

# Learning Deep Features for Discriminative Localization

Sanghyun Seo

2018-04-13

# Contents

1. Google Colaboratory
2. Paper Information
3. Introduction & Related Work
4. Class Activation Map
5. Weakly-supervised Object Localization
6. Conclusion

# Google Colaboratory

- Collaboration + jupyter notebook

안녕하세요, Colaboratory입니다.

파일 수정 보기 삽입 런타임 도구 도움말

코드 텍스트 업 셀 드라이브로 복사 연결 수정 중

Colaboratory에 오신 것을 환영합니다!

Colaboratory는 텍스트, 코드, 코드 출력을 하나의 공동작업 문서로 통합해 주는 데이터 분석 도구입니다.

```
[ ] print('Hello, Colaboratory!')
```

Hello, Colaboratory!

Colaboratory를 사용하면 클릭 한 번만으로 사용 중인 브라우저에서 텐서플로우 코드를 실행할 수 있습니다. 아래의 예에서는 두 개의 행렬을 추가합니다.

$$\begin{bmatrix} 1. & 1. & 1. \end{bmatrix} + \begin{bmatrix} 1. & 2. & 3. \end{bmatrix} = \begin{bmatrix} 2. & 3. & 4. \end{bmatrix}$$
$$\begin{bmatrix} 1. & 1. & 1. \end{bmatrix} + \begin{bmatrix} 4. & 5. & 6. \end{bmatrix} = \begin{bmatrix} 5. & 6. & 7. \end{bmatrix}$$

```
[ ] import tensorflow as tf
import numpy as np

with tf.Session():
    input1 = tf.constant(1.0, shape=[2, 3])
    input2 = tf.constant(np.reshape(np.arange(1.0, 7.0, dtype=np.float32), (2, 3)))
    output = tf.add(input1, input2)
    result = output.eval()

result
```

array([[ 2., 3., 4.],
 [ 5., 6., 7.]], dtype=float32)

Colaboratory에는 [matplotlib](#)와 같이 널리 사용되는 라이브러리가 포함되어 시각화를 단순화할 수 있습니다.

```
[ ] import matplotlib.pyplot as plt
x = np.arange(20)
y = [x_i + np.random.randn() for x_i in x]
a, b = np.polyfit(x, y, 1)
plt.plot(x, y, 'o', np.arange(20), a+np.arange(20)*b, '-');
```

Colaboratory는 Google Cloud BigQuery와 함께 사용할 수 있습니다.

[샘플 BigQuery 노트](#).

# Google Colaboratory

- Tesla K80(12.6gb GPU) ≥ GTX 1080TI(11gb GPU)

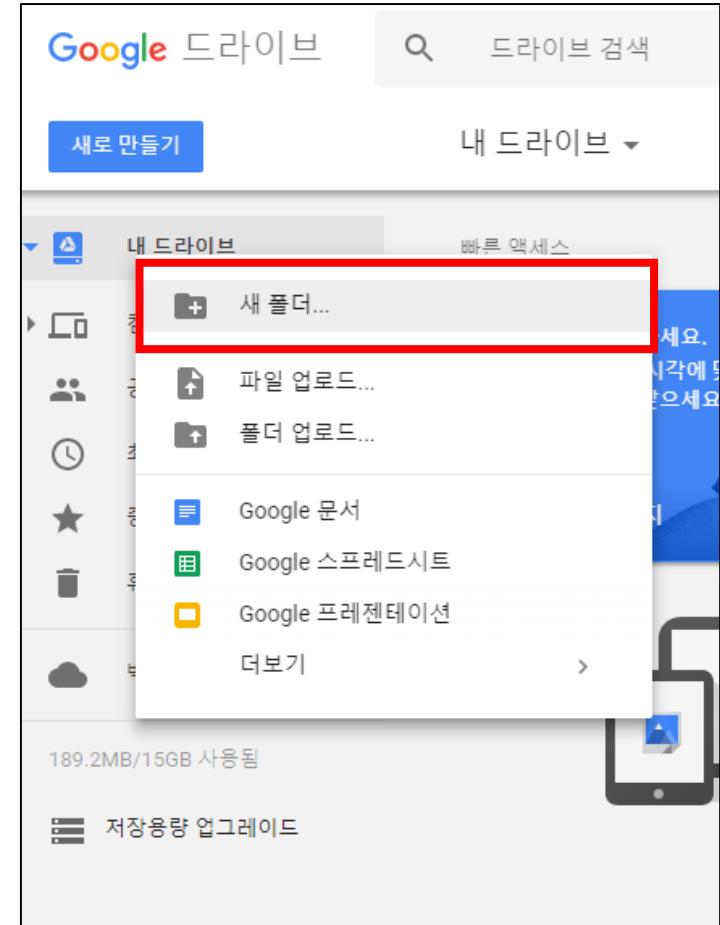
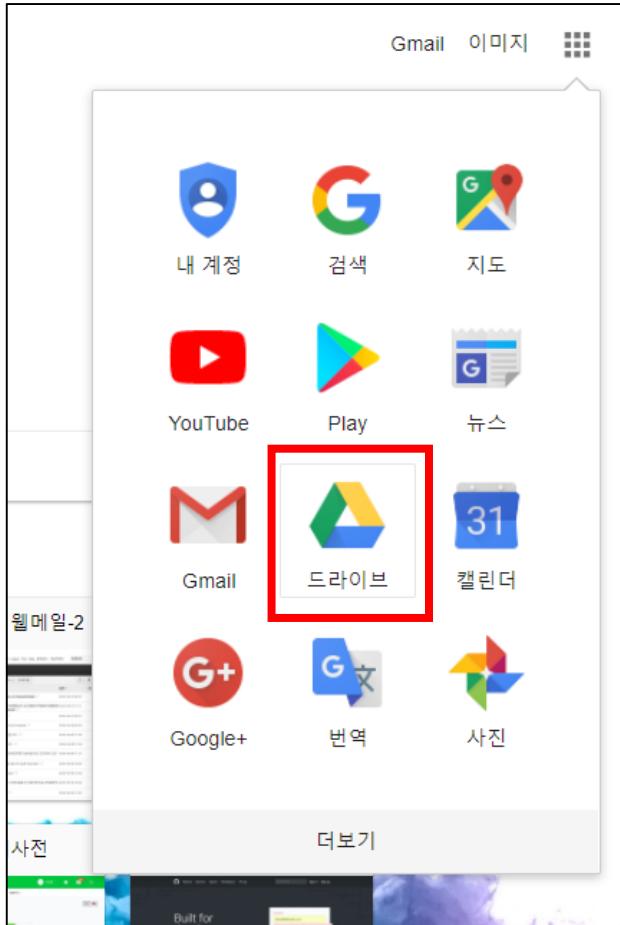
The screenshot shows a Google Colaboratory notebook titled "Untitled0.ipynb". The code cell contains Python code to check GPU memory usage. The output cell shows the results of running the code, indicating a Tesla K80 GPU with 12.6 GB of RAM.

```
# memory_footprint support libraries/code
!ln -sf /opt/bin/nvidia-smi /usr/bin/nvidia-smi
!pip install gputil
!pip install psutil
!pip install humanize
import psutil
import humanize
import os
import GPUtil as GPU
GPUs = GPU.getGPUs()
# XXX: only one GPU on Colab and isn't guaranteed
gpu = GPUs[0]
print(gpu.__dict__["name"])
def printm():
    process = psutil.Process(os.getpid())
    print("Gen RAM Free: " + humanize.naturalsize(psutil.virtual_memory().available), " | Proc size: " + humanize.naturalsize(process.memory_info().rss))
    print("GPU RAM Free: {0:.0f}MB | Used: {1:.0f}MB | Util {2:.0f}% | Total {3:.0f}MB".format(gpu.memoryFree, gpu.memoryUsed, gpu.memoryUtil, gpu.totalMemory))
    printm()

Requirement already satisfied: gputil in /usr/local/lib/python3.6/dist-packages
Requirement already satisfied: numpy in /usr/local/lib/python3.6/dist-packages (from gputil)
Requirement already satisfied: psutil in /usr/local/lib/python3.6/dist-packages
Requirement already satisfied: humanize in /usr/local/lib/python3.6/dist-packages
Tesla K80
Gen RAM Free: 12.6 GB | Proc size: 132.1 MB
GPU RAM Free: 11439MB | Used: 0MB | Util 0% | Total 11439MB
```

