

March 29, 2023

# Image Similarity: Understanding Embeddings



DSML: Computer Vision.

Embeddings Class starts  
@ 9:05 pm.

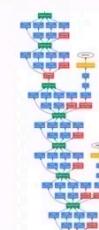
**What normal people see  
when they walk on street**



**What Computer Vision  
folks see**



**WHO WOULD WIN?**



**STATE OF THE ART  
NEURAL NETWORK**



**ONE NOISY BOI**

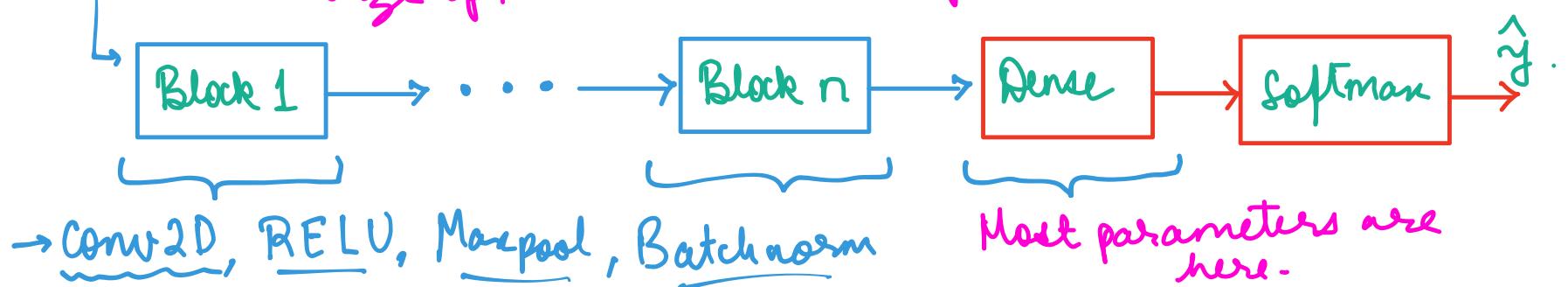
## Recap:

- \* Problem: Classification.
- \* Approach:
  - (a) Preprocess Image Data.
  - (b) Design an architecture.
  - (c) Train the architecture.
  - (d) Enjoy the awesome results given by the CNN.

## Architecture:

- { (a) AlexNet - CNNs with ReLU on a GPU.
- { (b) VGG - Convolution blocks with  $3 \times 3$  convolutions.

Size of the model  $\approx$  100s of MBs.



## Agenda:

- \* Study 2 new classification architectures:
  - Inception.
  - ResNet.
- \* Use a trained classifier for something else !!

### REVERSE IMAGE SEARCH:

given an image, find other images which look like this one.

Google ders

Pinterest

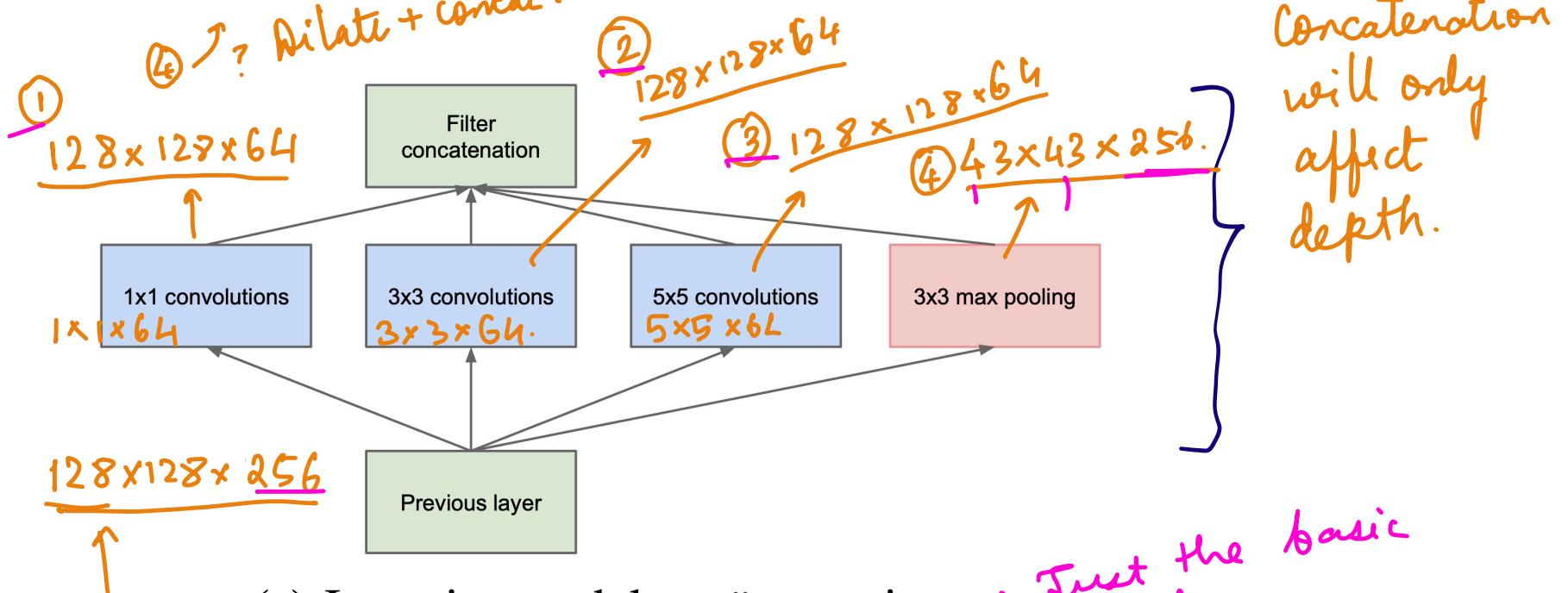
Celebs who look like me.

## Inception Network

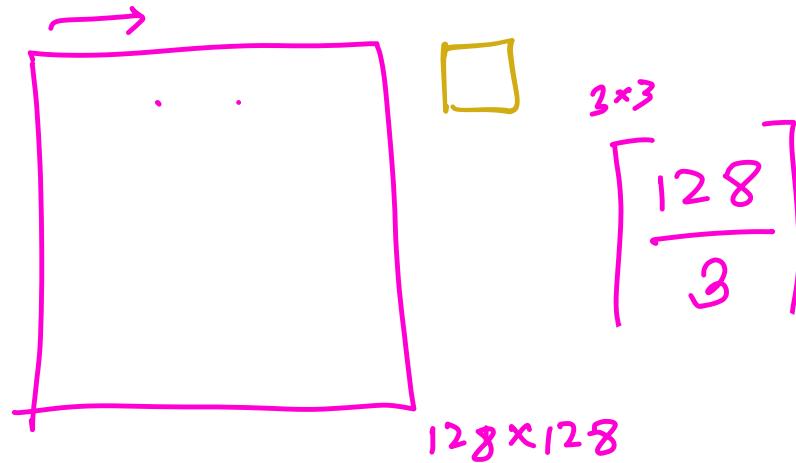
→ google → ImageNet 2014 challenge.

This network did slightly better than the VGG-19 network.

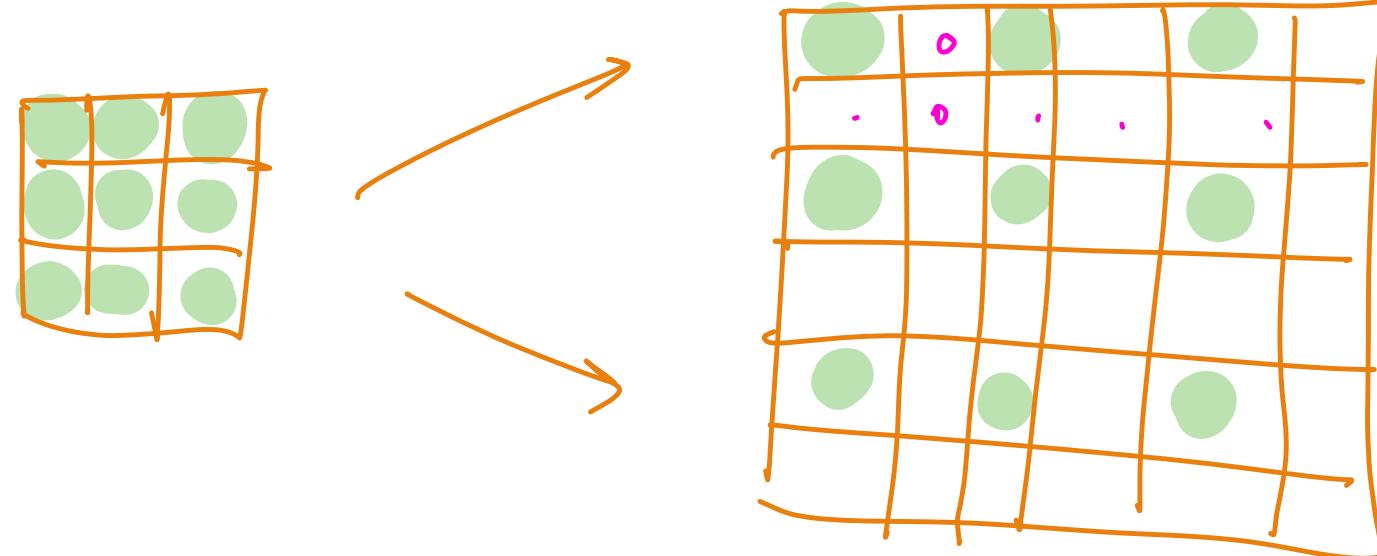
It used only 11M parameters while VGG-19 ~130M parameters.



