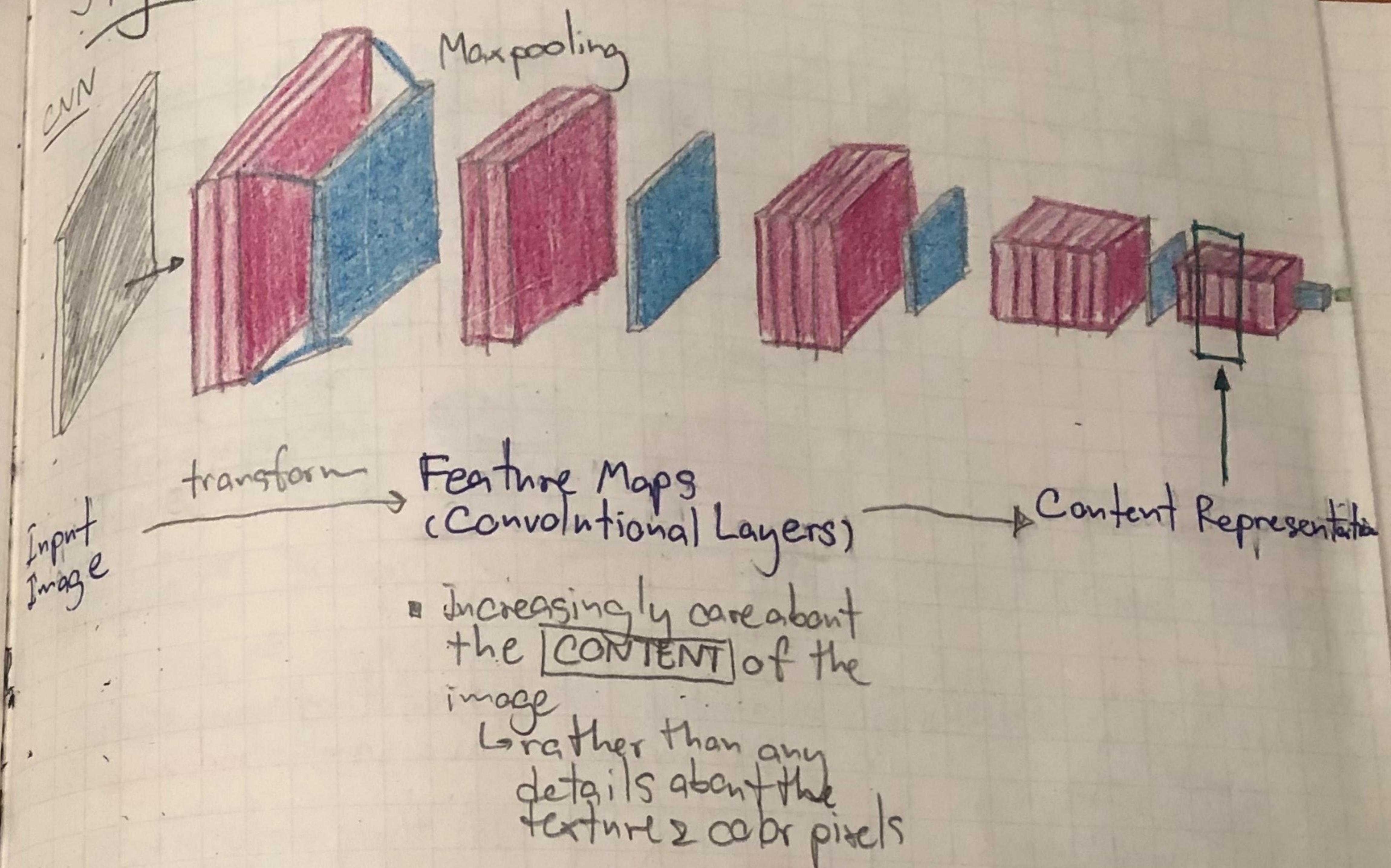


style Transfer



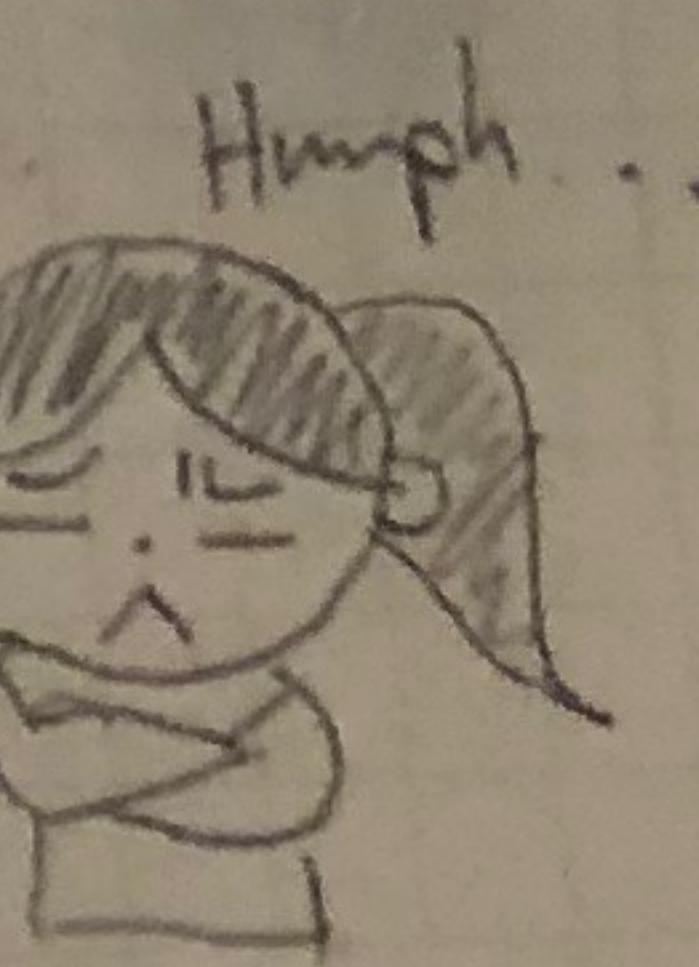
■ A trained CNN has already learned to represent the content of an image.

- How about the style.
- Style transfer

CONTENT
of
one image

STYLE
of
Another image

How can we isolate
only the style of
an image?