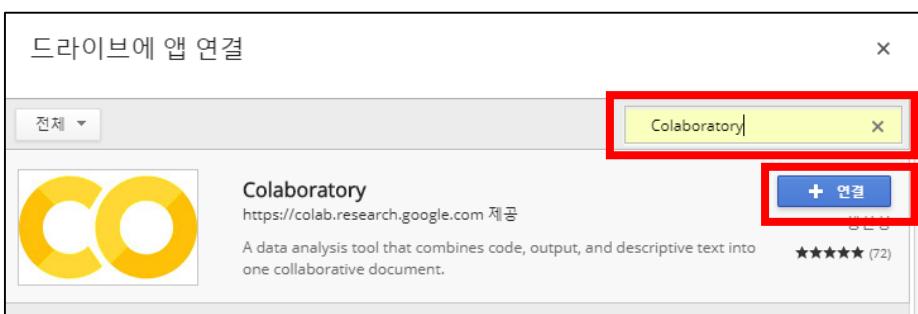
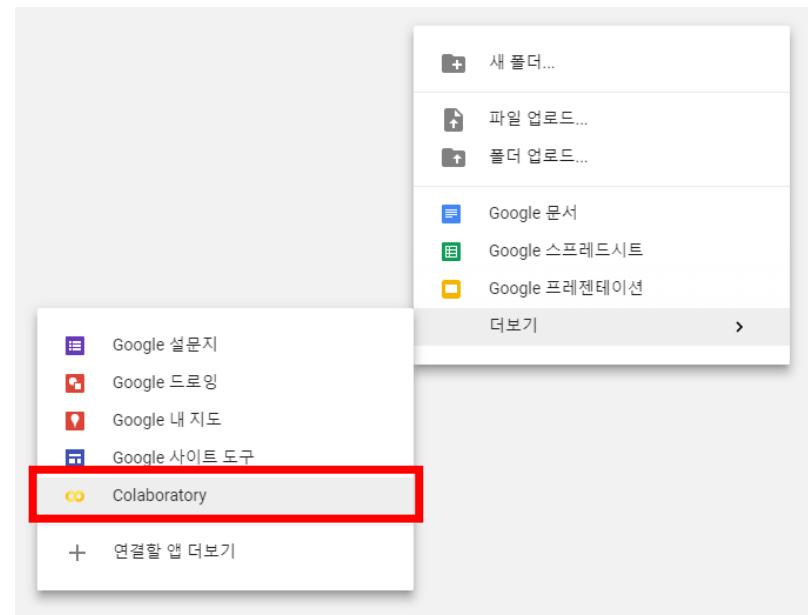
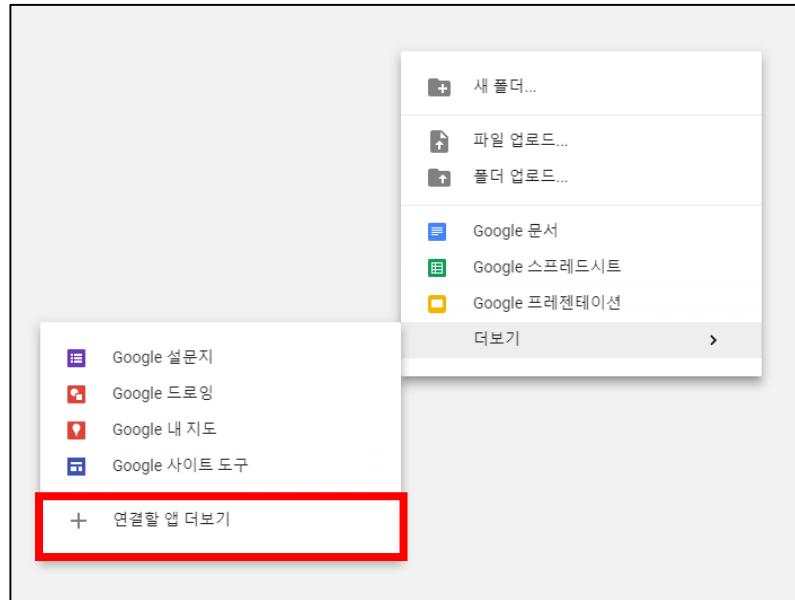
# Google Colaboratory