(a) Inception module, naive version → Just the basic idea.

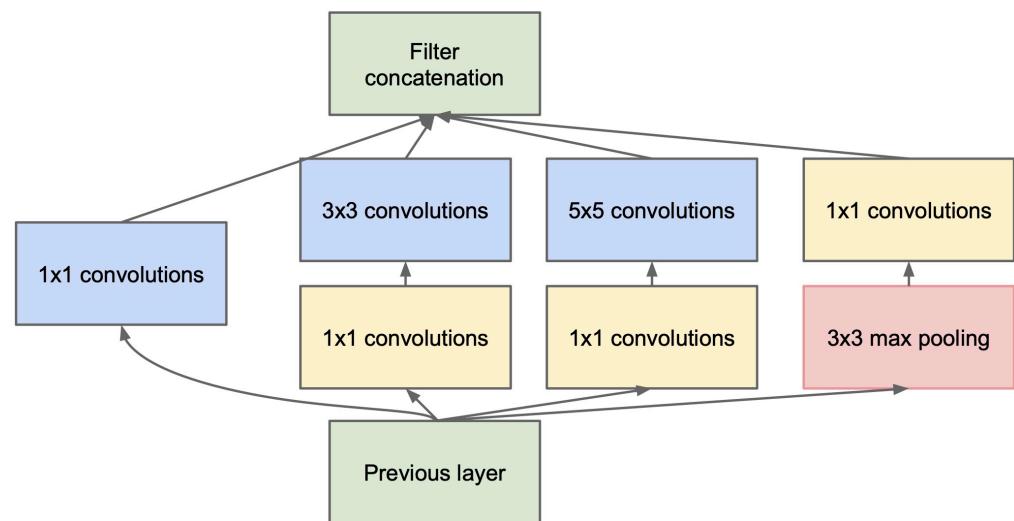
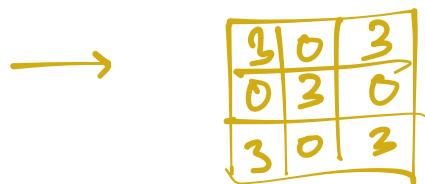
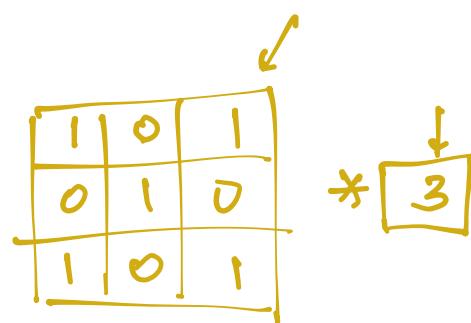
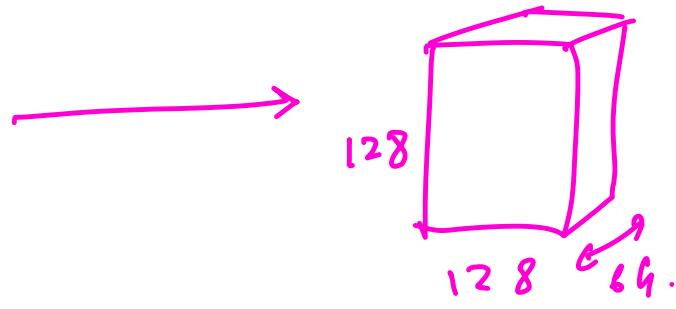
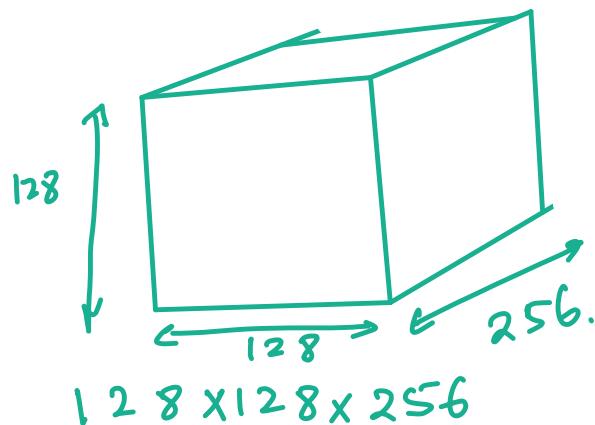


$$\begin{bmatrix} 4 & 2 & 6 & 6 \end{bmatrix} = 43 \times 43$$

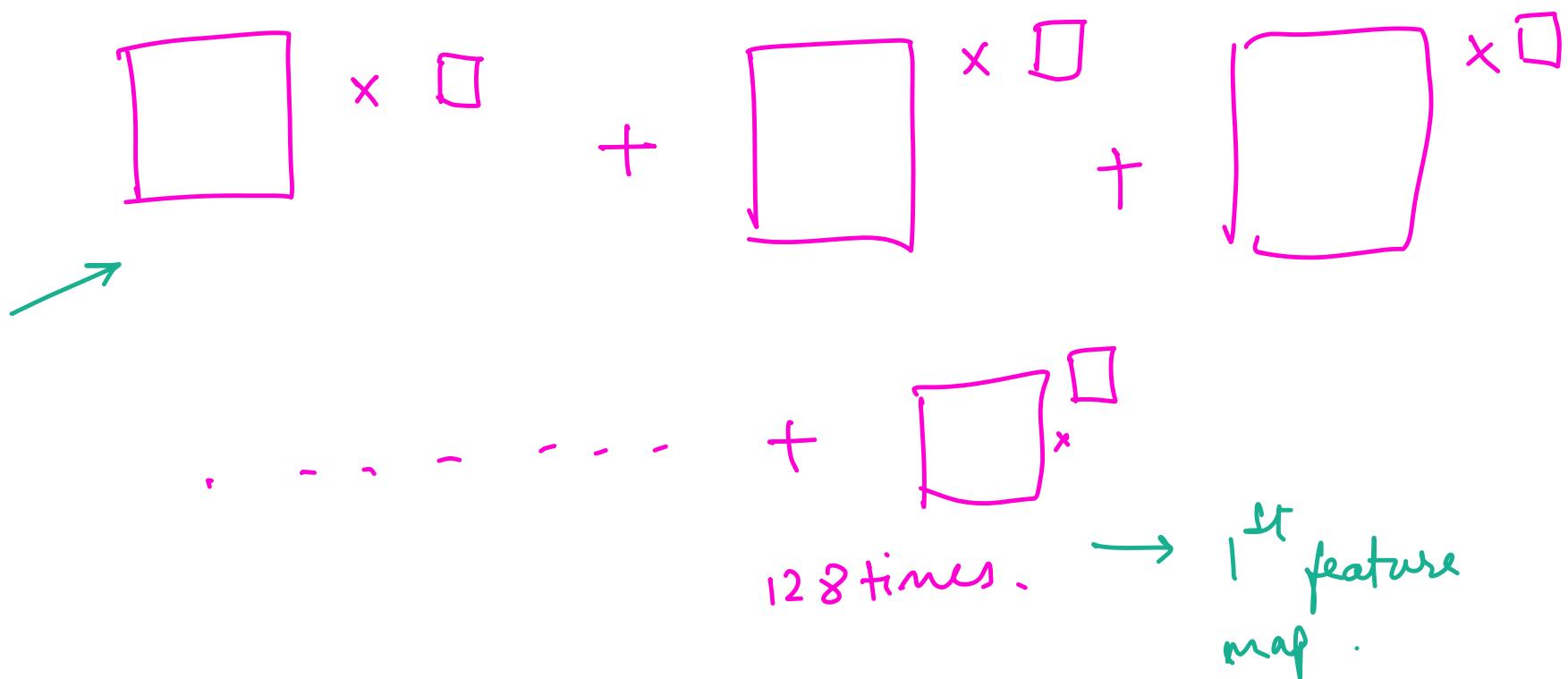
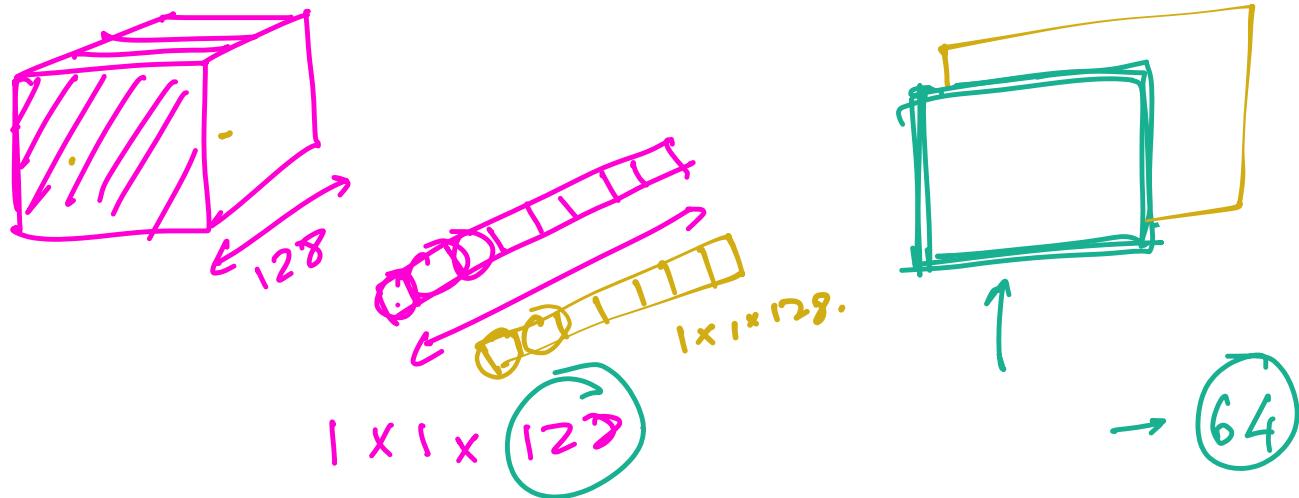