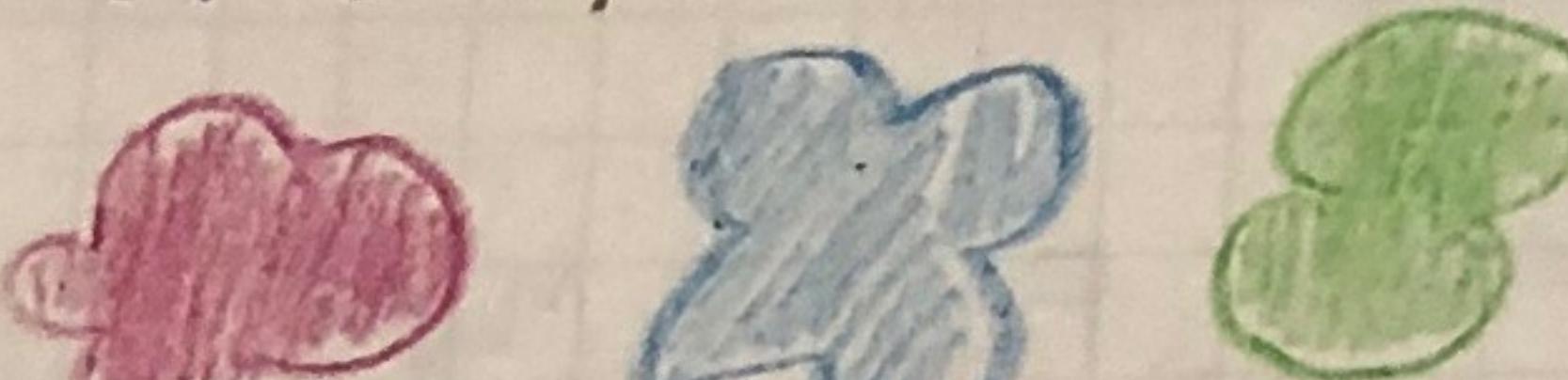


To represent the STYLE of an input image,

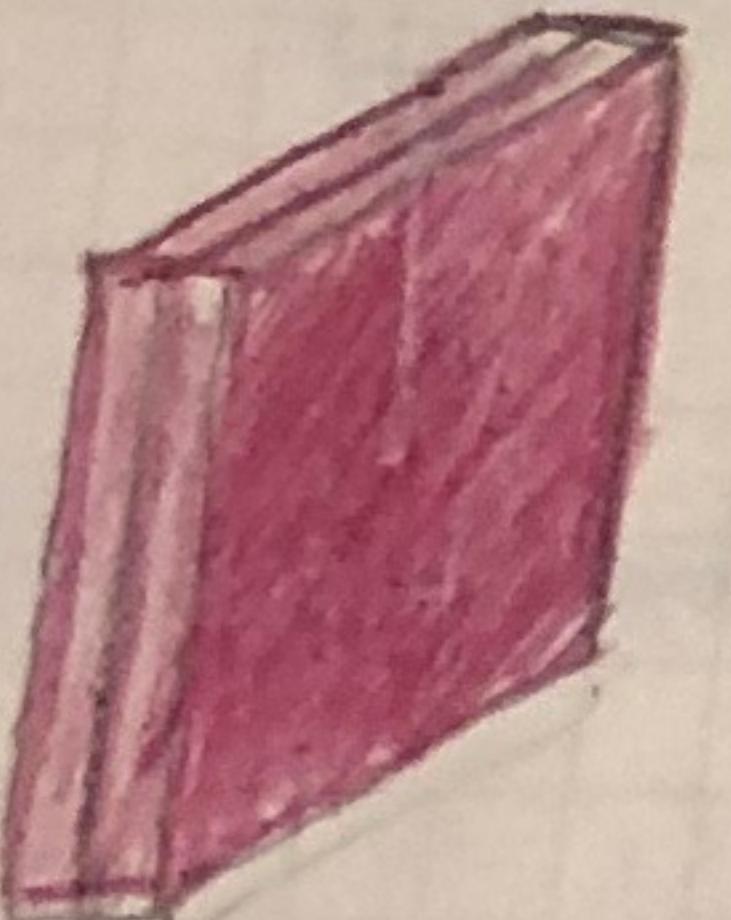
- Feature Space designed to capture texture and color information is used

↳ This space looks at spatial correlations within a layer of network

- How do the features in this layer relate to one another?



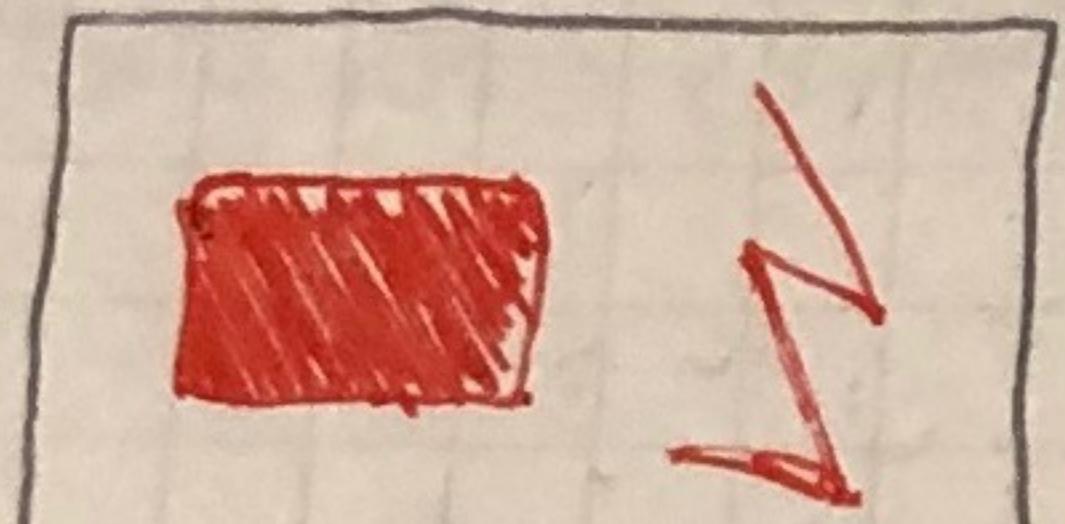
See which colors & shapes in a layer of feature maps are related & which are not



depth = 64

conv1
depth corresponding to feature maps in that layer

- For example, we detect that mini-feature maps in 1st conv layer have similar red edge features

