

Session -3

# CNN UNDER THE HOOD

March 15,  
2024

Me using Image Augmentation  
to increase the training data

\*Flips image horizontally\*



Deep learning model:



## AGENDA

- ① Intuition behind multi-channel Convolution
- ② Global Average Pooling
- ③ One x One Convolution
- ④ Dilated Convolution
- ⑤ Image Augmentation.

## Receptive Field

		<u>400 x 400 x 1</u>	
$3 \times 3 L_1$	$\rightarrow$	$3 \times 3 \times 32$ <u>398</u>	( <u>288</u> )
$5 \times 5 L_2$	$\rightarrow$	$3 \times 3 \times 64$ <u>326</u>	( <u>576</u> )
$7 \times 7 L_3$	$\rightarrow$	$3 \times 3 \times 128$ <u>394</u>	( <u>1152</u> )
$9 \times 9 L_4$	$\rightarrow$	$3 \times 3 \times 256$ <u>392</u>	( <u>2304</u> )
$11 \times 11 L_5$	$\rightarrow$	$3 \times 3 \times 512$ <u>396</u>	( <u>4608</u> )
<u>Max - Pooling</u>			
		<u>39012</u>	

# INTUITION - MULTI CHANNEL CONV

CASE-I Multispectral Imaging in agriculture.

→ Capture Image data at specific wavelength across Electromagnetic spectrum.

→ Eg. Following spectrum

→ RGB spectrum

\* For visual information what humans see

→ Infrared Spectrum

\* Healthy plants reflect more near infrared

NDVI (Normalized Diff)  
Vegetation Index

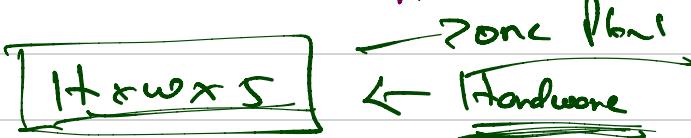
(NIR  
light.)

→ Shortwave Infrared (SWIR)

\* For water content in soil.

\* Water absorbs SWIR signals strongly

so it appears darker.



→ For what kernels might look for.

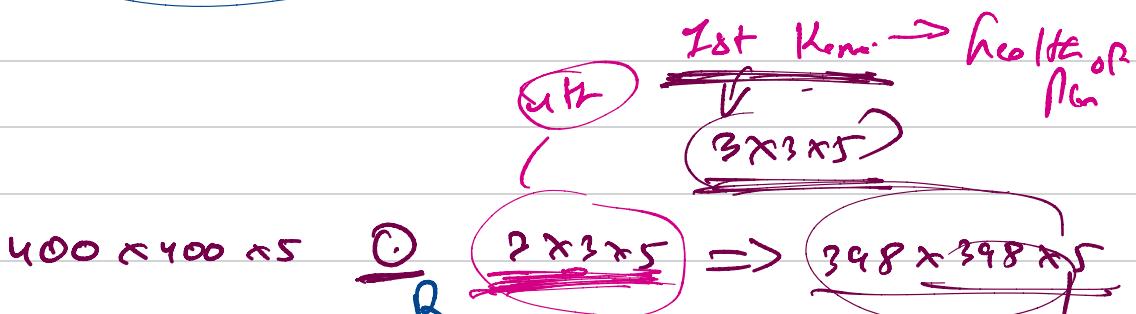
(i) Signs of Stress

(ii) Optimizing water use

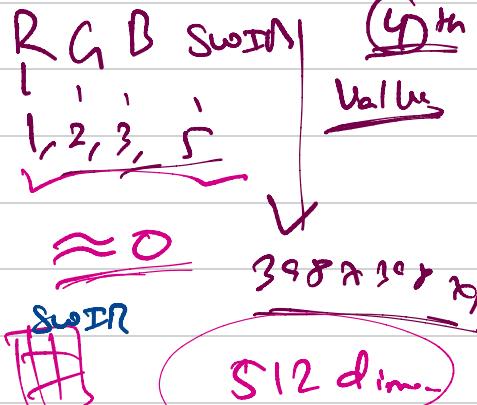
(iii) Crop health

(iv) Yield prediction.