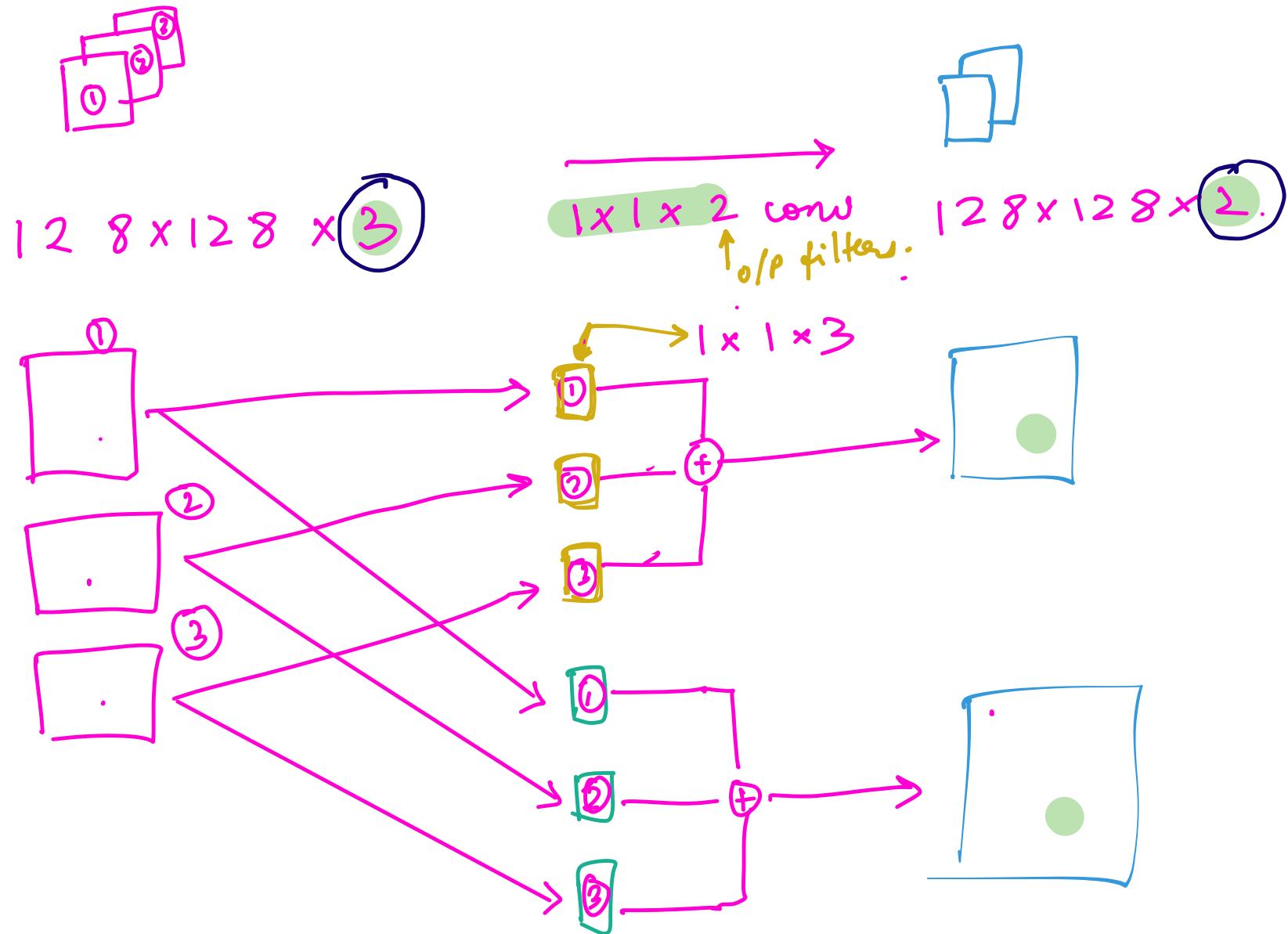
What are  $1 \times 1$  convolutions?

→  $1 \times 1$  convolutions help us reduce the filter depth  
information loss is minimized.

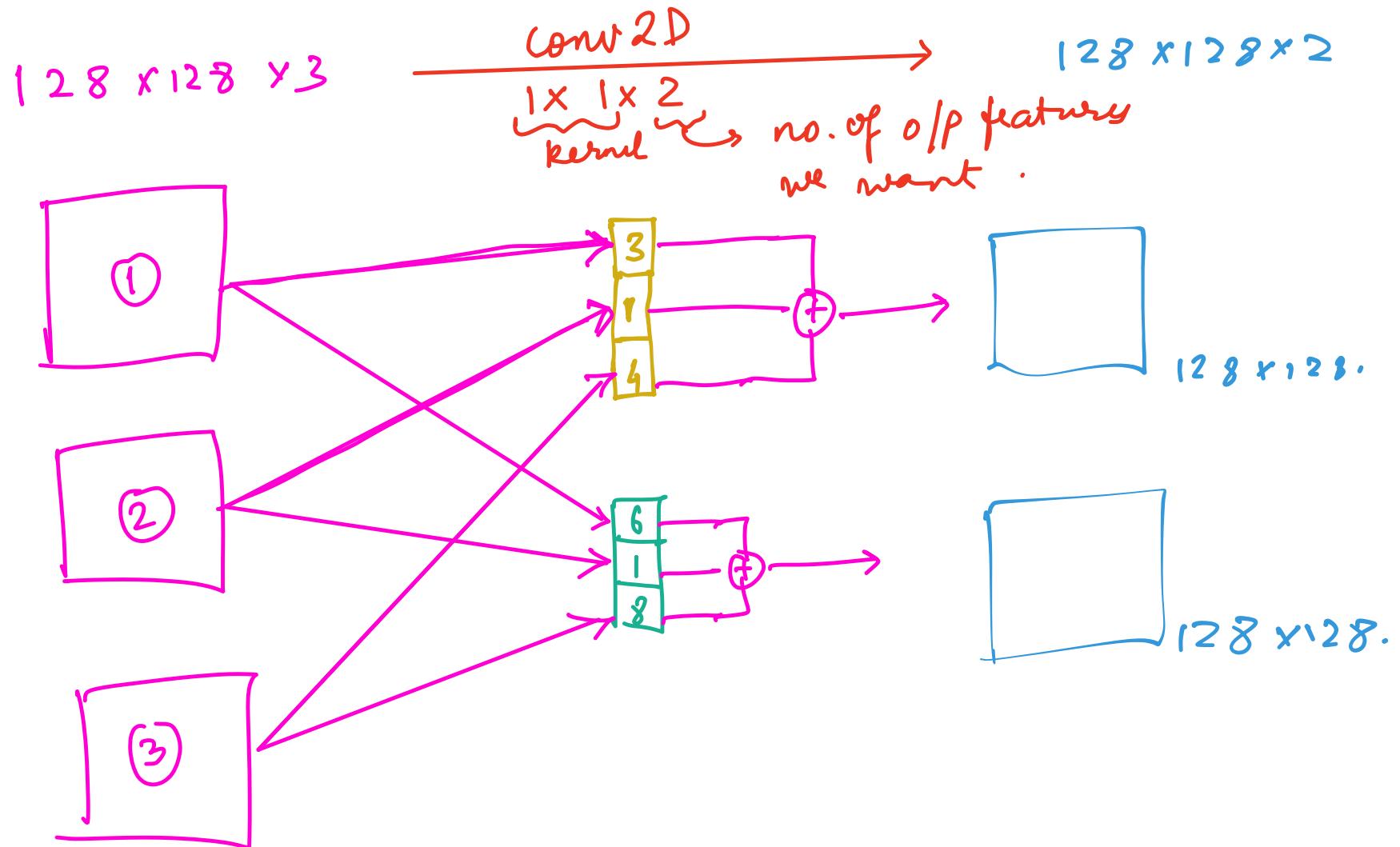


(b) Inception module with dimension reductions





$1 \times 1$  convolutions are nothing but weighted averages!!