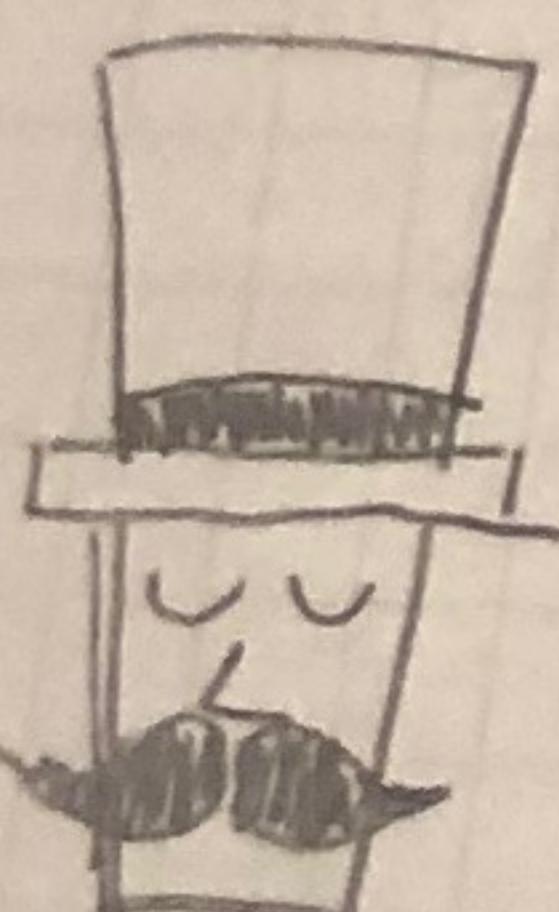


If there are COMMON colors and shapes among the feature maps

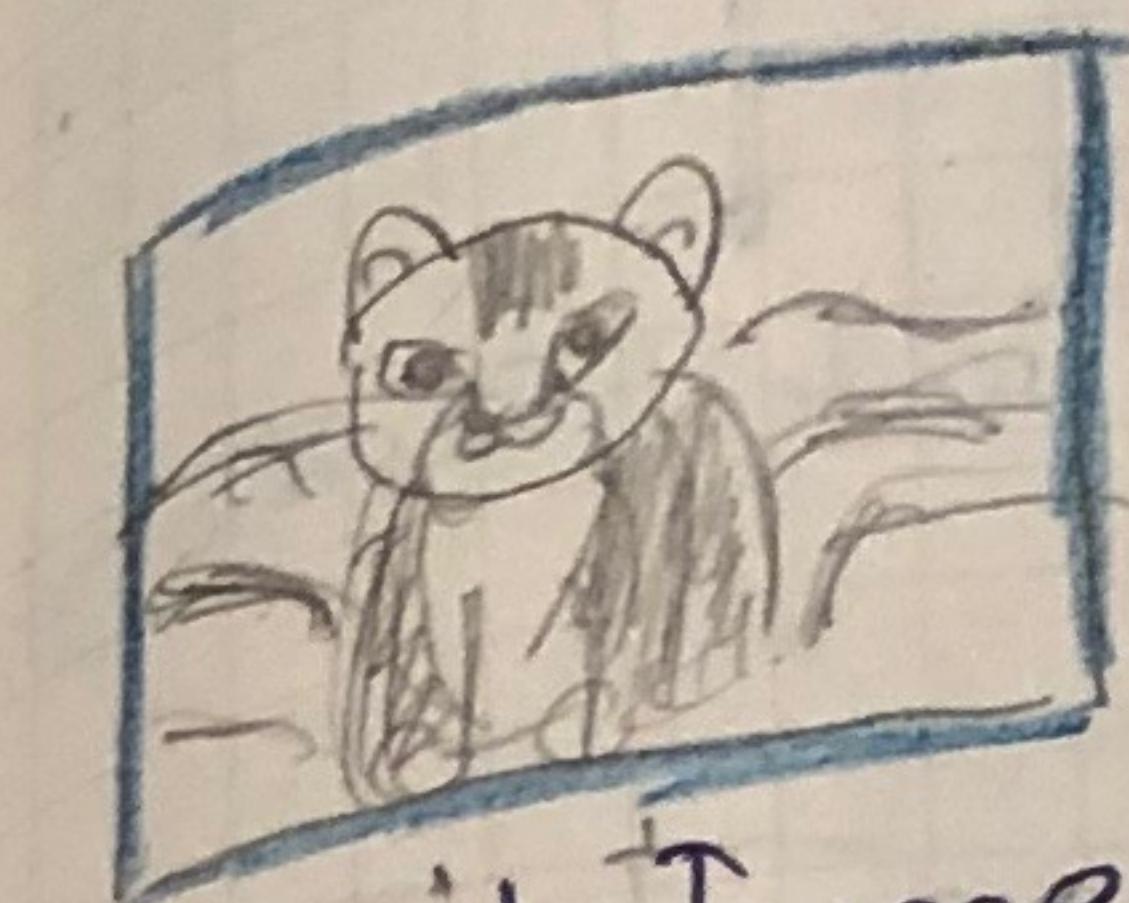
- Then, this can be thought of as part of that image style

- So the Similarities & Differences between features in a layer should give us some information about the texture and color information found in an image

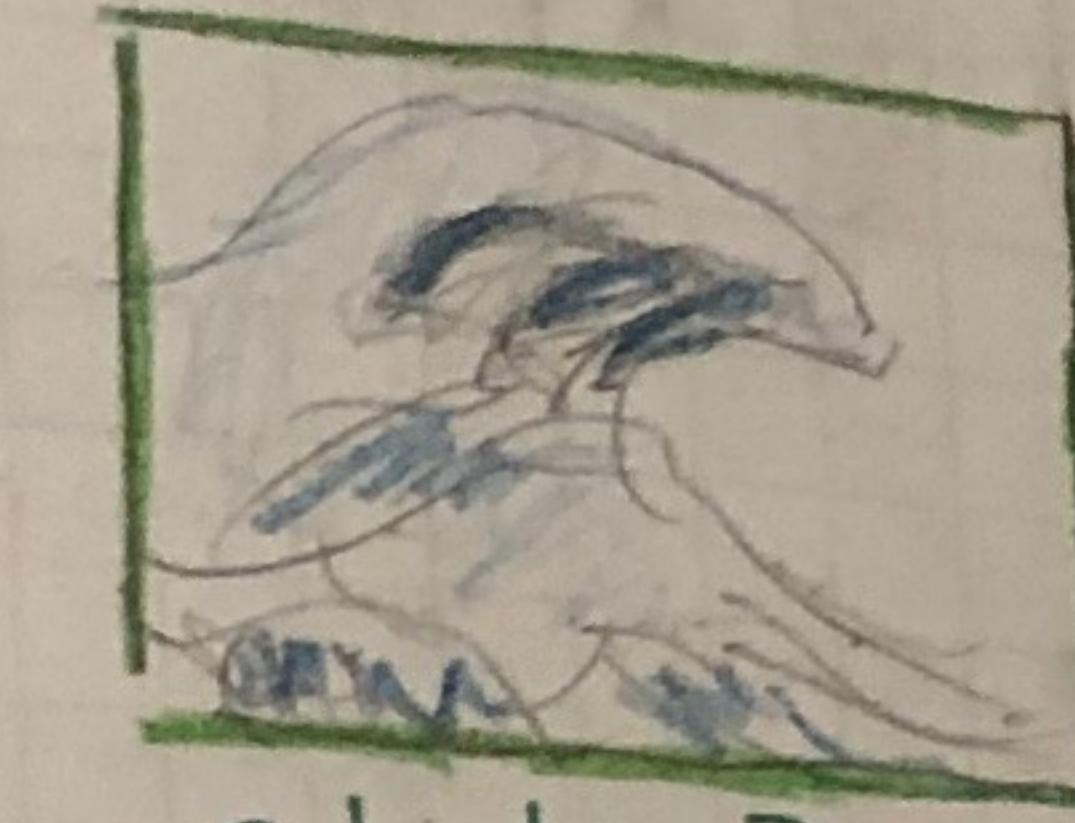
- It SHOULD NOT tell us anything about the identity or placement of objects in an image