398x398x32



0.1	0.7	0.2
5	7	8
1	2	3

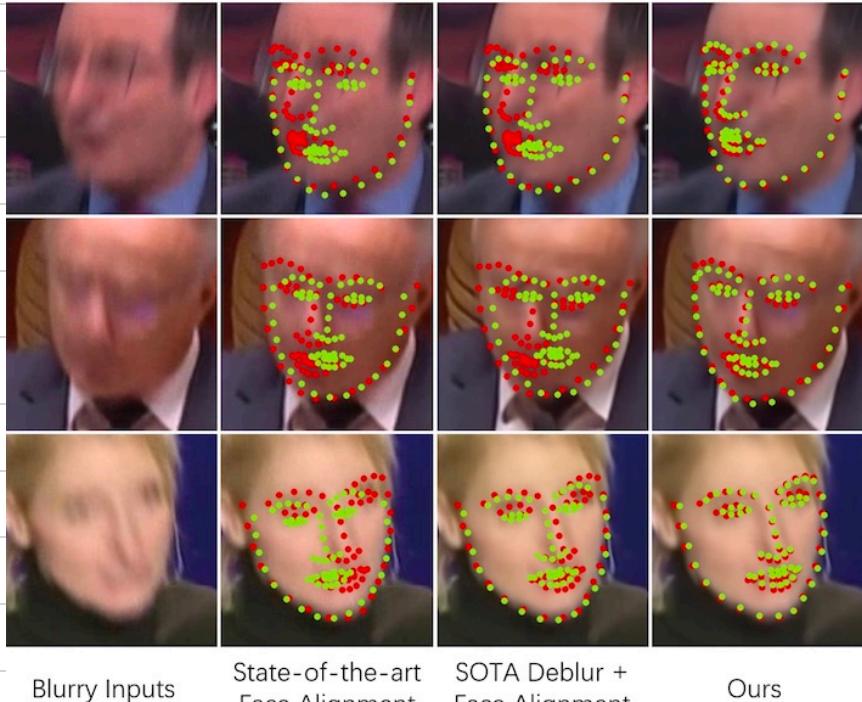


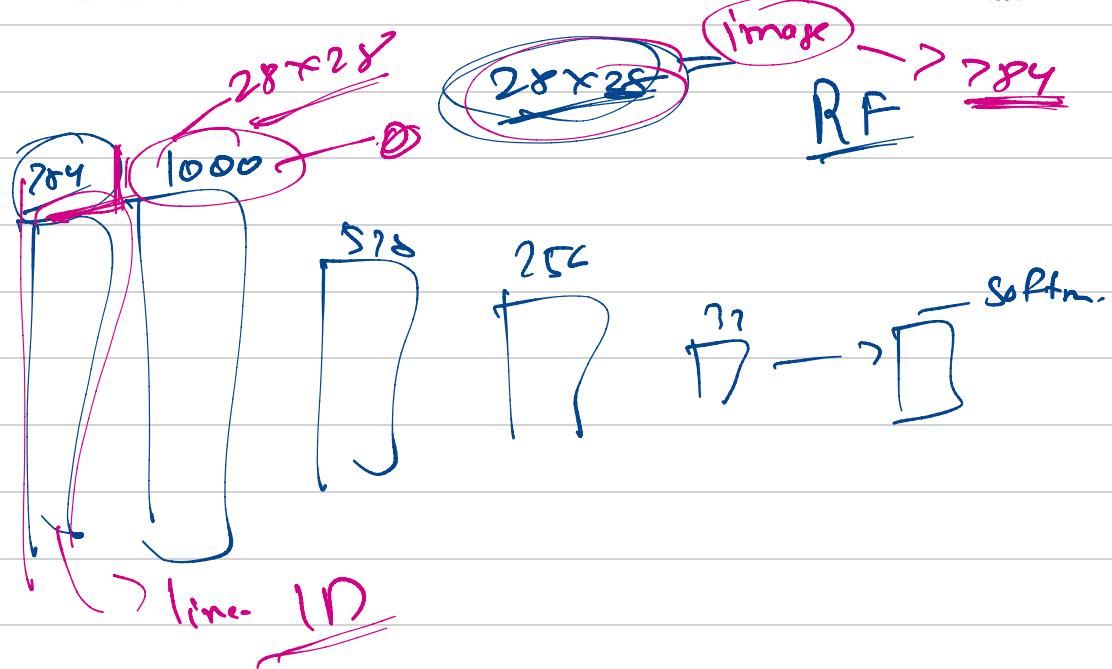
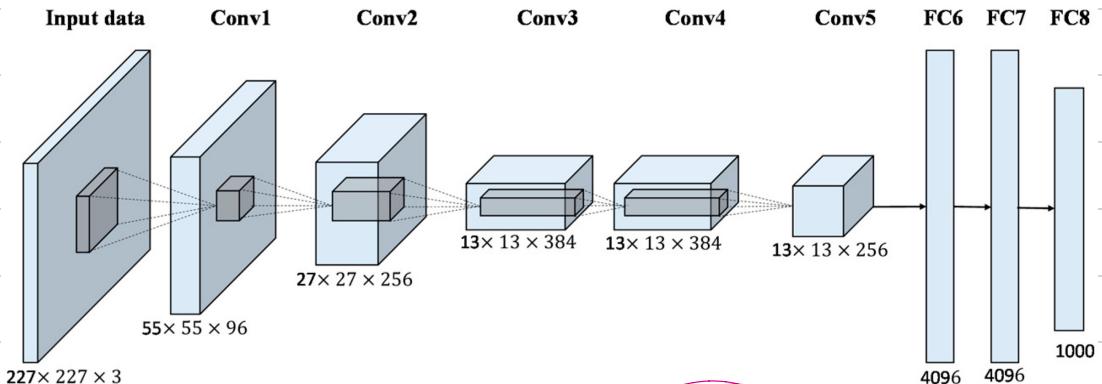
(L) 5 7