Q] Why do  $1 \times 1$  convolutions when  
the same reduction can be  
done with a  $3 \times 3$  kernel?

$128 \times 128$

Input:  $m$  channels

Output:  $n$  channels.

No. of  
parameters

No. of mul + add  
operations.

$3 \times 3$

$9mn$

$128 \times 128 \times 9mn$

$1 \times 1$

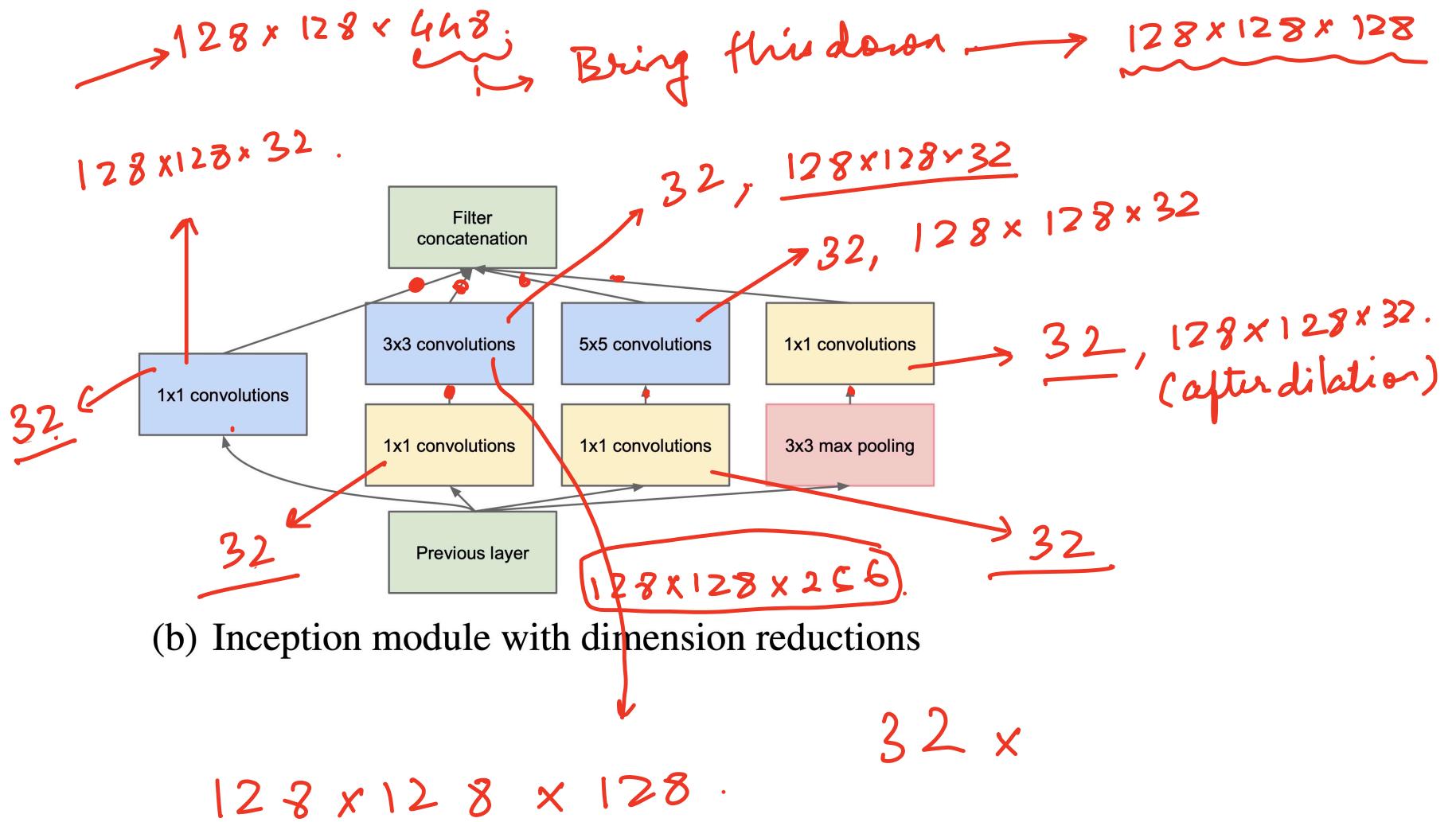
$mn$



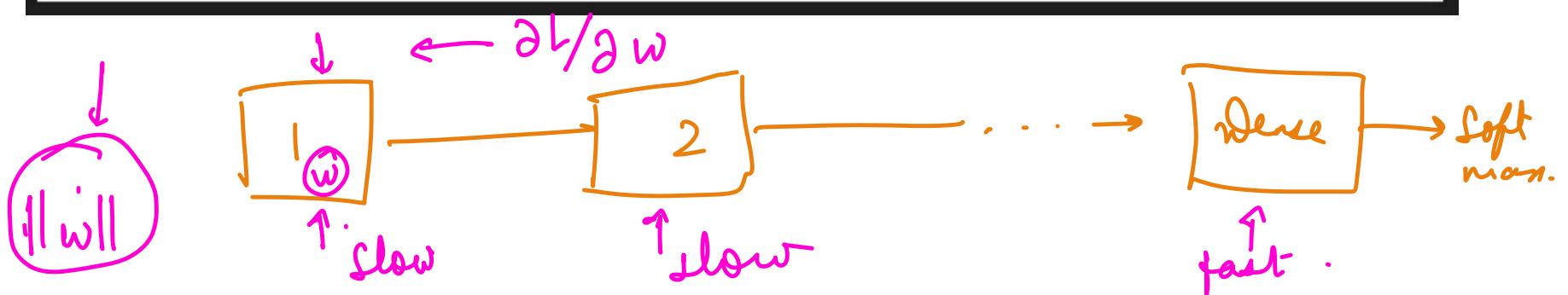
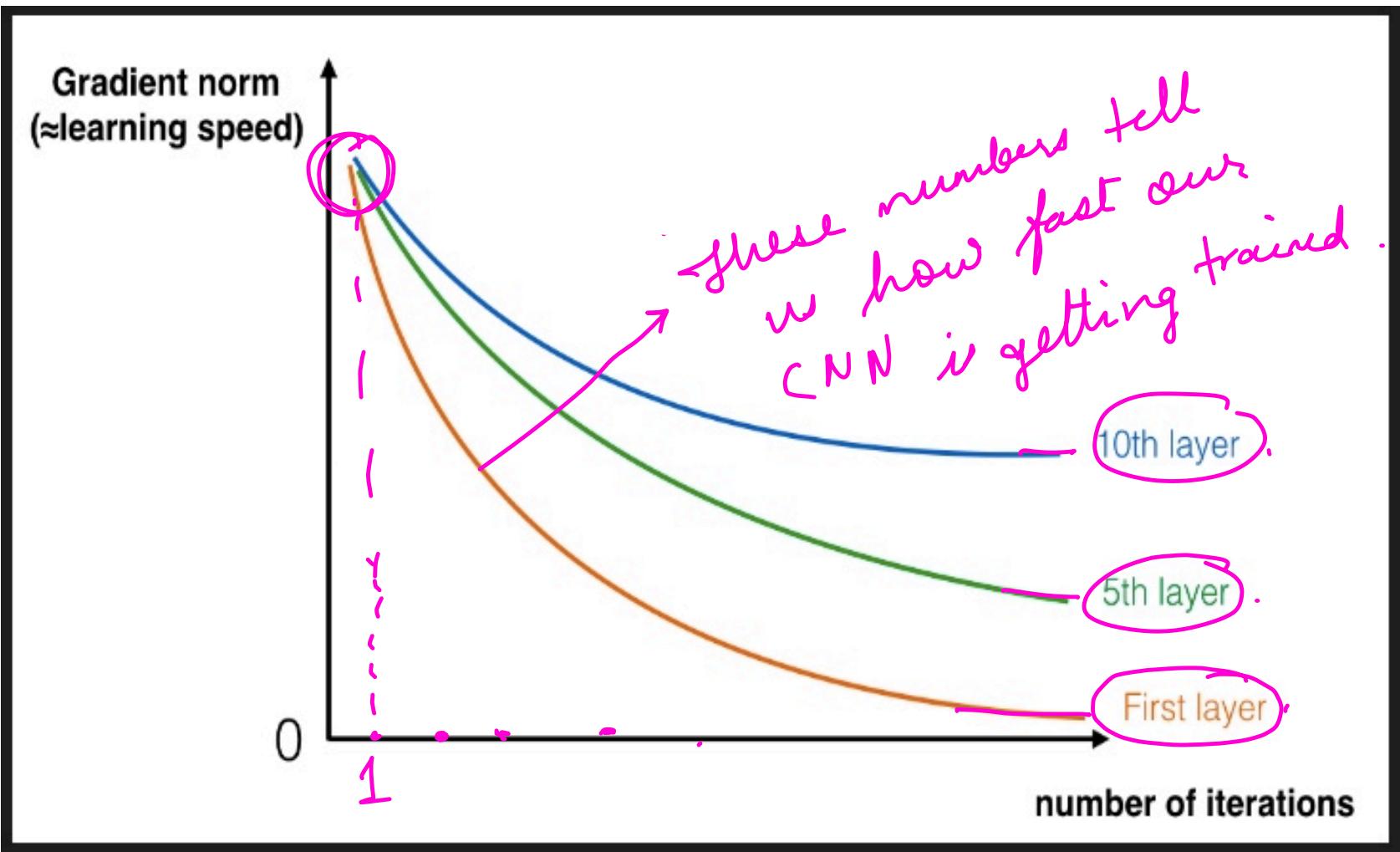
$128 \times 128 \times mn$