Style Transfer



Content Image



Style Image



Target Image

object & shape arrangement

Colors & textures

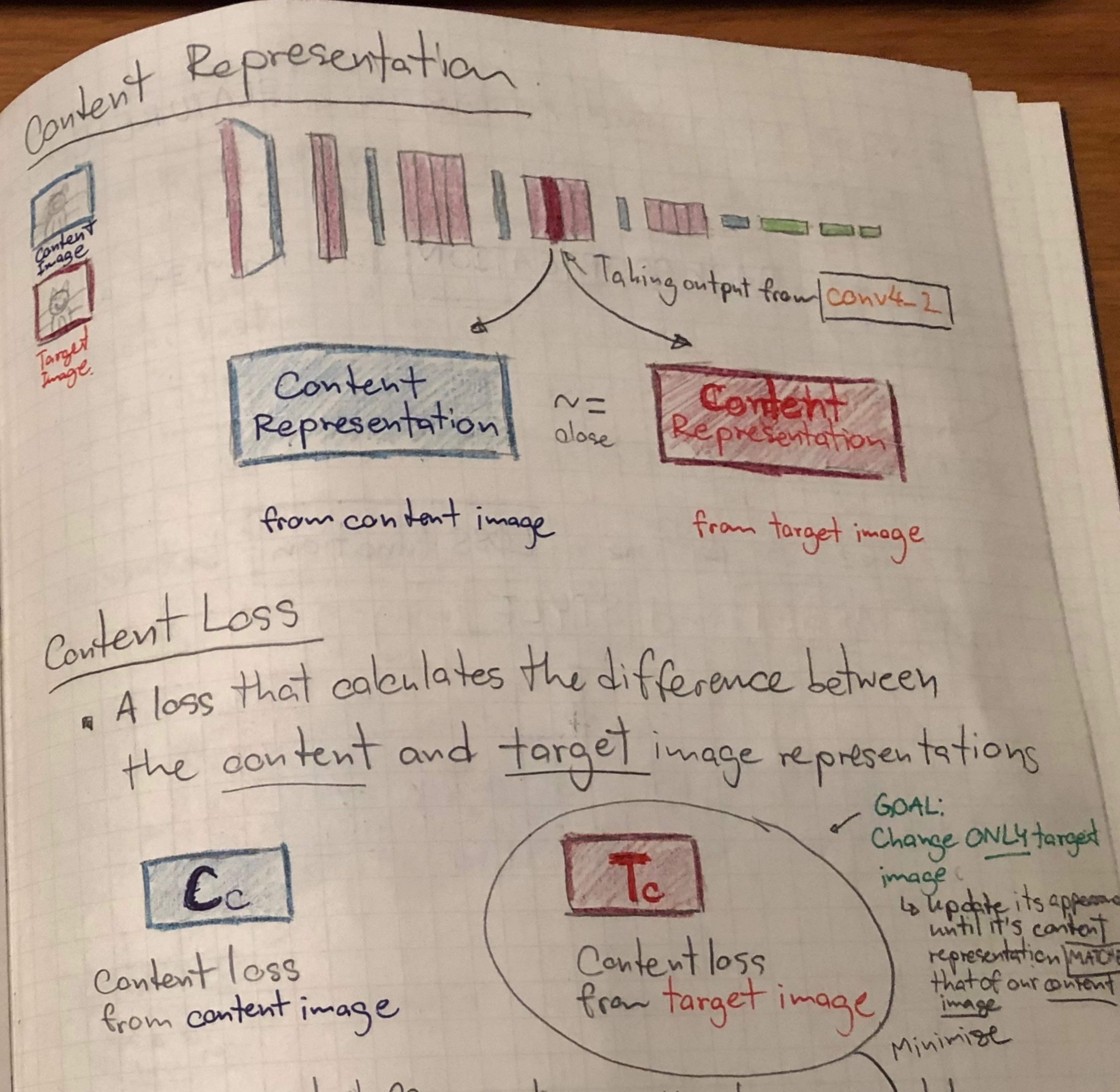
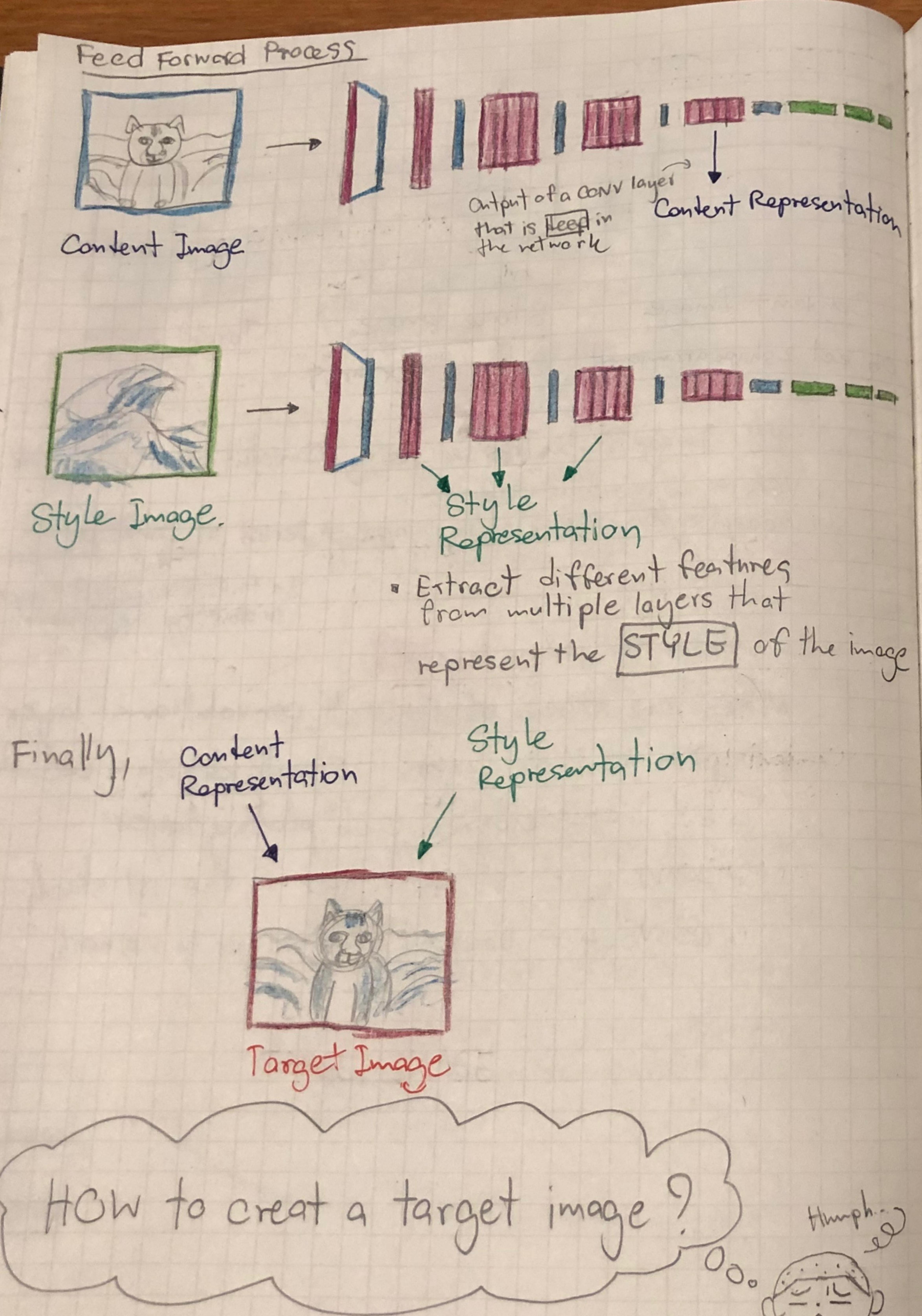
From paper: Image Style Transfer Using Convolutional Neural Networks