512 dim

● Ground Truth

● Predict Results





# Global Average Pooling

$224 \times 224 \times 3$

Layer (type)	Output Shape	Param #
conv2d (Conv2D)	(None, 224, 224, 32)	896
conv2d_1 (Conv2D)	(None, 224, 224, 32)	9248
max_pooling2d (MaxPooling2D)	(None, 112, 112, 32)	0
conv2d_2 (Conv2D)	(None, 112, 112, 64)	18496
conv2d_3 (Conv2D)	(None, 112, 112, 64)	36928
max_pooling2d_1 (MaxPooling2D)	(None, 56, 56, 64)	0
conv2d_4 (Conv2D)	(None, 56, 56, 128)	73856
conv2d_5 (Conv2D)	(None, 56, 56, 128)	147584
max_pooling2d_2 (MaxPooling2D)	(None, 28, 28, 128)	0
conv2d_6 (Conv2D)	(None, 28, 28, 256)	295168
conv2d_7 (Conv2D)	(None, 28, 28, 256)	590080
max_pooling2d_3 (MaxPooling2D)	(None, 14, 14, 256)	0
conv2d_8 (Conv2D)	(None, 14, 14, 512)	1180160
conv2d_9 (Conv2D)	(None, 14, 14, 512)	2359808
max_pooling2d_4 (MaxPooling2D)	(None, 7, 7, 512)	0
flatten (Flatten)	(None, 25088) $\xrightarrow{7 \times 7 \times 512}$	$0 + 1024$
dense (Dense)	(None, 1024)	25691136
dense_1 (Dense)	(None, 1024)	1049600
dense_2 (Dense)	(None, 10)	10250
<hr/>		
Total params:	31463210 (120.02 MB)	
Trainable params:	31463210 (120.02 MB)	
Non-trainable params:	0 (0.00 Byte)	

Dog vs Cat

Zt ~ Dog

2n ~ Cat

$25088 \times 1024$

Dog's Pic

7x 7x 512      (Dog's img)  
                ↓  
Cat's nose

For Cat's Features Kernel ??  $\rightarrow 0$

7x 7x 512  $\rightarrow$  1x1 x 512  
                        ↓  
                        0