Better space, time.  
complexity



$$\begin{array}{c}
 32 \\
 \diagdown \\
 4 \\
 \diagup \\
 128
 \end{array}$$



The problem we see is because of

Vanishing gradients

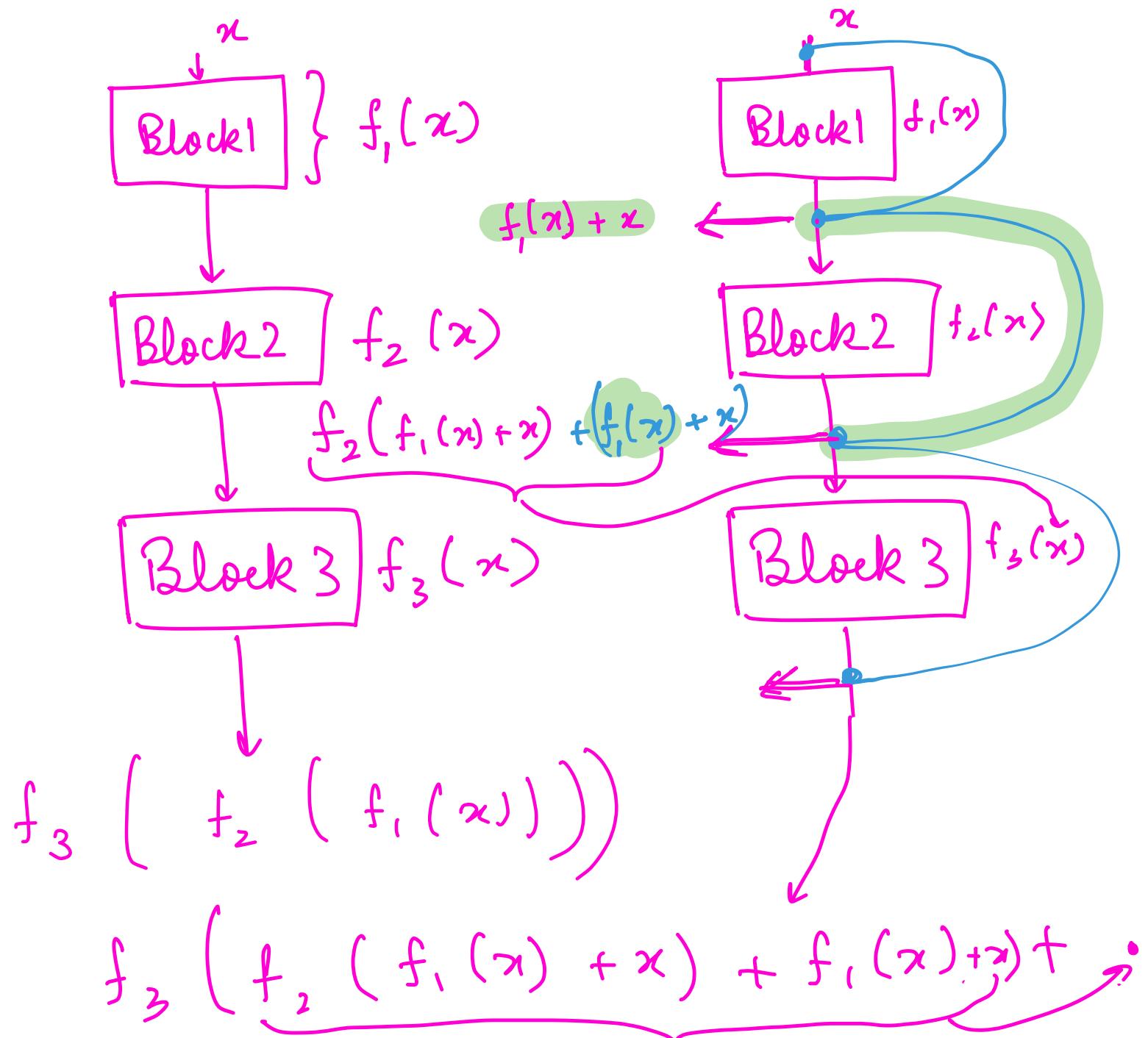
10

Layer 1's update =

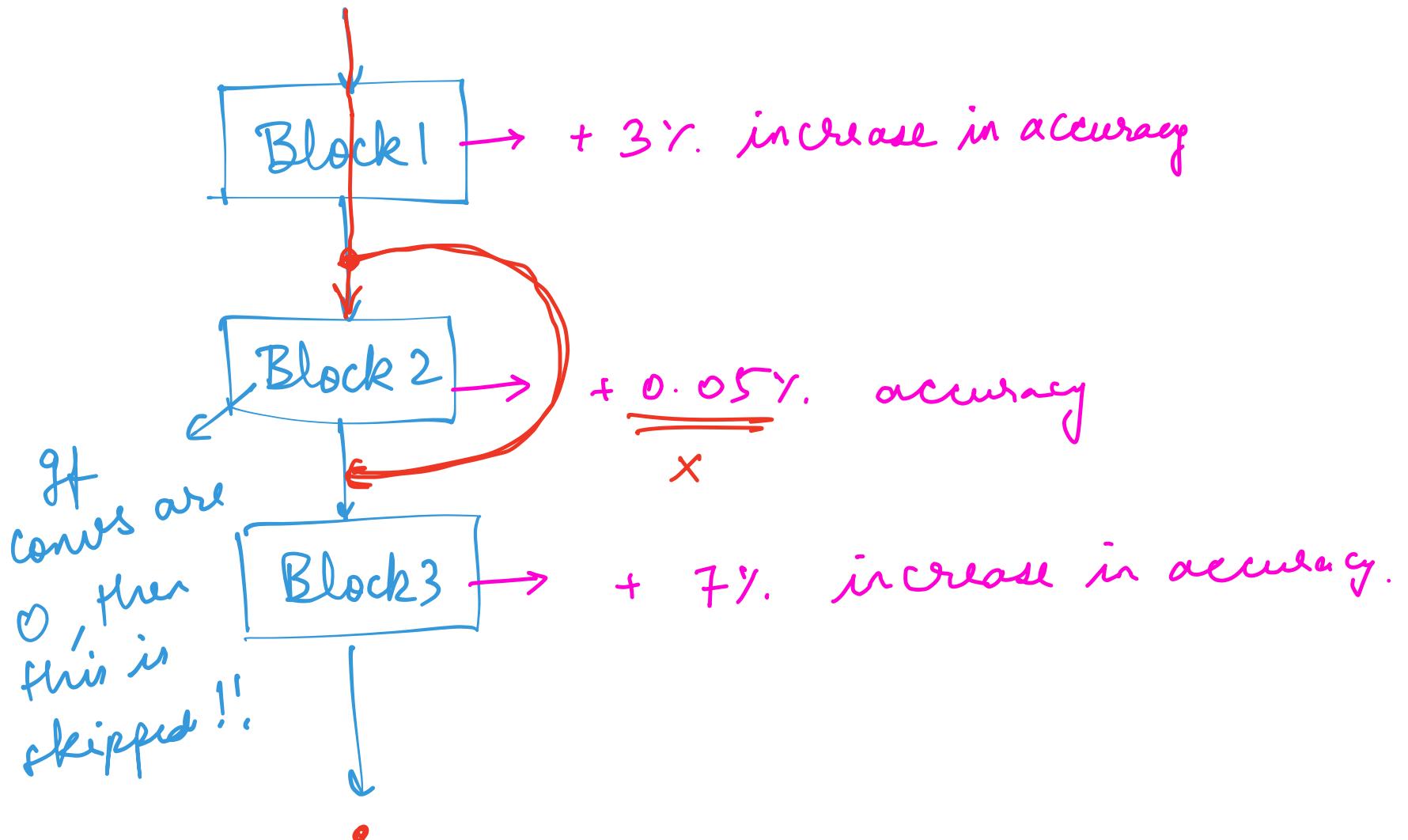


layer 10 =

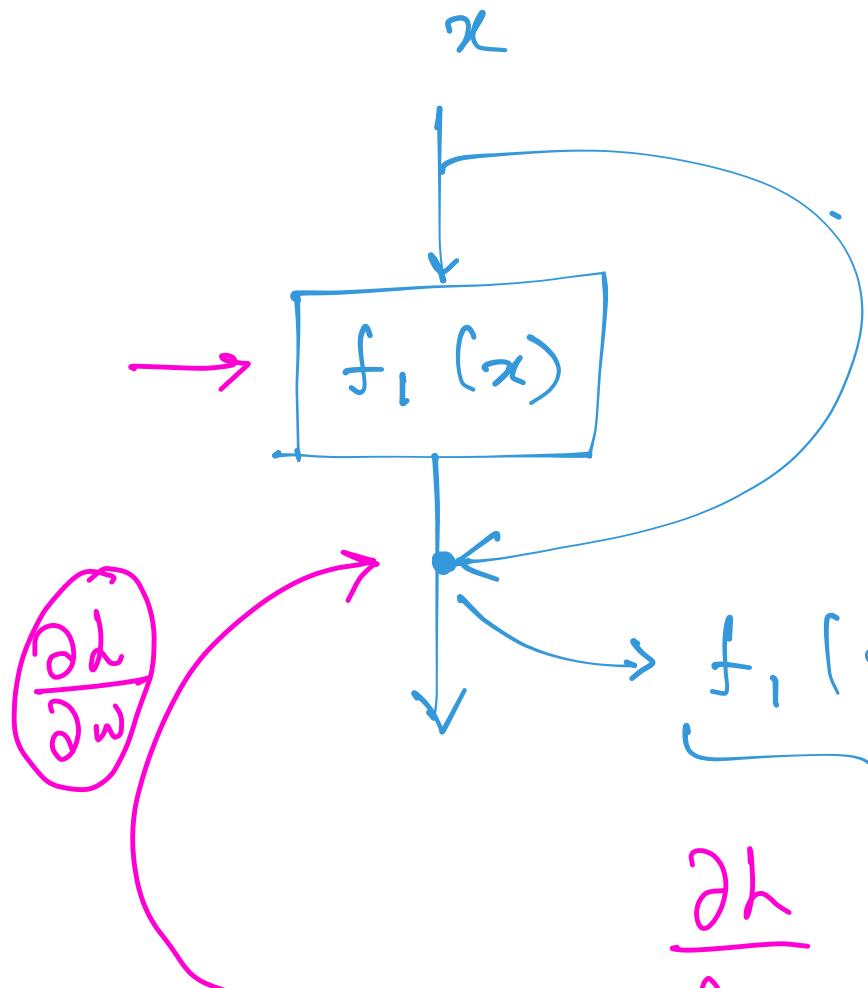




## Forward pass :



Backward pass:



$$h'(x) = f'_i(x) + \frac{1}{w}$$

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial w} \times \frac{\partial w}{\partial x}$$