- Login google → go to the google drive → my drive → new folder(Optional)



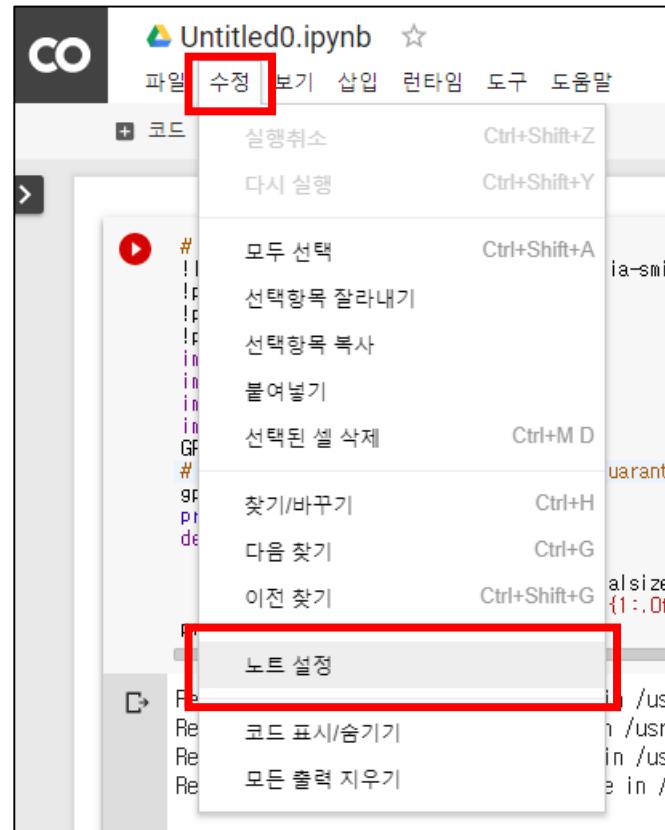
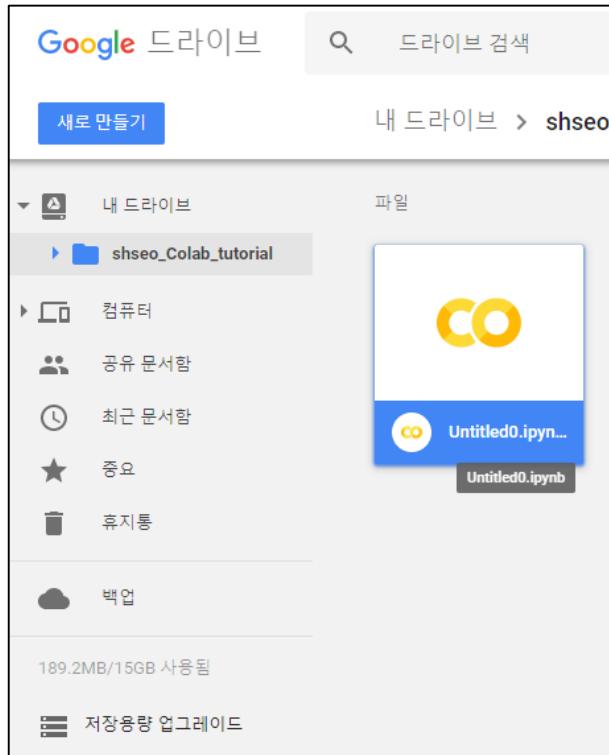
# Google Colaboratory

- Mouse right click → See more → See more app to connect → search the Colaboratory → Connect
- Mouse right click → See more → **Colaboratory**



# Google Colaboratory

- You can make .ipynb file
- Edit → Setting note → GPU → Save



# Google Colaboratory

- Simple test
- Import tensorflow

The screenshot shows the Google Colaboratory interface. At the top, there's a toolbar with a 'CO' logo, a file menu (파일), and various tool icons. Below the toolbar is a code editor area. A code cell contains the following Python script:

```
import tensorflow as tf
print("tensorflow version: %s" % tf.__version__)

x = tf.constant(1)
y = tf.constant(2)
z = x + y

sess = tf.Session()
res = sess.run(z)

print("x + y = %d" % res)
```

The output cell below the code shows the results of the execution:

```
tensorflow version: 1.6.0
x + y = 3
```

The output text is highlighted with a red box.

# Google Colaboratory

- Import torch, but there are no module named 'torch'
- Install TORCH → click insert → run new cell

목차 코드 스니펫 X

import torch

Install [pytorch](http://pytorch.org/) ➔

Install [pytorch](http://pytorch.org/) 삽입

```
# http://pytorch.org/
from os import path
from wheel.pep425tags import get_abbr_impl, get_implementation_version_tag
platform = '{0}-{1}'.format(get_abbr_impl(), get_implementation_version_tag())
accelerator = 'cu80' if path.exists('/opt/nvidia/cuda/include') else 'cpu'
!pip install -q http://download.pytorch.org/whl/{accelerator}/torch-0.3.0.post4-{platform}-py3-none-any.whl
import torch
```

[7] import torch

ModuleNotFoundError Traceback (most recent call last)  
<ipython-input-7-eb42ca6e4af3> in <module>()  
----> 1 import torch

ModuleNotFoundError: No module named 'torch'

NOTE: If your import is failing due to a missing package, you can manually install dependencies using either !pip or !apt.

To install torch, click the button below.