Dog's nose  $\rightarrow$  Dog  $\gg 0$

{ 5, 6, 78, 0, 0.2, 0.01 }

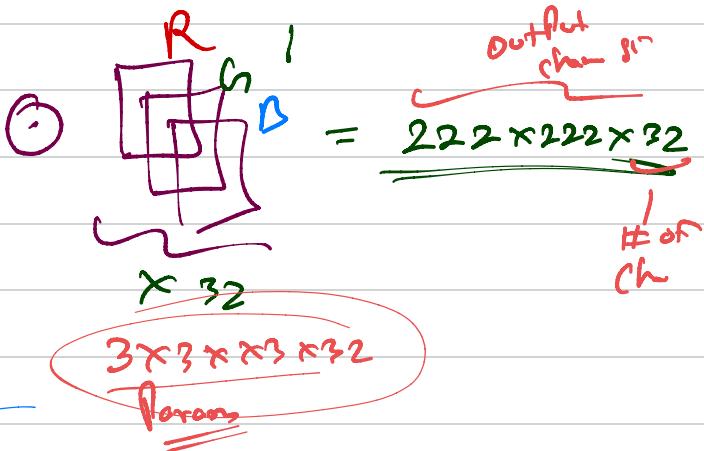
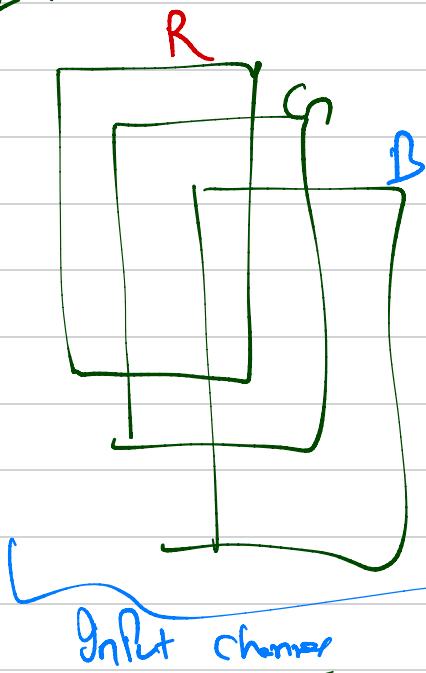
512  
+  
softmax  
Feat  
Dense (10, softmax)  
a, b, c - - - - - 512

$$I_{\text{image}} = H, W, \# \text{ of channels}$$

$H = 224, W = 224, \# \text{ of channels} = 3$

Z im.

Everyth.  $\rightarrow$  channel



$$400 \times 400 \times 3 \quad (3 \times 3 \times 3) \times 37$$



$$\underline{\underline{398 \times 398 \times 32}}$$

$32, \underline{32}$

$4,6 / Y, B$

# One x One Convolution

400x400x3 | (3x3x3)x32 | 398x398x32 RF of 3x3

398x398x32 | (3x3x32)x64 | 396x396x64 RF of 5X5

396x396x64 | (3x3x64)x128 | 394x394x128 RF of 7X7

394x394x128 | (3x3x128)x256 | 392x392x256 RF of 9X9

392x392x256 | (3x3x256)x512 | 390x390x512 RF of 11X11

## MaxPooling

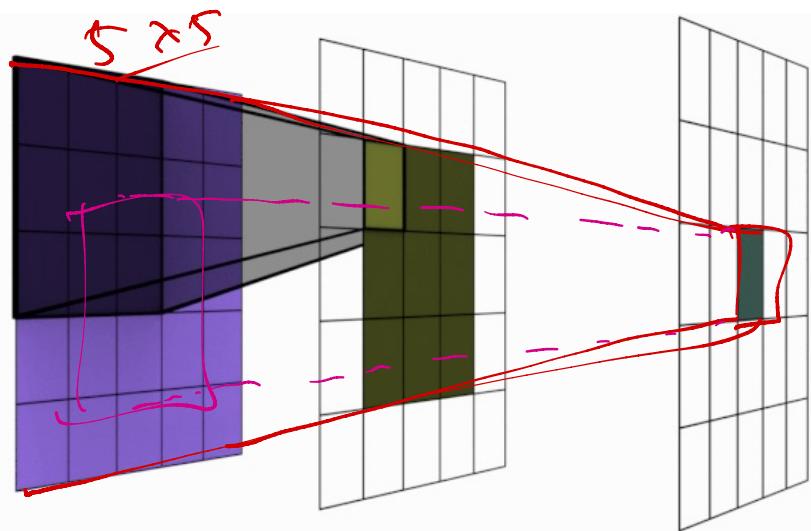
195x195x512 | (?x?x512)x32 | ?x?x32 RF of 22x22