Getting  
from backprop.

$$f_i(x) + \frac{1}{w} = h(x)$$

is an  
addition.

$$\frac{\partial L}{\partial x} = \frac{\partial d}{\partial w} \times \frac{\partial f_i}{\partial x} + \frac{\partial L}{\partial w}$$

*Note: The term  $\frac{\partial f_i}{\partial x}$  is crossed out with a large red X.*

## Takeaways :

① CNNs have many types of architectures.

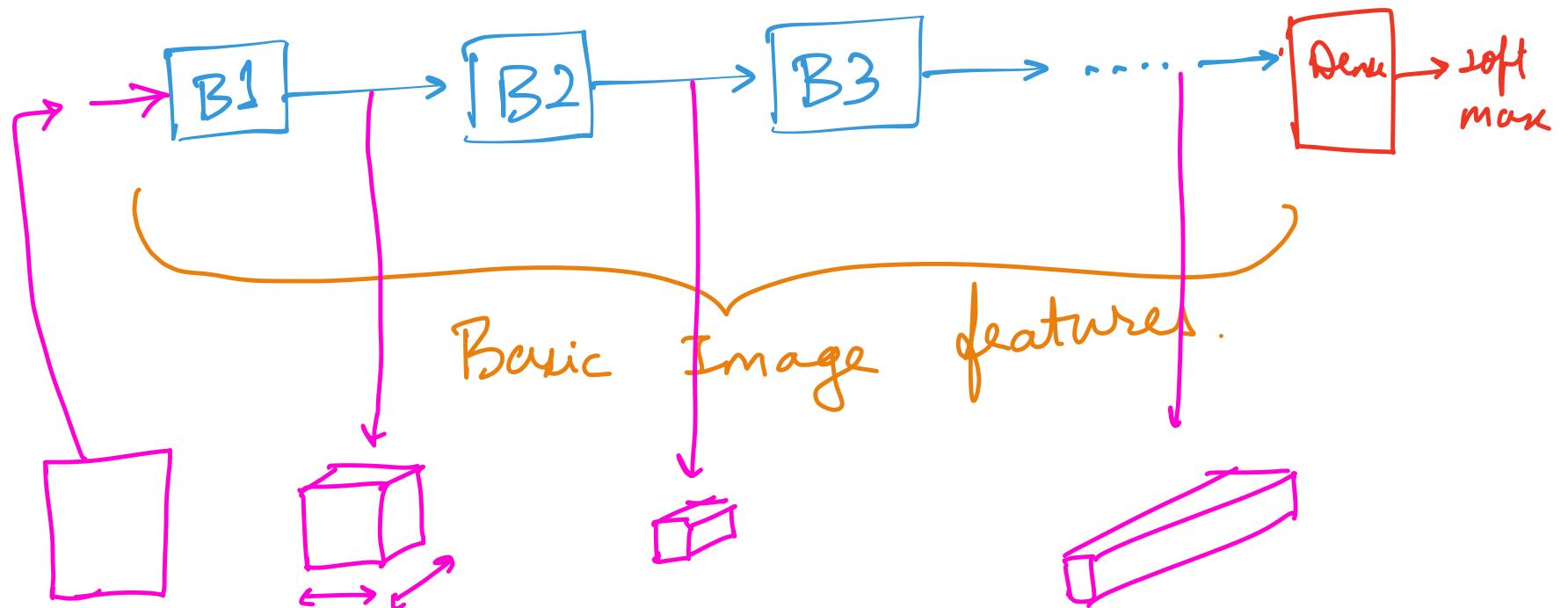
- Alexnet :
- VGG :
- Inception :
- Resnet :

↓  
increasing  
order  
of performance .