INSTALL TORCH SEARCH STACK OVERFLOW + 양식

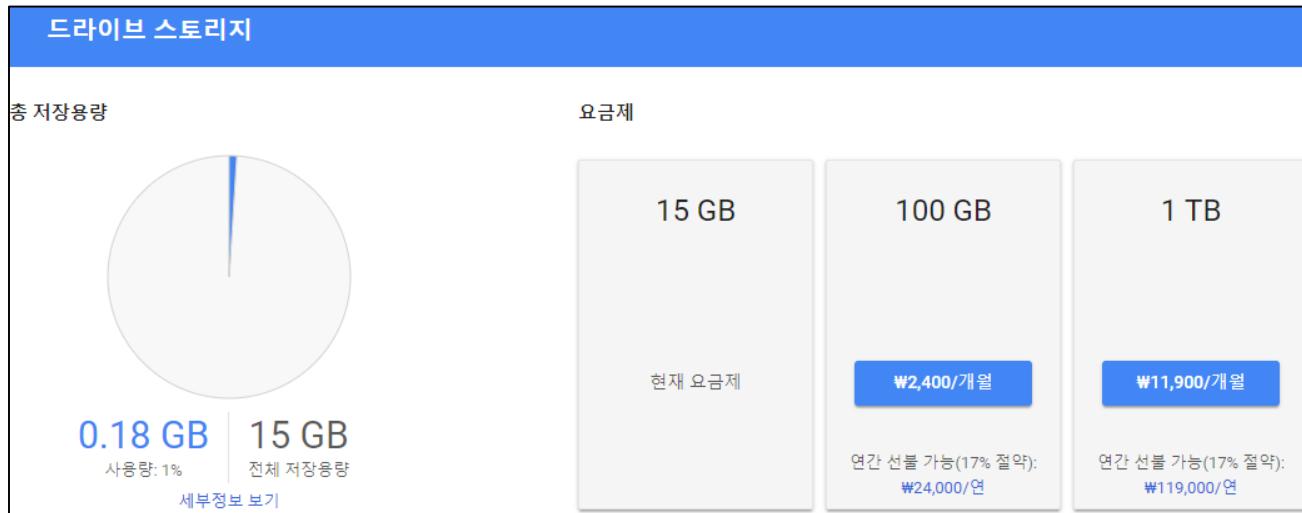
# http://pytorch.org/
from os import path
from wheel.pep425tags import get\_abbr\_impl, get\_implementation\_version\_tag
platform = '{0}-{1}'.format(get\_abbr\_impl(), get\_implementation\_version\_tag())
accelerator = 'cu80' if path.exists('/opt/nvidia/cuda/include') else 'cpu'
!pip install -q http://download.pytorch.org/whl/{accelerator}/torch-0.3.0.post4-{platform}-py3-none-any.whl
import torch

▶ import torch
print("torch version: %s" % torch.\_\_version\_\_)

torch version: 0.3.0.post4

# Google Colaboratory

- It appears that the instance will last 12 hours
- It is possible to upload and download files using the google drive
- For example, for every learning epoch, save the .ckpt file on the google drive, and continue to use the saved .ckpt file after 12 hours
- Default is 15gb cloud drive
- By investing about 2,400 won or 11,900 won a month, we can get Tesla- K80(12gb GPU).



# Paper Information

- **Title:** Learning Deep Features for Discriminative Localization
- **Author:** Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, Antonio Torralba
- **Affiliation:** Computer Science and Artificial Intelligence Laboratory, MIT
- **Proceeding:** (CVPR 2016) Computer Vision and Pattern Recognition 2016
- <http://cnnlocalization.csail.mit.edu/>
- <https://github.com/metalbubble/CAM>

**Learning Deep Features for Discriminative Localization**

Bolei Zhou   Aditya Khosla   Agata Lapedriza   Aude Oliva   Antonio Torralba  
Massachusetts Institute of Technology

In this work, we revisit the global average pooling layer and shed light on how it explicitly enables the convolutional neural network to have remarkable localization ability despite being trained on image-level labels. While this technique was previously proposed as a means for regularizing training, we find that it actually builds a generic localizable deep representation that can be applied to a variety of tasks. Despite the apparent simplicity of global average pooling, we are able to achieve 37.1% top-5 error for object localization on ILSVRC 2014, which is remarkably close to the 34.2% top-5 error achieved by a fully supervised CNN approach. We demonstrate that our network is able to localize the discriminative image regions on a variety of tasks despite not being trained for them.

[Source code and pre-trained models are available.](#)



Download CVPR'16 Paper Supplementary Materials

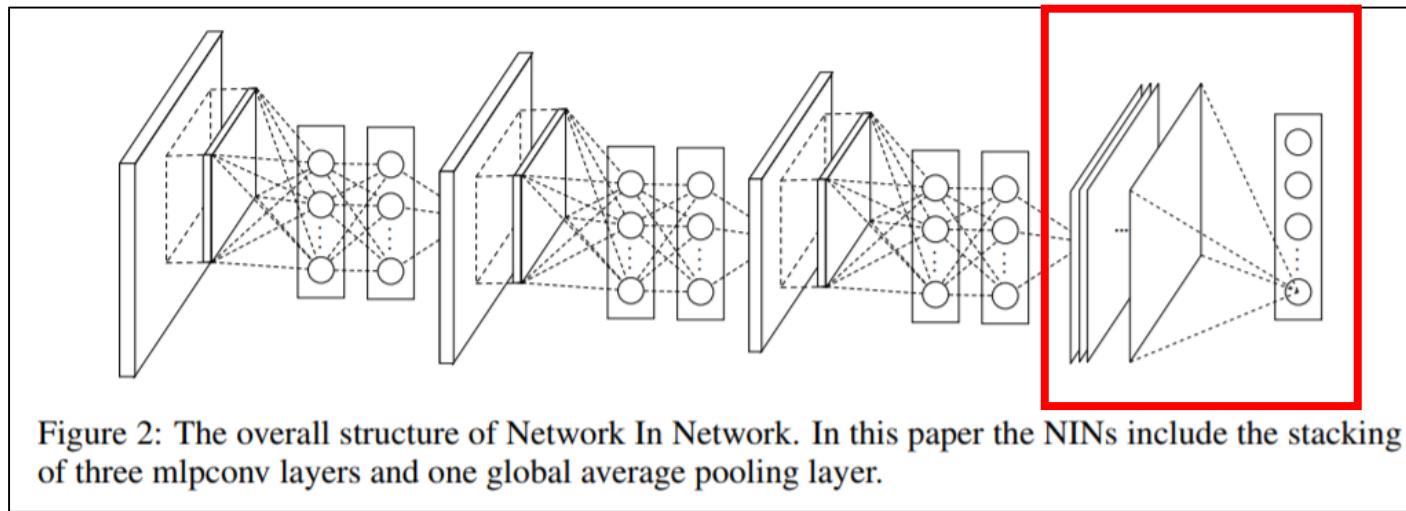
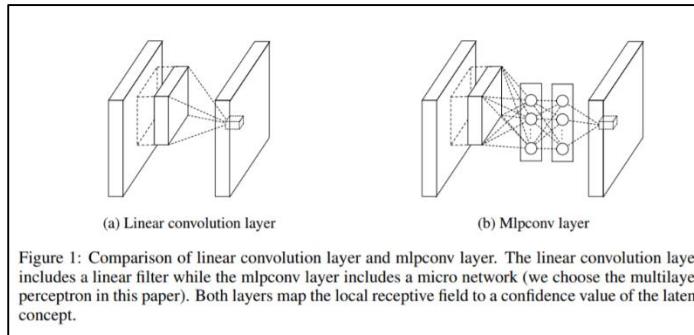
# Introduction & Related Works