↳ VGG 19 network is used

↳ Accepts a $224 \times 224 \times 3$ color image \rightarrow Series of conv & pooling layers

3x Fully connected layers to classify the passed image

- ↳ In between the pooling layers,
 - there are stacks of 2 or 4 convolutional layers
 - Depth of these layers is standard within each stack
 - but increase after each pooling layer
 - conv1-1 \rightarrow 1st conv layer in the 1st stack
 - conv5-4 \rightarrow deepest conv layer in network



$$L_{\text{content}} = \frac{1}{2} \sum (T_c - C_c)^2$$

- Measures how far away these two representations are from one another.
 - To create the BEST target image

GOAL: Minimize Content Logs (Lcontent)

- Using a Pretrained network (VGG19) as a **FEATURE Extractor**
 - different from a traditional way → NOT training it to produce a specific output
- Using **BACK PROPAGATION** to minimize a defined loss function between our target & content images.

For style 

- We need to define a **LOSS function** between our **TARGET** and **STYLE** images



Produce an image w/
our **DESIRED** style.



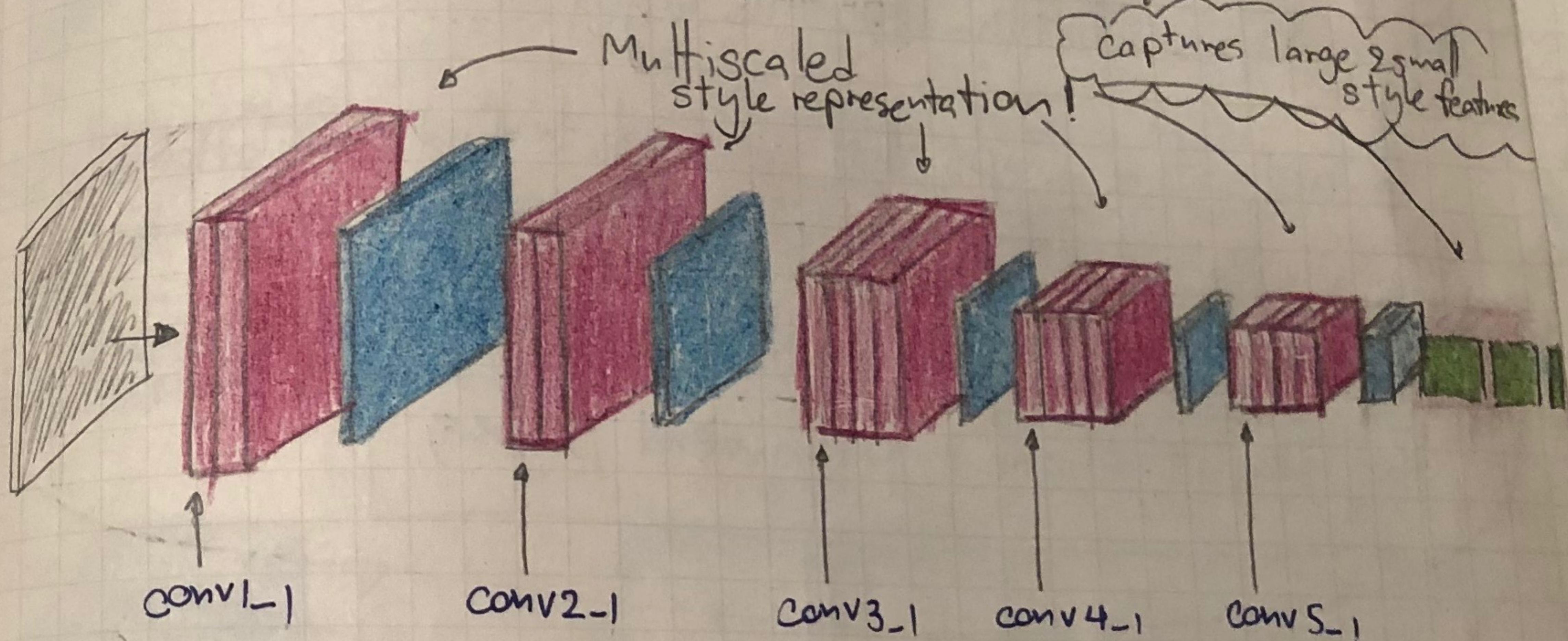
Style Representation

- Relyes on looking at **CORRELATIONS** between the features in individual layers of VGG19
 - looking at how similar the features in a **single** layer are