② The best architecture uses :

- \* Conv + ReLU + GPU.
- \* same kernel size ( $3 \times 3$  are okay).
- \*  $1 \times 1$  comes to reduce the no. of params .
- \* skip connections -

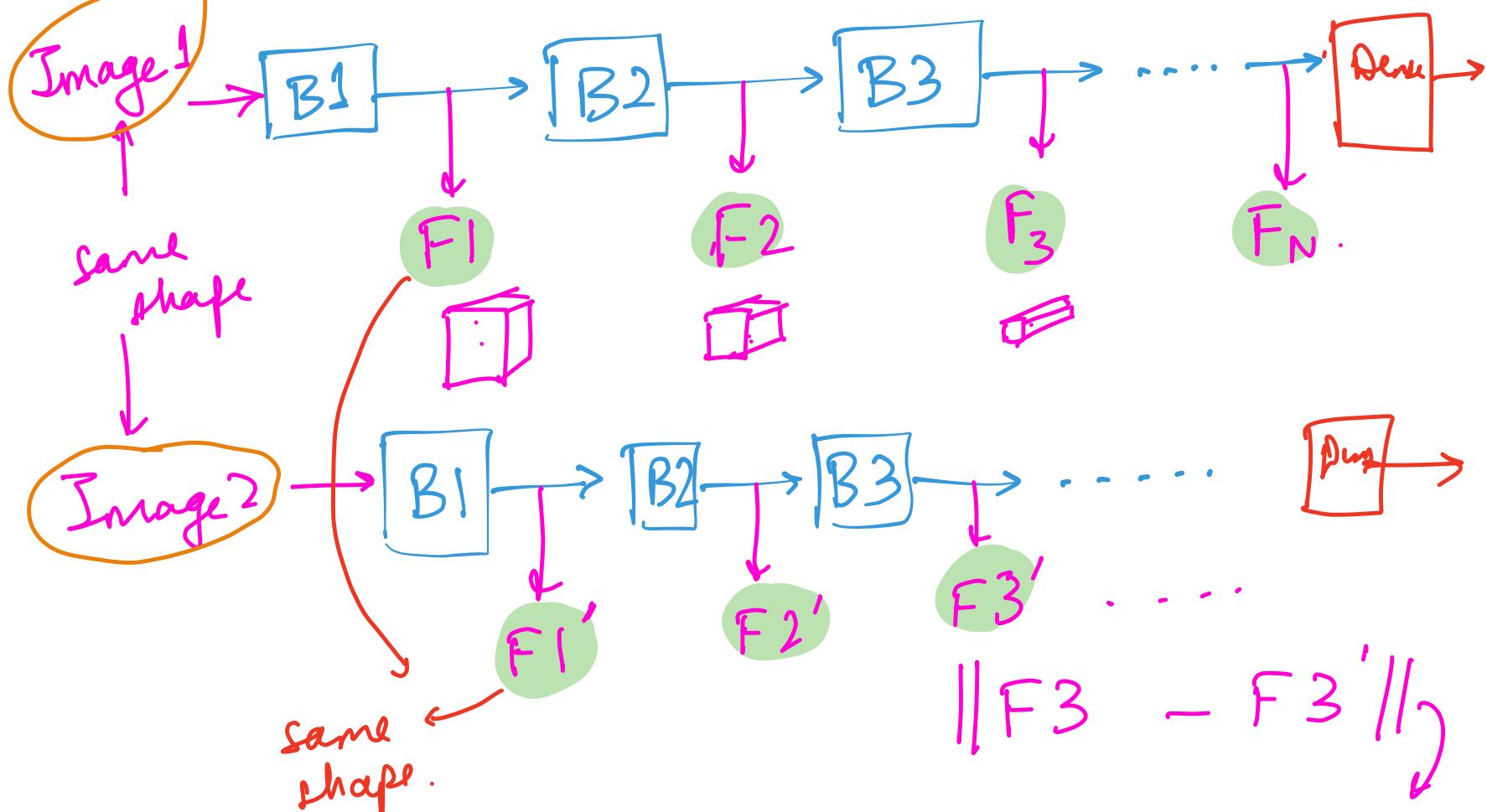
## Basic CNN architecture:



How to check if 2 images are similar?

As you go deeper,

the intermediate representations of similar images becomes closer.



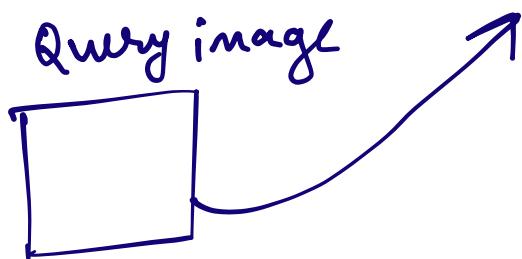
$$\|\text{Image}^1 - \text{Image}^2\| \rightarrow \text{very high}$$

much lower.

## Image similarity problem:

(a) Start with an already trained network.

→ ResNet → F 37 from ResNet50.



F 37 picked after some experimentation.

(b) Compute the F 37 for all images in the dataset.

(c) Return the images which are closest to the query image in terms of F 37.

Embedding from a CNN



Take one of the  
features, and flatten it!

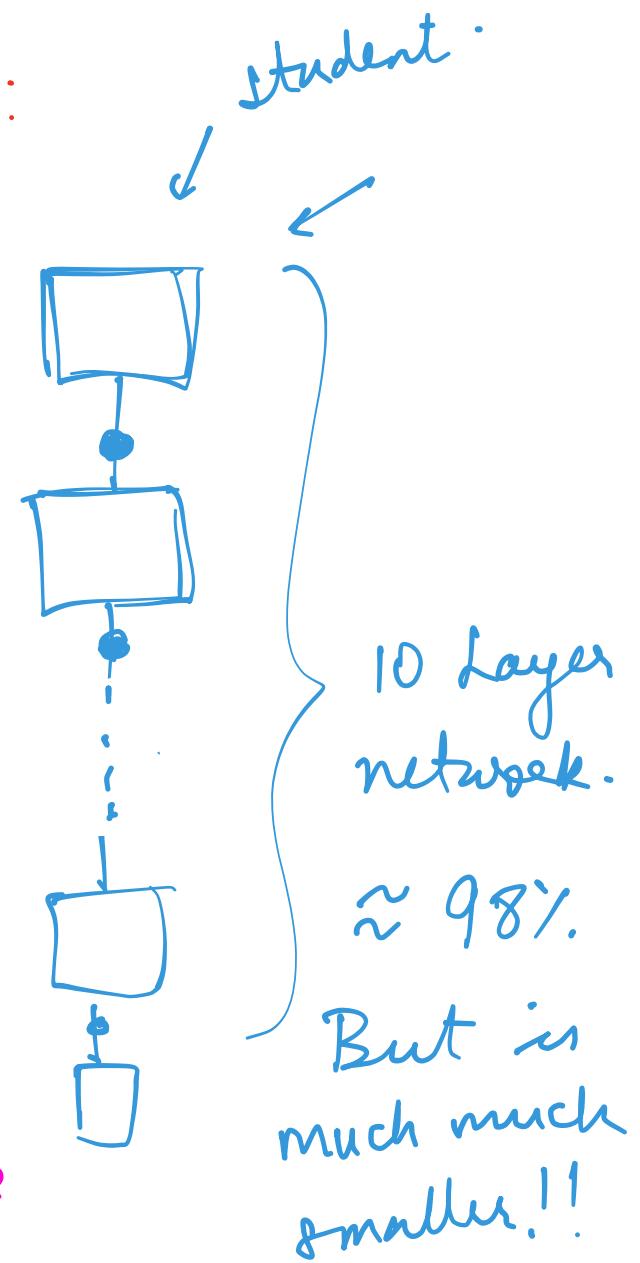
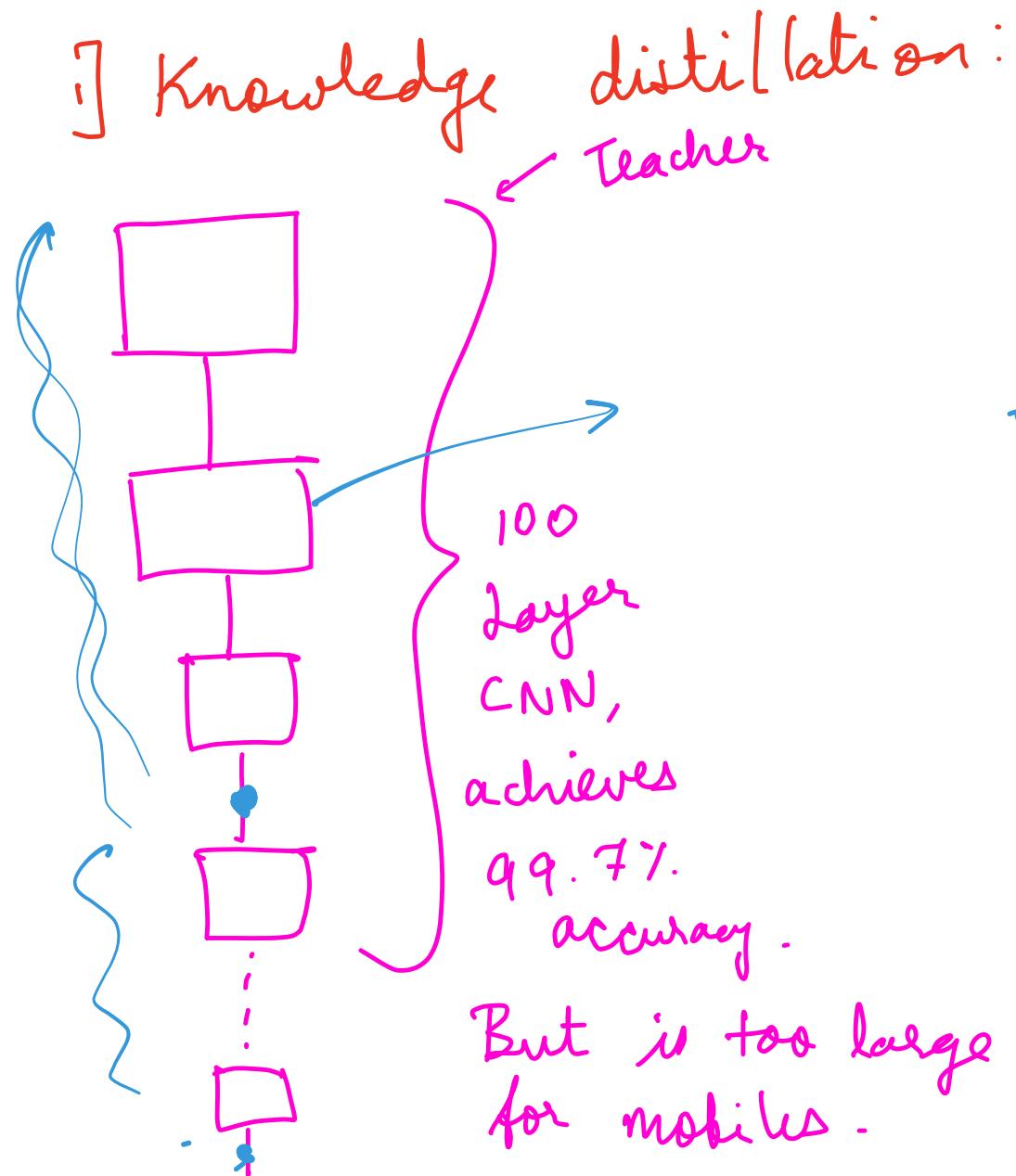
Time for 1 img  $\rightarrow$  99.8ms  
 $\approx 100\text{ ms} = 0.1\text{ s}$ .