- Demo video



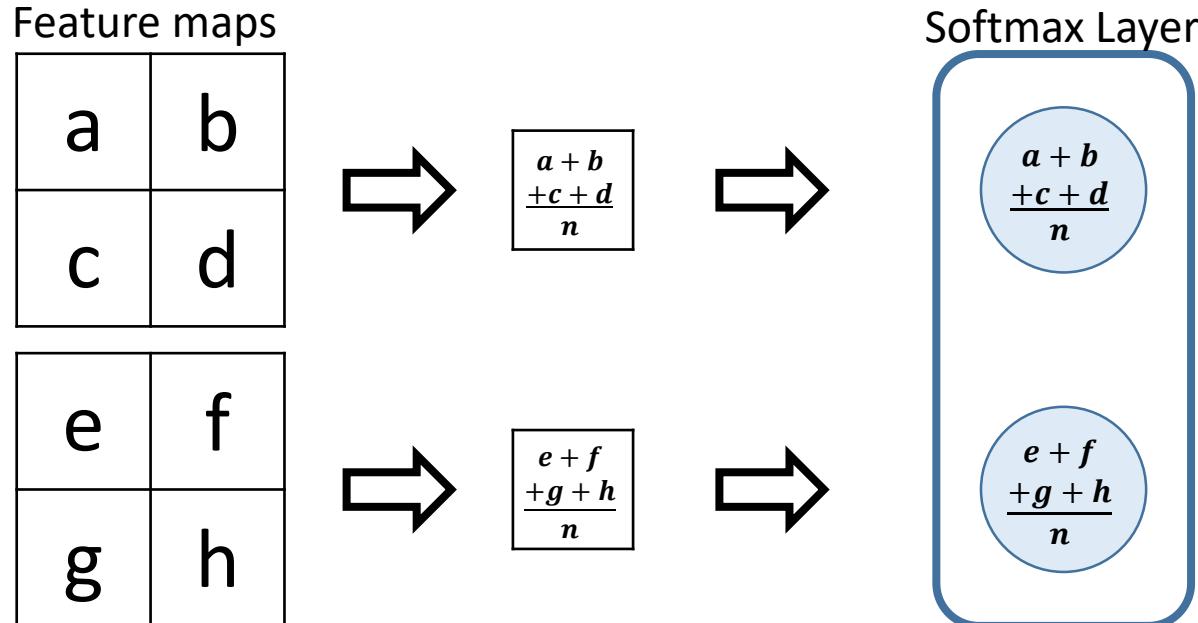
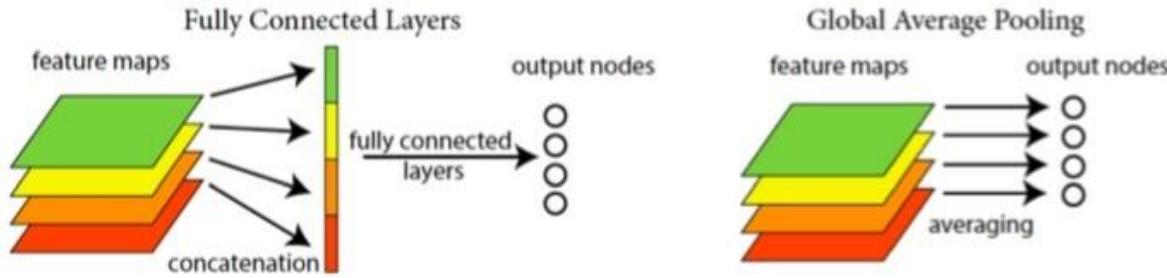
# Introduction & Related Works

- Global Average Pooling
  - Lin, Min, Qiang Chen, and Shuicheng Yan. "Network in network." *arXiv preprint arXiv:1312.4400* (2013).