.. 3x3x32x64 RF of 24x24

.. 3x3x64x128 RF of 26x26

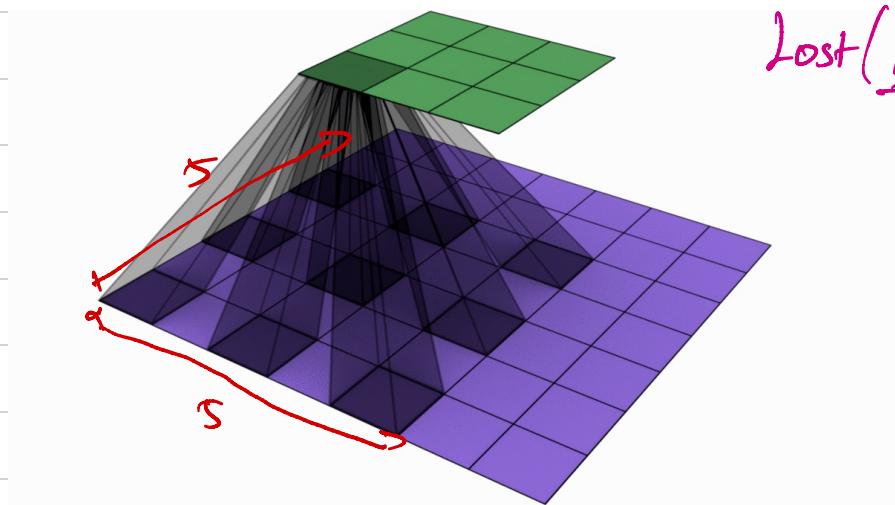
.. 3x3x128x256 RF of 38x28

.. 3x3x256x512 RF of 30x30



# DILATED CONVOLUTION

Q. How to get RF of  $5 \times 5$  with one  $3 \times 3$  Kernel!



Lost/Filtering

- Dilated convolution increases receptive field of the network exponentially and linear parameter accretion
- It's used in application where we care more about integrating knowledge of the broader content with less cost.

→ Key application the dilated conv authors have in mind  
is **dense prediction**: vision applications where predicted object has a similar size & structure to input image.

## Semantic Segmentation vs. Instance Segmentation vs. Panoptic Segmentation



(a) Image



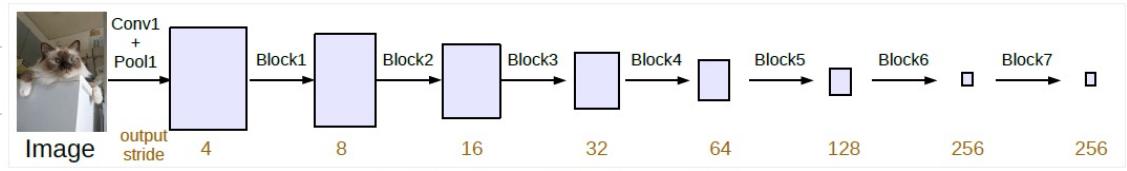
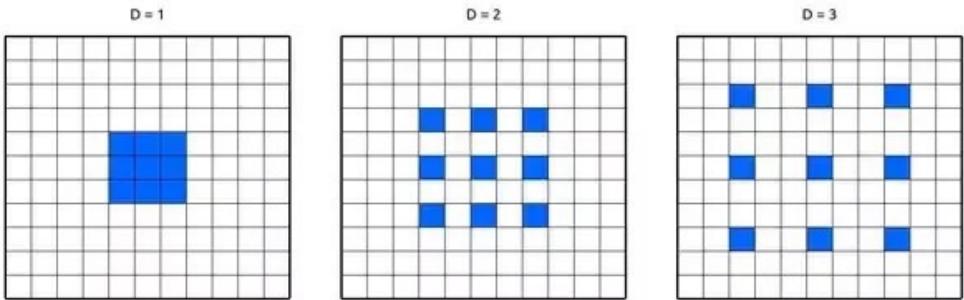
(b) Semantic Segmentation



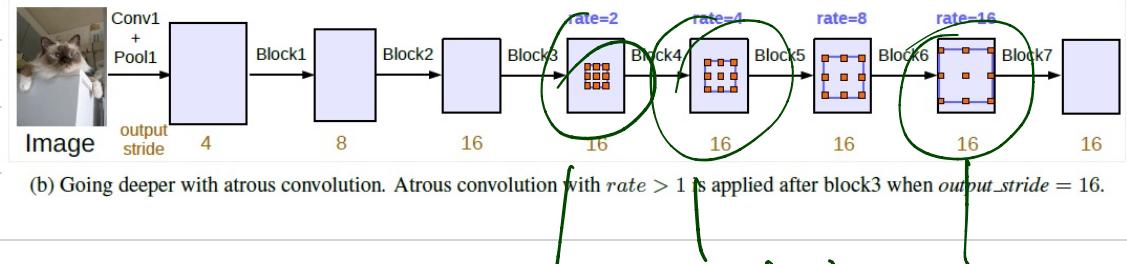
(c) Instance Segmentation



(d) Panoptic Segmentation



(a) Going deeper without atrous convolution.