SIMILARITIES

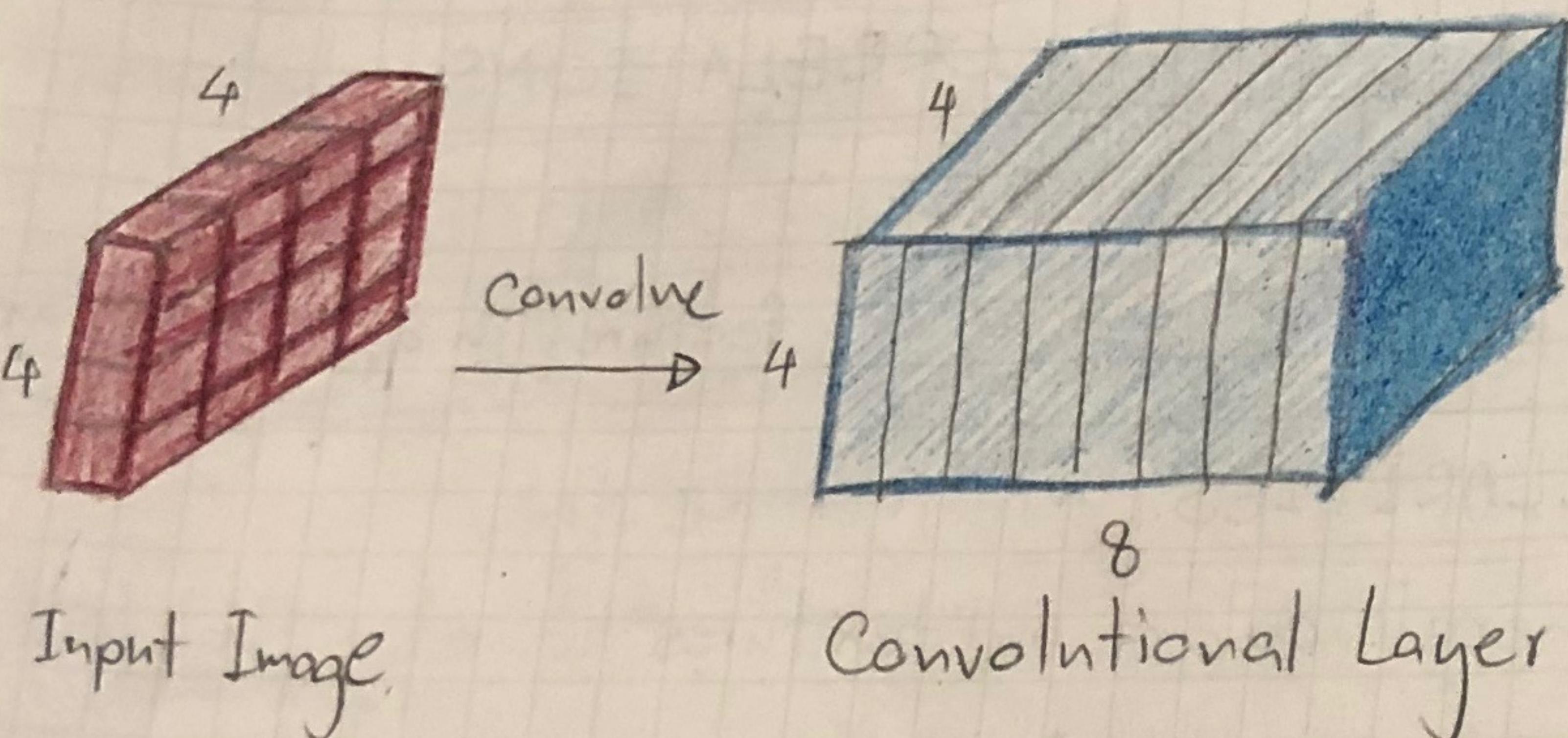
will include the general colors and textures found in that layer

- we typically find the similarities between features in multiple layers in the network



- Style representation is calculated as an image passes through the network at 1st conv layer in all five stacks
- Correlations** at each layer are given by a **GRAM MATRIX**

