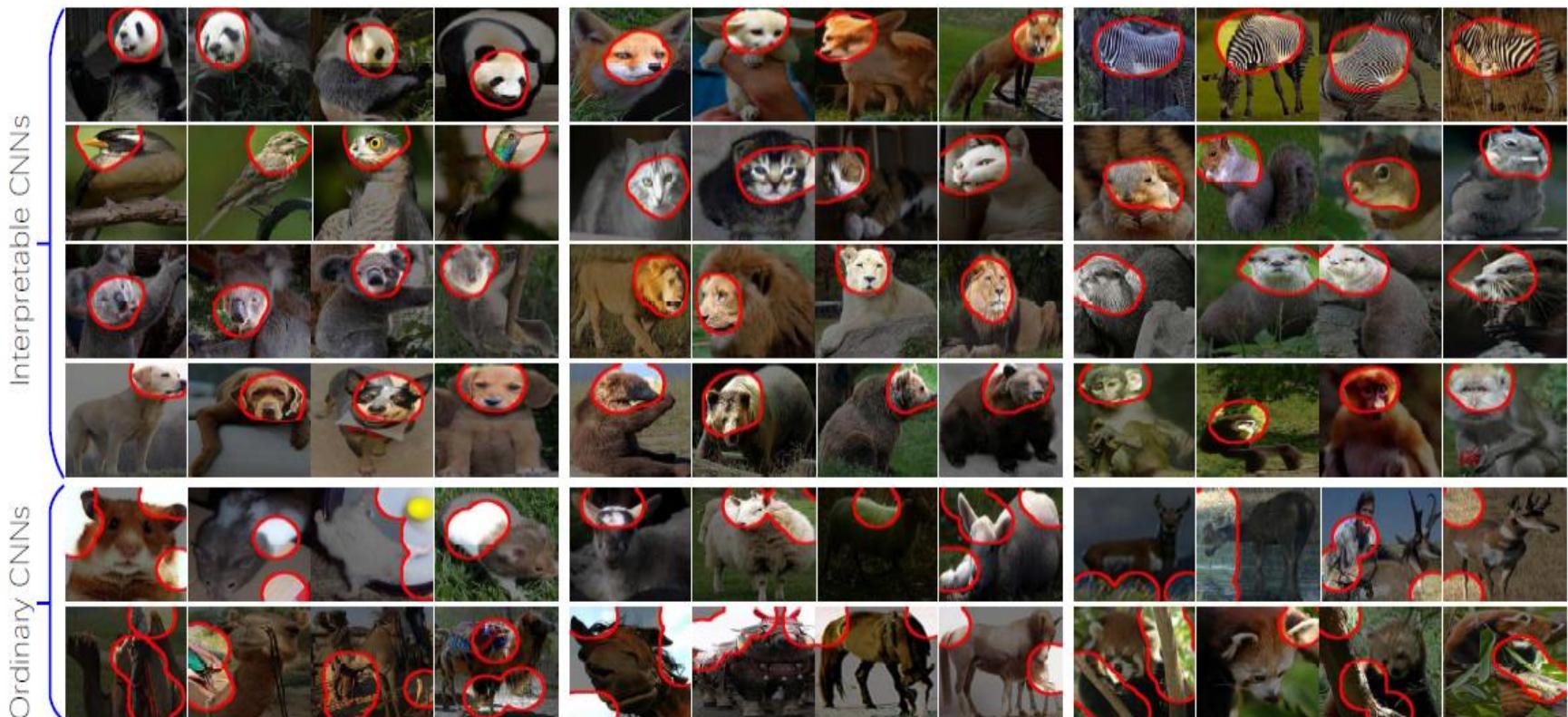
- Classification Accuracy = Single category is bad, Multi-category is better
- 명확한 이유에 대해서는 저자도 언급이 없었음