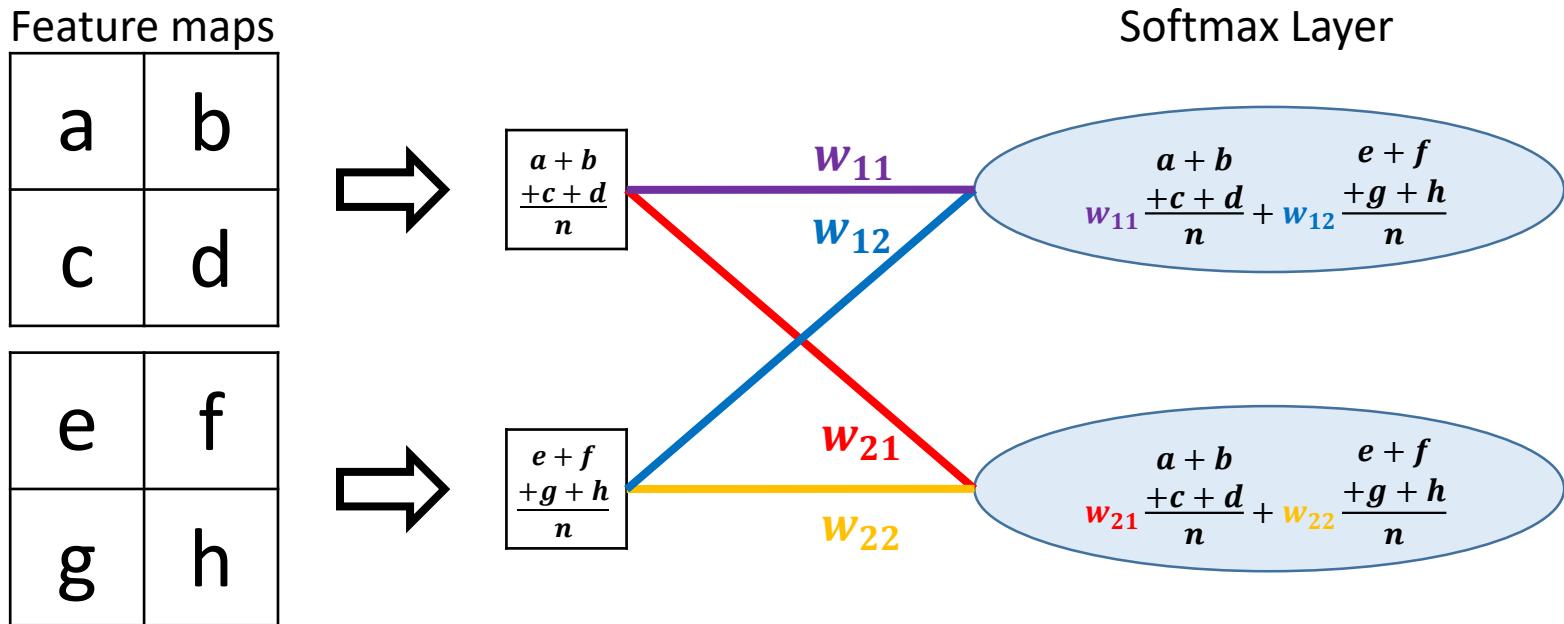
# Introduction & Related Works

- Global Average Pooling
  - Without FC layer?



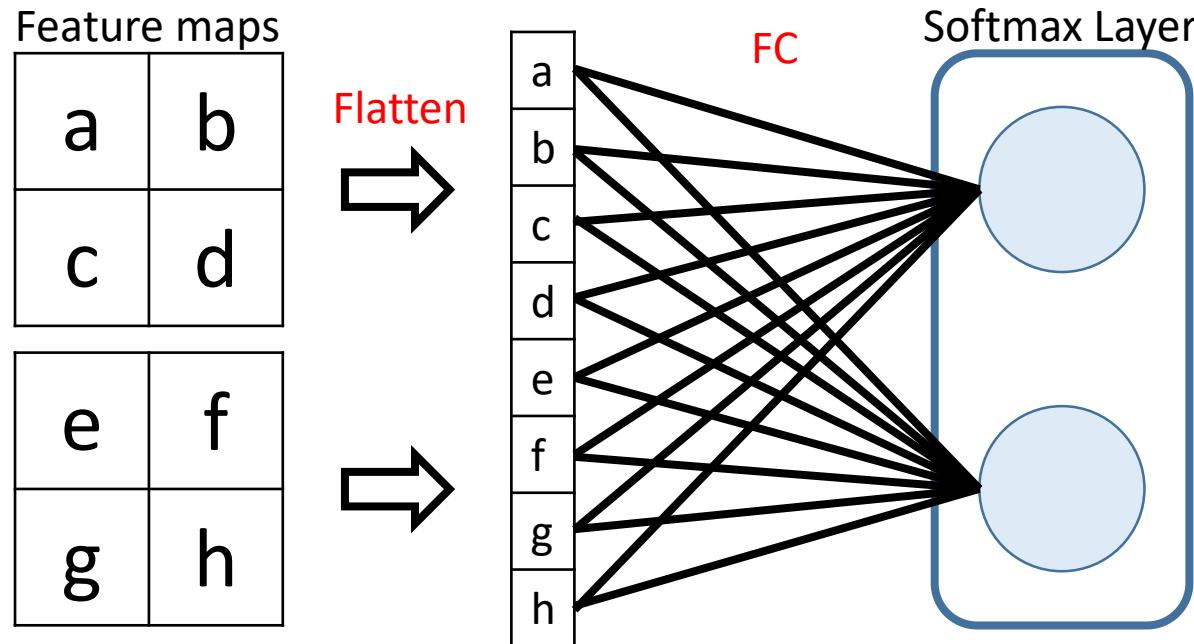
# Introduction & Related Works

- Global Average Pooling
  - With FC layer
  - The size of feature map:  $[2, 2, 2] \rightarrow [1, 1, 2] = 2$
  - The number of parameter:  $[2, 2] = 4$



# Introduction & Related Works

- Without Global Average Pooling
- The size of feature map:  $[2, 2, 2] \rightarrow [1, 8]$
- The number of parameter:  $[8, 2] \rightarrow 16$