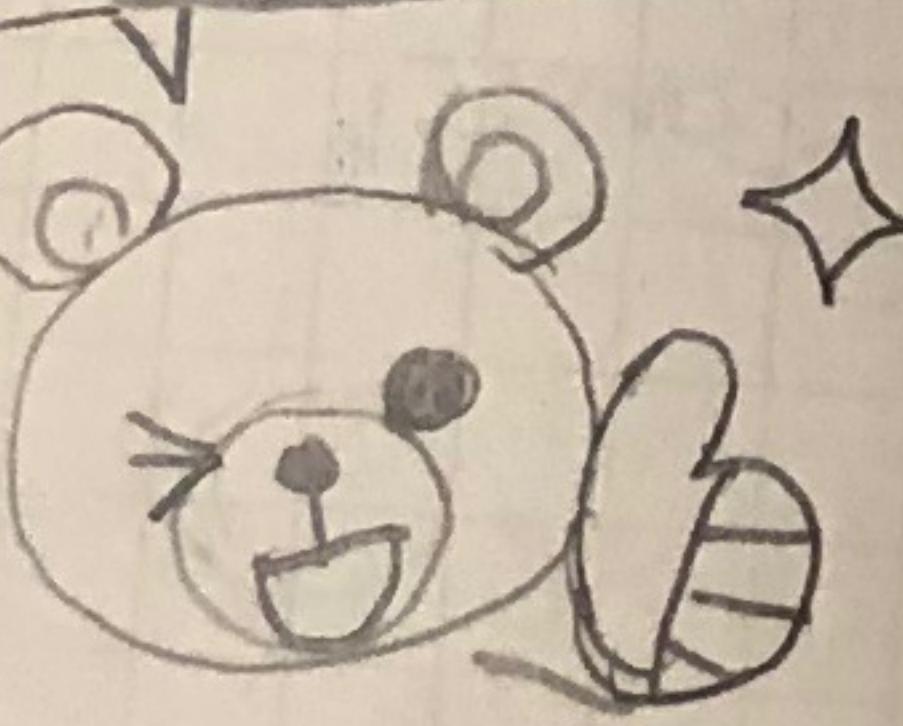
GRAM MATRIX



Style Representation for this layer

↳ 8 Feature maps that we want to find the relationships between

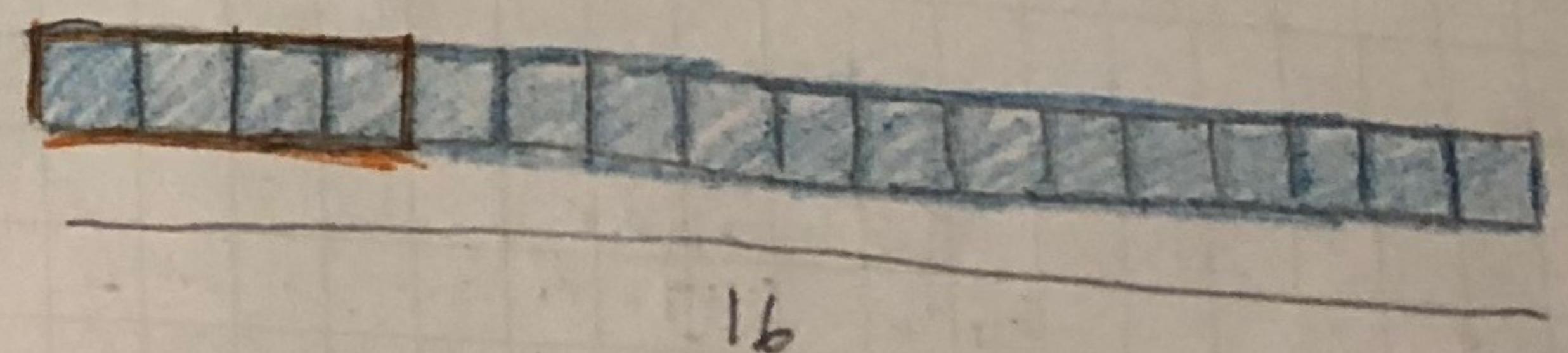
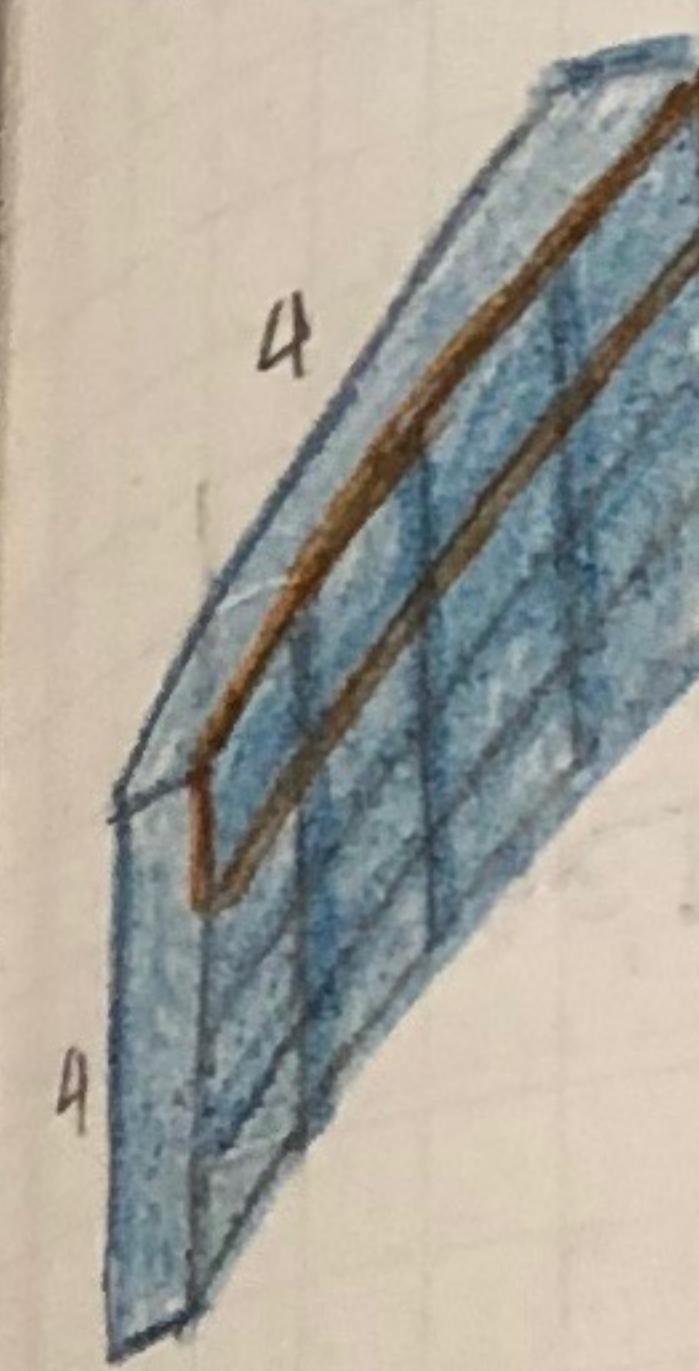
- GRAM Matrix is a mathematical way of representing the idea of share in prominent styles.