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

# Interpretable Convolutional Neural Network

- ❖ 실험 결과(Image Data)
  - 일반적인 CNN에 비해 이미지의 특징을 더 잘 추출하고 있음



# Interpretable Convolutional Neural Network

## ❖ 실험 결과(Image Data)

- Top conv-layer만 사용 했지만 중요 feature를 잘 선별하고 있었음

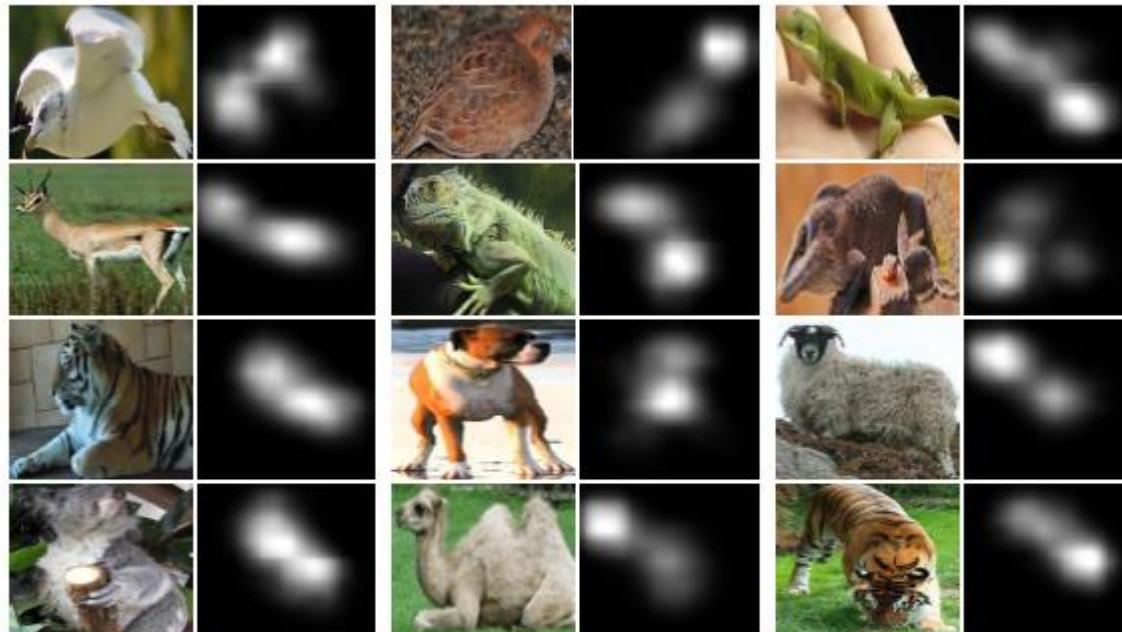


Figure 6. Heat maps for distributions of object parts that are encoded in interpretable filters. We use all filters in the top conv-layer to compute the heat map.

---

# 5. Conclusion

---

# Conclusion

---

## ❖ Conclusion

- ① 기본 CNN 구조를 크게 변경한 것이 아니기 때문에 다른 구조에도 적용 가능
- ② 특별한 주석도 필요 없다. Filter 값들이 그것을 표현
- ③ Loss function의 변경이나 Training Set의 대한 변경도 필요 없음
- ④ 하지만 Single Category 사례처럼 모델의 성능을 조금 떨어뜨릴 수 있음

## ❖ Critic

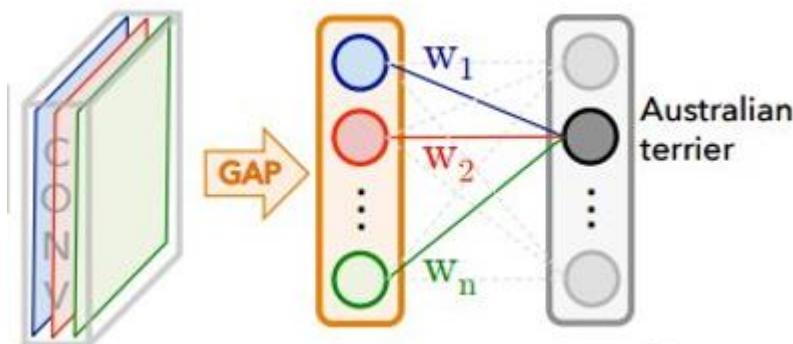
- ① 사실 기본적인 구조를 변경하지 않았다 것을 제외하고는 큰 의미 찾지 못함
- ② CAM과 대비했을 때 효율적인지 의문스러움
- ③ Classification의 성능이 좋아지고, 나빠진 근거가 불충분
- ④ 개인적으로는 Cite 횟수에 대해 Cross-Check 해야겠다고 생각함

**감사합니다.**

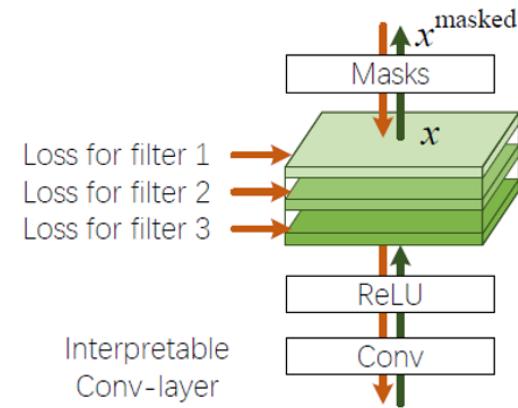
# Appendix.

## ❖ CAM vs Interpretable CNN

- CAM은 마지막 Gap Layer에 Weight \* Channel을 합친 값
- Interpretable CNN은 기본 CNN의 구조 변경 없이 Filter에 Loss 만 계산



[CAM Last Layer]



[Interpretable CNN Last Layer]

## ❖ 출처

- Figure: Basic architecture of Convolutional Neural Networks(<https://pydeeplearning.weebly.com/blog/basics-of-convolutional-neural-networks>)
- <https://bcho.tistory.com/1149>(조대협의 블로그)
- <https://blog.ees.guru/50> (EeS의 연구실)
- Quanshi Zhang, Ying Nian Wu, and Song-Chun Zhu(2018). Interpretable Convolutional Neural Networks. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (pp. 8827-8836)
- <http://taewan.kim/post/cnn/>
- <https://medium.com/@kamwohng/interpretable-convolutional-neural-network-3f7ef6c9b7ae>