microscale. normalcy  
Telescope

less non-linear combination  
of features

In +  $2 \times 0$

# Image Segmentation

Lower Level Segmentation

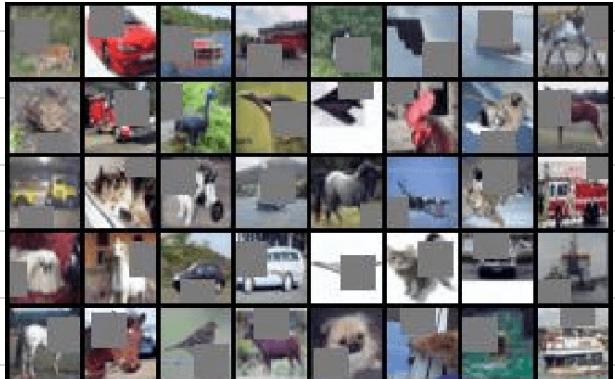
most of it is  
in built in keras

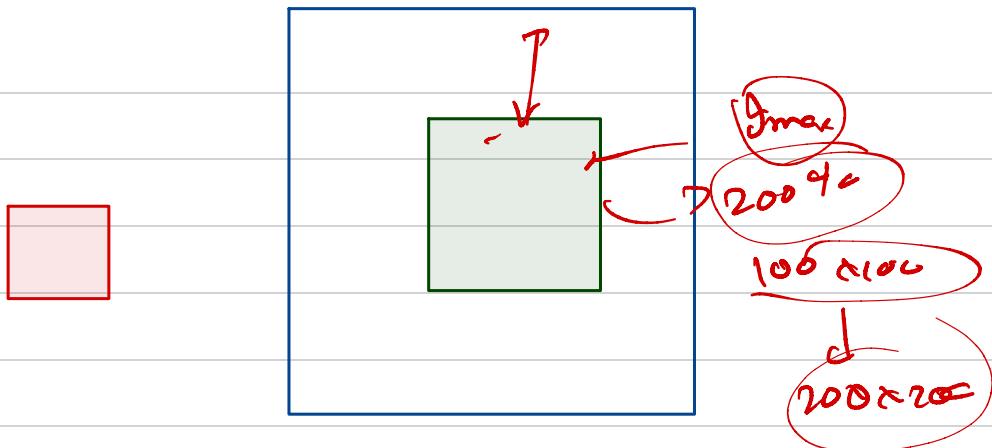


Cut Off - Higher Levels

V.V. Effective.

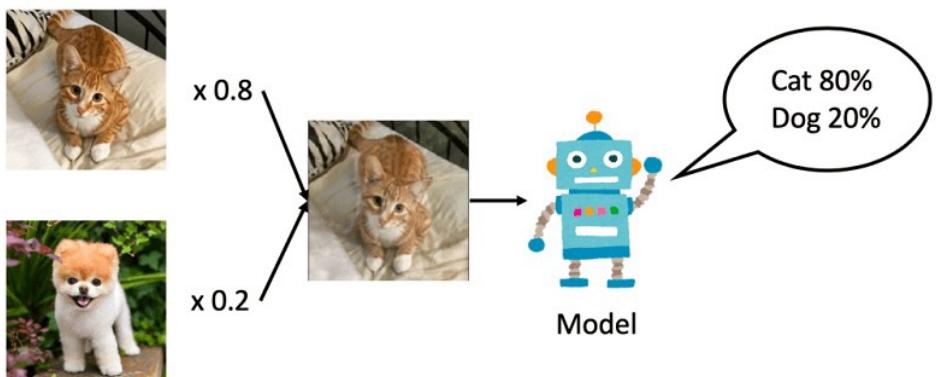
All bumentatio





(Scientist with more facial Expressions)

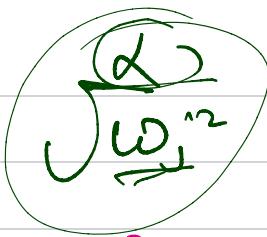
Mix UP - Higher levels



Custom Implementation needed.

Tough but Effective.

Adam



Reduce LOR or Plot

5 epochs

if val-loss ??

$\rho_{\text{final}} = 0$

$L_P = L_{R=0.3}$