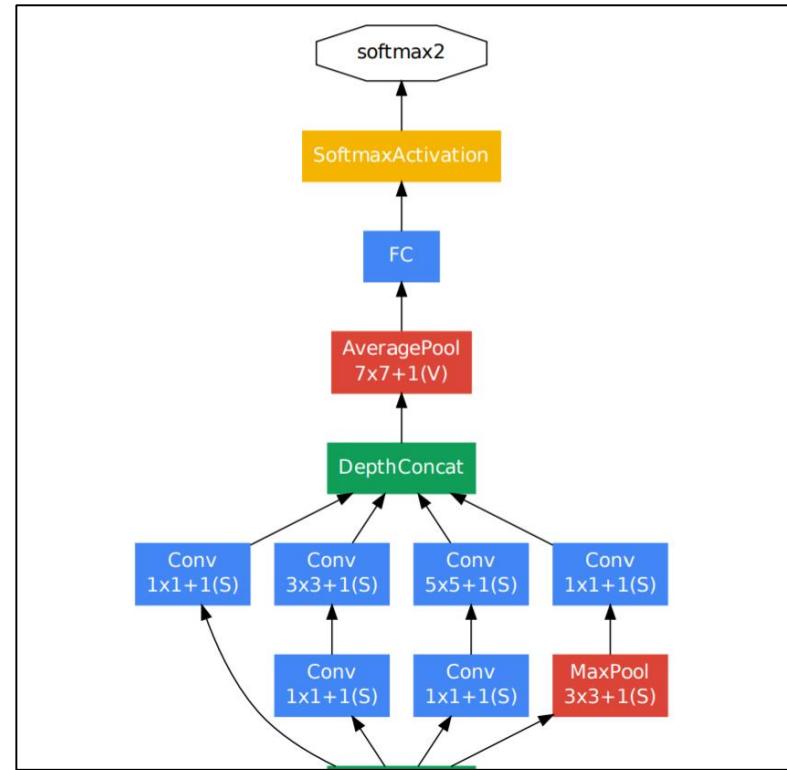
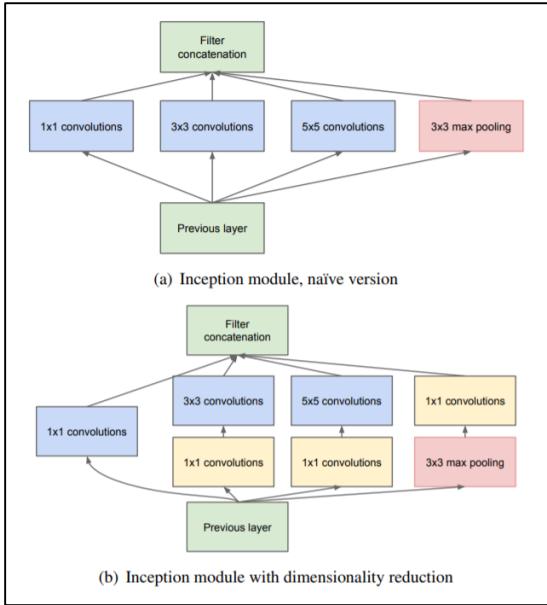
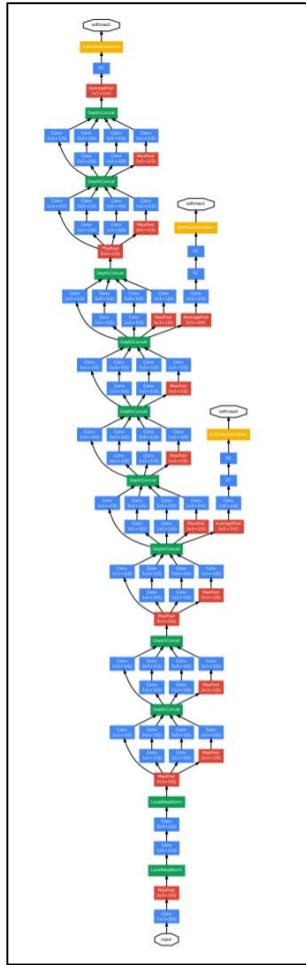


# Introduction & Related Works

- GAP vs GAP with FC vs Without GAP
- GAP
  - Feature map:  $[W, H, C] \rightarrow [1, 1, C] = F$
  - The number of parameter: 0
- GAP with FC
  - Feature map:  $[W, H, C] \rightarrow [1, 1, C] = F$
  - The number of parameter:  $F * \# \text{ of output}$
- Without GAP
  - Feature map:  $[W, H, C] \rightarrow [1, W*H*C] = F$
  - The number of parameter:  $F * \# \text{ of output}$
- **GAP with FC : without GAP  $\propto 1: W*H$**
- If  $W=10, H=10 \rightarrow$  the number of Parameters are 100 Times Different

# Introduction & Related Works

- Google Net & GAP with FC



type	patch size/ stride	output size	depth	#1x1	#3x3 reduce	#3x3	#5x5 reduce	#5x5	pool proj	params	ops
avg pool	7×7/1	1×1×1024	0								
dropout (40%)		1×1×1024	0								
linear		1×1×1000	1							1000K	1M
softmax		1×1×1000	0								

Table 1: GoogLeNet incarnation of the Inception architecture.

# Class Activation Mapping

- Feature map → GAP → FC → Class Activation Map

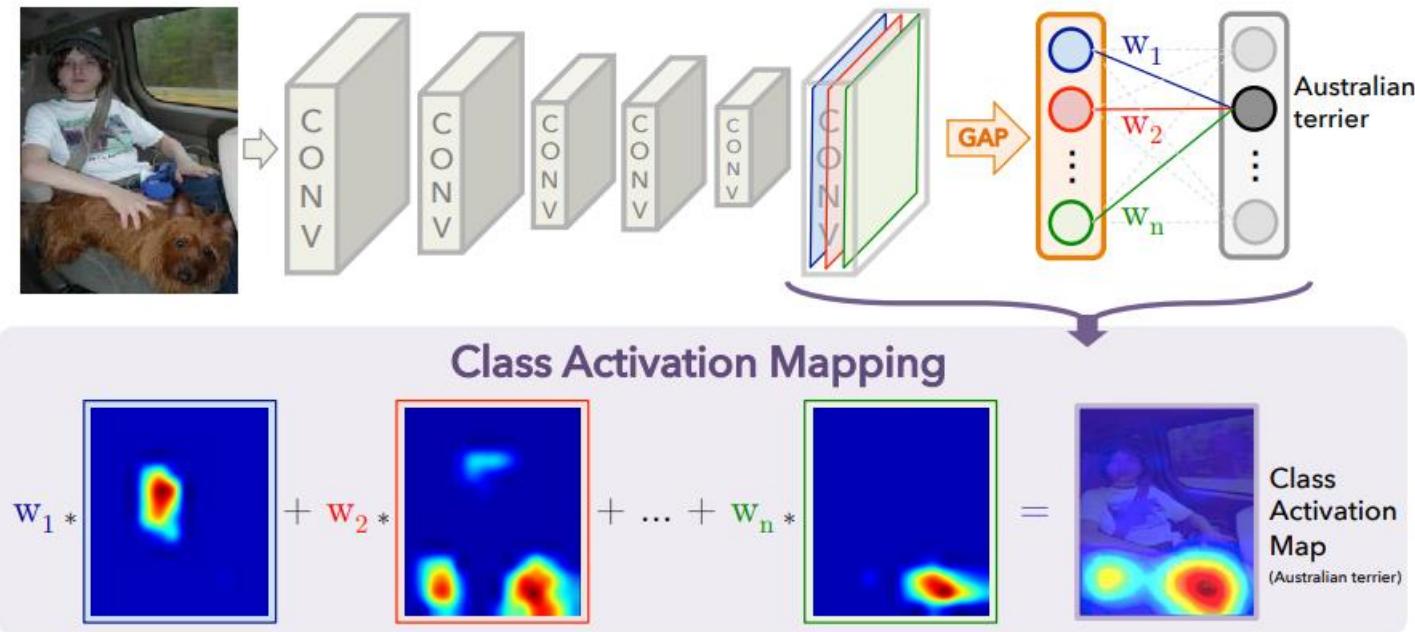


Figure 2. Class Activation Mapping: the predicted class score is mapped back to the previous convolutional layer to generate the class activation maps (CAMs). The CAM highlights the class-specific discriminative regions.