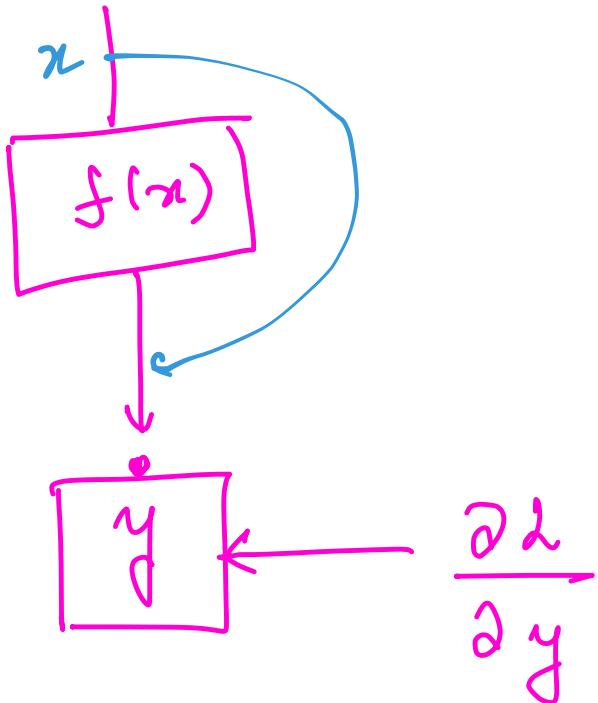
10,000,000  $\rightarrow$  Can get done pretty fast  
if we have multiple GPUs.

Once we have all f.v.s, apply KNN  
↓  
to return results!!

K-Nearest  
Neighbours.

## Doubts :





$$y = f(x) + \epsilon$$

$$\frac{\partial y}{\partial x} = \frac{\partial f}{\partial x} + 1$$

$$\begin{aligned}\frac{\partial L}{\partial x} &= \frac{\partial L}{\partial y} \cdot \frac{\partial y}{\partial x} \\ \frac{\partial L}{\partial x} &= \frac{\partial L}{\partial y} \left( \frac{\partial f}{\partial x} + 1 \right) \\ &= \underbrace{\frac{\partial L}{\partial y}}_{\text{even if this goes to 0,}} \cdot \frac{\partial f}{\partial x} + \frac{\partial L}{\partial y}.\end{aligned}$$

↑

this will update my weights.

10,000,000  
 $(227 \times 227 \times 3)$

$F_{50}$  (of all).

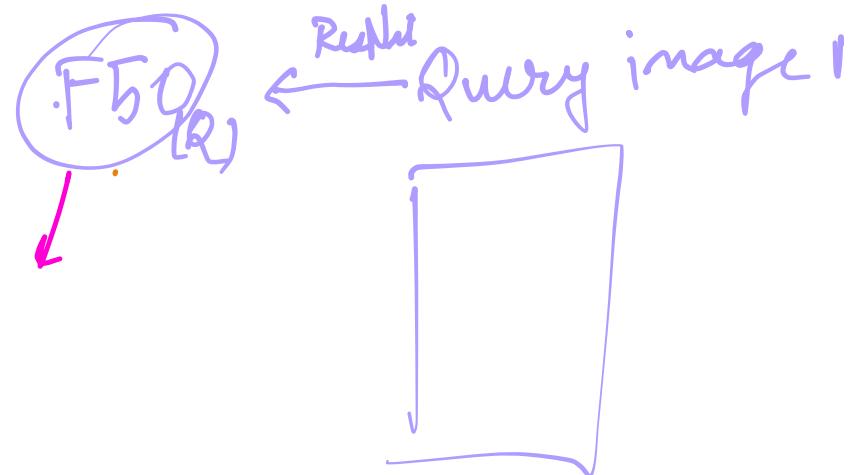
(2048).

Diagram illustrating the storage of image features:

A table with two columns: "Image path" and  $F_{50}$ . The "Image path" column contains four rows, each represented by a blue bracket. The  $F_{50}$  column contains four rows, each represented by a red bracket. Orange arrows point from the "Image path" rows to the corresponding  $F_{50}$  rows.

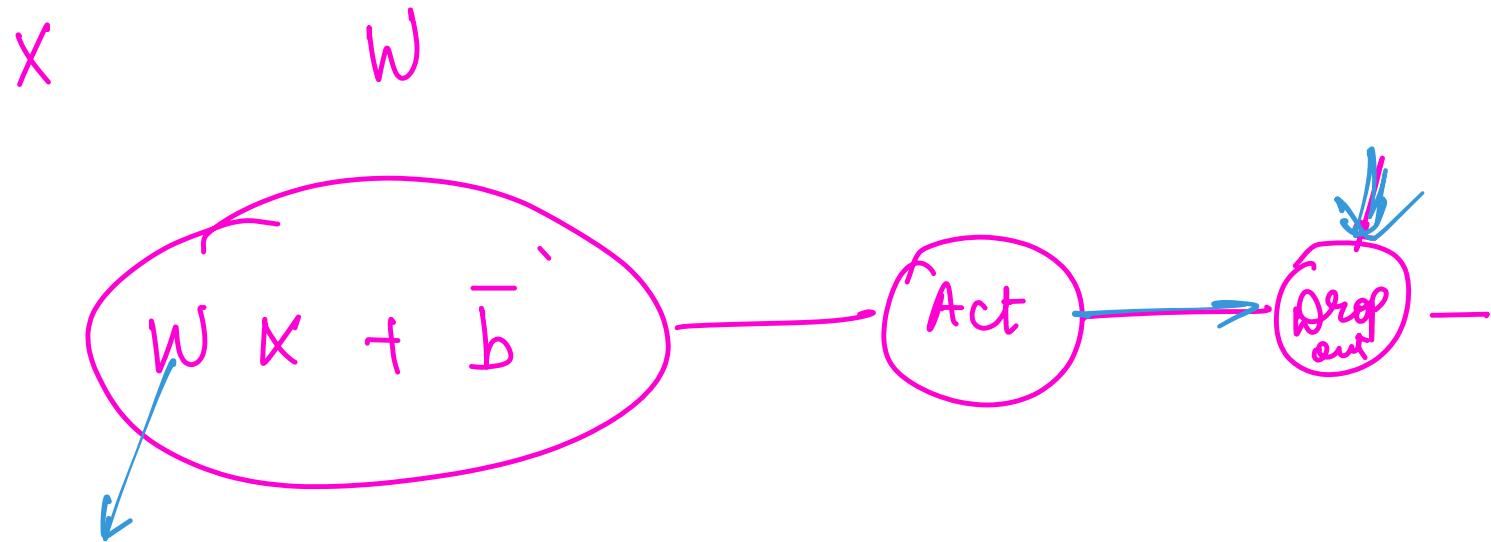
Image path	$F_{50}$

dSH  $\rightarrow$  Locality Sensitive Hashing.



$(F_{50})_R \leftarrow$  Query image 2

$\underbrace{10,000,000}_{}$  calculations for each query point.



W

$$\begin{bmatrix} 0 & 0 & 1 & 0 \\ 0 & 1 & 1 & 1 \\ 1 & 0 & 0 & 1 & 0 \end{bmatrix}$$