# Class Activation Mapping

For a given image, let  $f_k(x, y)$  represent the activation of unit  $k$  in the last convolutional layer at spatial location  $(x, y)$ . Then, for unit  $k$ , the result of performing global average pooling,  $F^k$  is  $\sum_{x,y} f_k(x, y)$ . Thus, for a given class  $c$ , the input to the softmax,  $S_c$ , is  $\sum_k w_k^c F_k$  where  $w_k^c$  is the weight corresponding to class  $c$  for unit  $k$ . Essentially,  $w_k^c$  indicates the *importance* of  $F_k$  for class  $c$ . Finally the output of the softmax for class  $c$ ,  $P_c$  is given by  $\frac{\exp(S_c)}{\sum_c \exp(S_c)}$ .

Here we ignore the bias term: we explicitly set the input bias of the softmax to 0 as it has little to no impact on the classification performance.

By plugging  $F_k = \sum_{x,y} f_k(x, y)$  into the class score,  $S_c$ , we obtain

$$S_c = \sum_k w_k^c \sum_{x,y} f_k(x, y) = \sum_{x,y} \sum_k w_k^c f_k(x, y). \quad (1)$$

We define  $M_c$  as the class activation map for class  $c$ , where each spatial element is given by

$$M_c(x, y) = \sum_k w_k^c f_k(x, y). \quad (2)$$

Thus,  $S_c = \sum_{x,y} M_c(x, y)$ , and hence  $M_c(x, y)$  directly indicates the importance of the activation at spatial grid  $(x, y)$  leading to the classification of an image to class  $c$ .

# Weakly-supervised Object Localization

- Classification

Table 1. Classification error on the ILSVRC validation set.

Networks	top-1 val. error	top-5 val. error
VGGnet-GAP	33.4	12.2
GoogLeNet-GAP	35.0	13.2
AlexNet*-GAP	44.9	20.9
AlexNet-GAP	51.1	26.3
GoogLeNet	31.9	11.3
VGGnet	31.2	11.4
AlexNet	42.6	19.5
NIN	41.9	19.6
GoogLeNet-GMP	35.6	13.9

- Localization

Table 2. Localization error on the ILSVRC validation set. *Backprop* refers to using [23] for localization instead of CAM.

Method	top-1 val.error	top-5 val. error
GoogLeNet-GAP	<b>56.40</b>	<b>43.00</b>
VGGnet-GAP	57.20	45.14
GoogLeNet	60.09	49.34
AlexNet*-GAP	63.75	49.53
AlexNet-GAP	67.19	52.16
NIN	65.47	54.19
Backprop on GoogLeNet	61.31	50.55
Backprop on VGGnet	61.12	51.46
Backprop on AlexNet	65.17	52.64
GoogLeNet-GMP	57.78	45.26

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 [23]	weakly	46.4
GoogLeNet [25]	full	26.7
OverFeat [22]	full	29.9
AlexNet [25]	full	34.2

# Weakly-supervised Object Localization

- Object detection

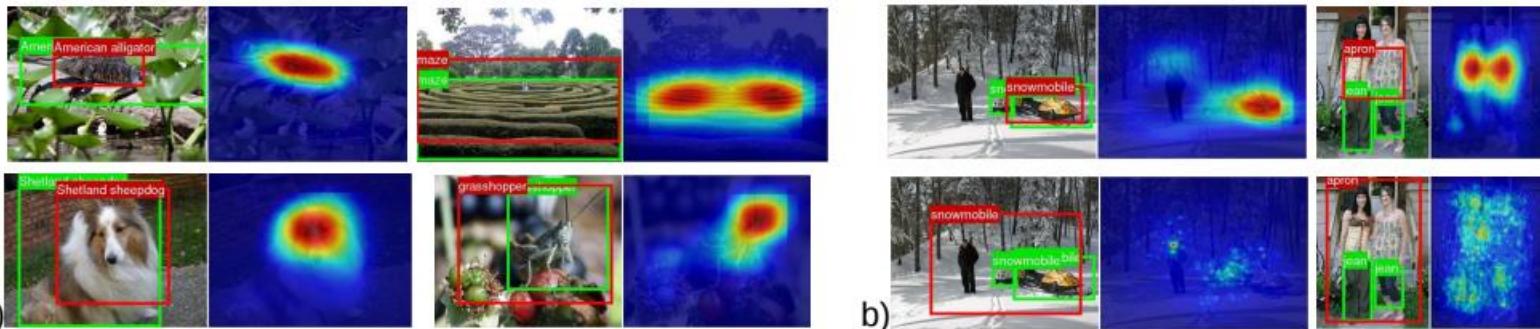


Figure 6. a) Examples of localization from GoogleNet-GAP. b) Comparison of the localization from GooleNet-GAP (upper two) and the backpropagation using AlexNet (lower two). The ground-truth boxes are in green and the predicted bounding boxes from the class activation map are in red.

- Weakly supervised text detector

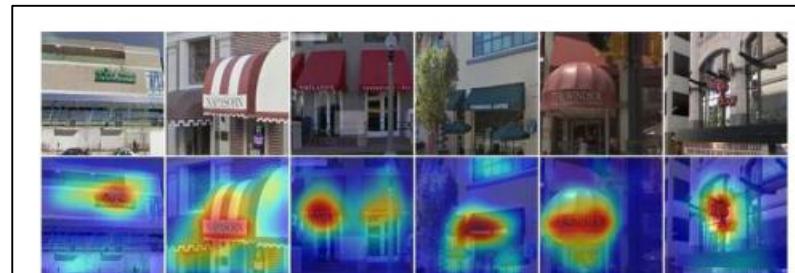


Figure 11. Learning a weakly supervised text detector. The text is accurately detected on the image even though our network is not trained with text or any bounding box annotations.

# Conclusion

- Class Activation Mapping(CAM) with Global Average Pooling(GAP)
- Weakly-supervised object localization → CNNs trained for classification can do object localization
- Various experiments are very important!
- Reducing uncertainty in deep learning (NRF Global ph.d research topic)

Thank you.  
Q & A