STEPS to calculate the Gram Matrix

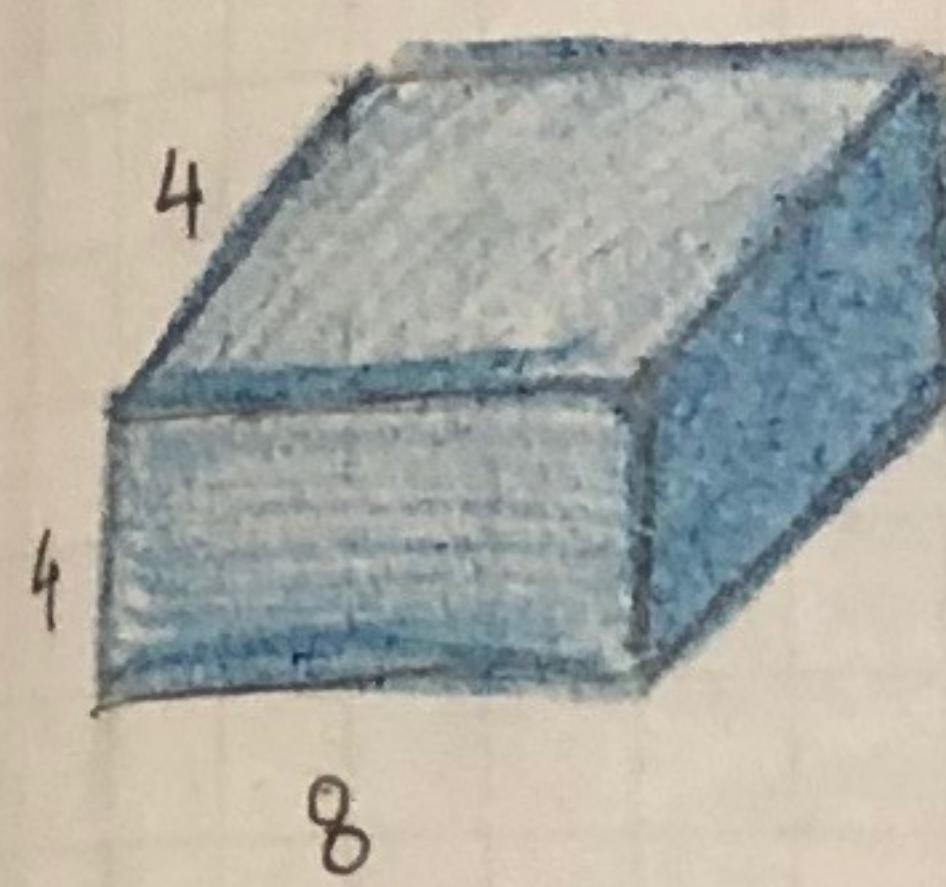
① Vectorize the values in this layer

For One Feature Map



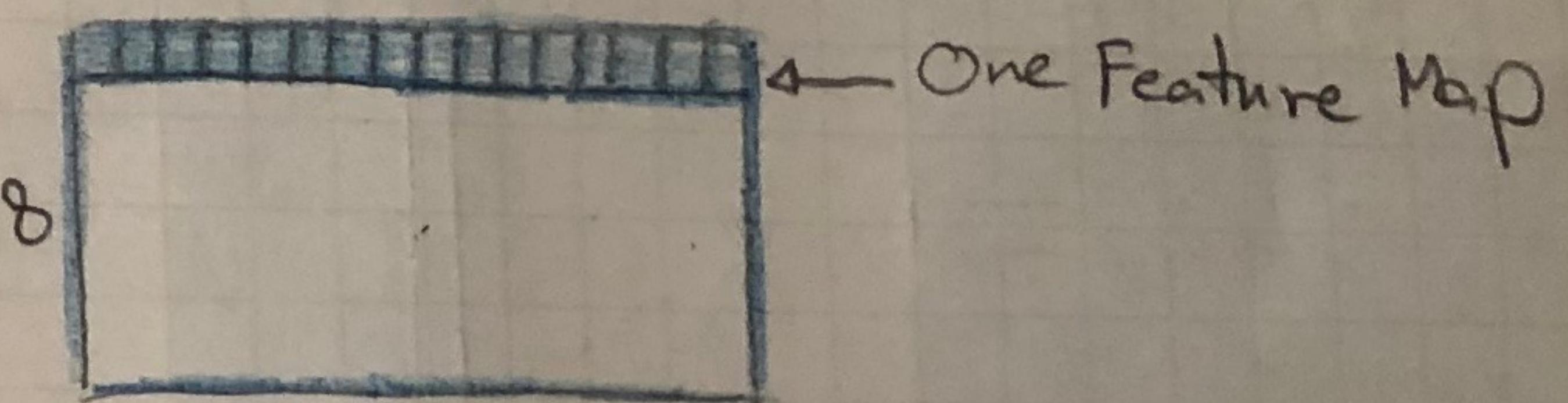
For N Feature Map

Converting 3D conv layer into a 2D matrix of values



8 Feature Maps

$$4 \times 4 = 16$$



Vectorized Feature Maps

② Multiply this matrix (Vectorized Feature Maps) by its transpose

↳ Essentially, we multiply the features in each map to get the GRAM matrix

↳ This operation treats each value in the feature map as an individual sample, unrelated in space to other values

⇒ Resultant Gram Matrix contains **NON-LOCALIZED** Information about the layer

Information that would still be there even if an image was shuffled around in space

⇒ Squared GRAM Matrix whose **VALUES** indicate the **SIMILARITIES** between the layers

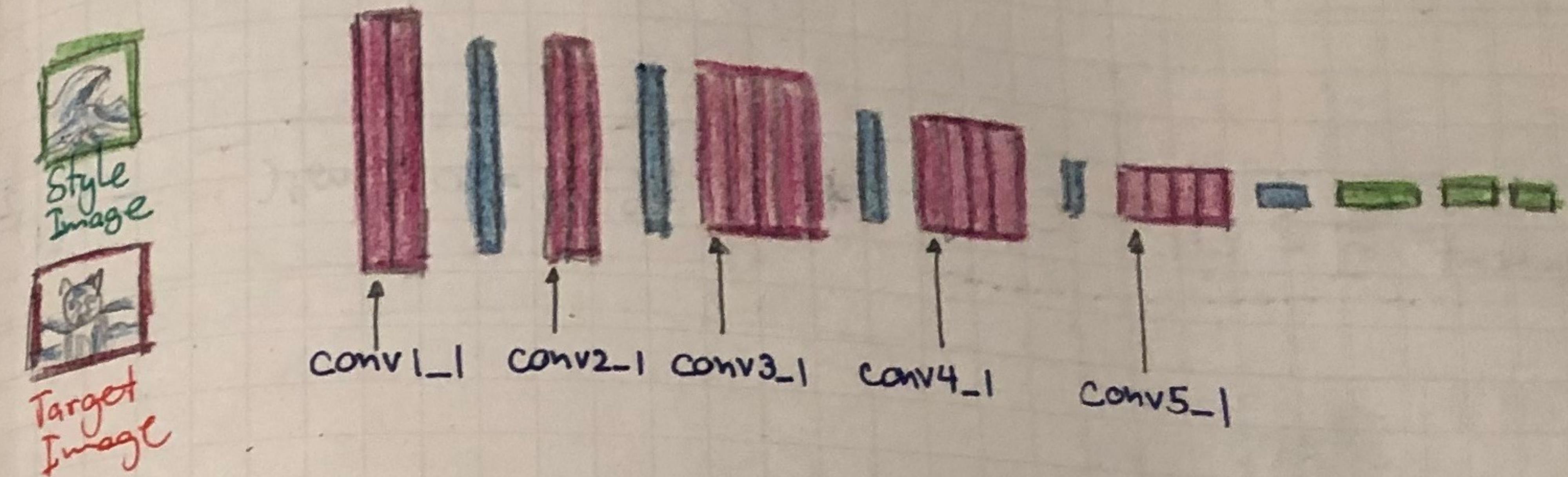
$$\begin{matrix} \begin{matrix} 4 \times 4 = 16 \\ 8 \end{matrix} & \times & \begin{matrix} 16 \\ 8 \end{matrix} \\ \begin{matrix} \text{Transpose} \end{matrix} & & \end{matrix} = \begin{matrix} \begin{matrix} 2 \\ 4 \end{matrix} \rightarrow \boxed{1} \\ \begin{matrix} 2 \\ 4 \end{matrix} \rightarrow \boxed{2} \\ \dots \\ \dots \end{matrix} \quad \begin{matrix} 8 \\ 8 \end{matrix}$$

Gram Matrix, G

↳ $G(4,2) \rightarrow$ Indicates the **similarity** between the 4th and 2nd feature maps in a layer

↳ Dimension will only relate to * feature map in this CONV layer

Style Loss



List of Gram Matrices

from style image

List of Gram Matrices

from target image

To compute style loss

▪ Find MEAN SQUARED Distance between the **style** and **target** Gram matrices

↳ All 5 pairs that are computed at each layer in our predefined list conv1_1 up to conv5_1

S_s
Style loss
from Style image

The only value that changes!
 T_s
Style loss
from Target image.

Style Loss

$$L_{\text{style}} = a \sum_i w_i (T_{s,i} - S_{s,i})^2$$

where a : constant that accounts for * values in each layer
 w : style weight → give more or less weight to the calculated style loss at each layer

Total Loss

$$L_{\text{content}} = \frac{1}{2} \sum (T_c - C_c)^2$$

$$L_{\text{style}} = \alpha \sum w_i (T_{s,i} - S_{s,i})^2$$

Steps:

- Compute TOTAL LOSS
- Then, use typical back propagation and optimization to reduce this loss
 - by iteratively changing the TARGET image to match our desired CONTENT and STYLE



We have values for the Content & style loss but, No!!!, because they're calculated differently, These values will be pretty different



We want our TARGET IMAGE to take both into account fairly equally (total loss reflects an equal balance) OK!

$$\alpha L_{\text{content}} + \beta L_{\text{style}}$$

Content weight Style weight

IN PRACTICE \Rightarrow

Often Much Larger!

$$\frac{\alpha}{\beta}$$



Bigger or smaller ratio have an impact to the target image

The more smaller $\frac{\alpha}{\beta}$, the more stylistic effect you'll see

Putting it all together.

- Extract features from the "content" image
- Calculate the "gram" matrices for each layer in our style presentation from the "STYLE" image
- Create a third "target" image
 - It is a good idea to start off w/ the target as a copy of our "content" image

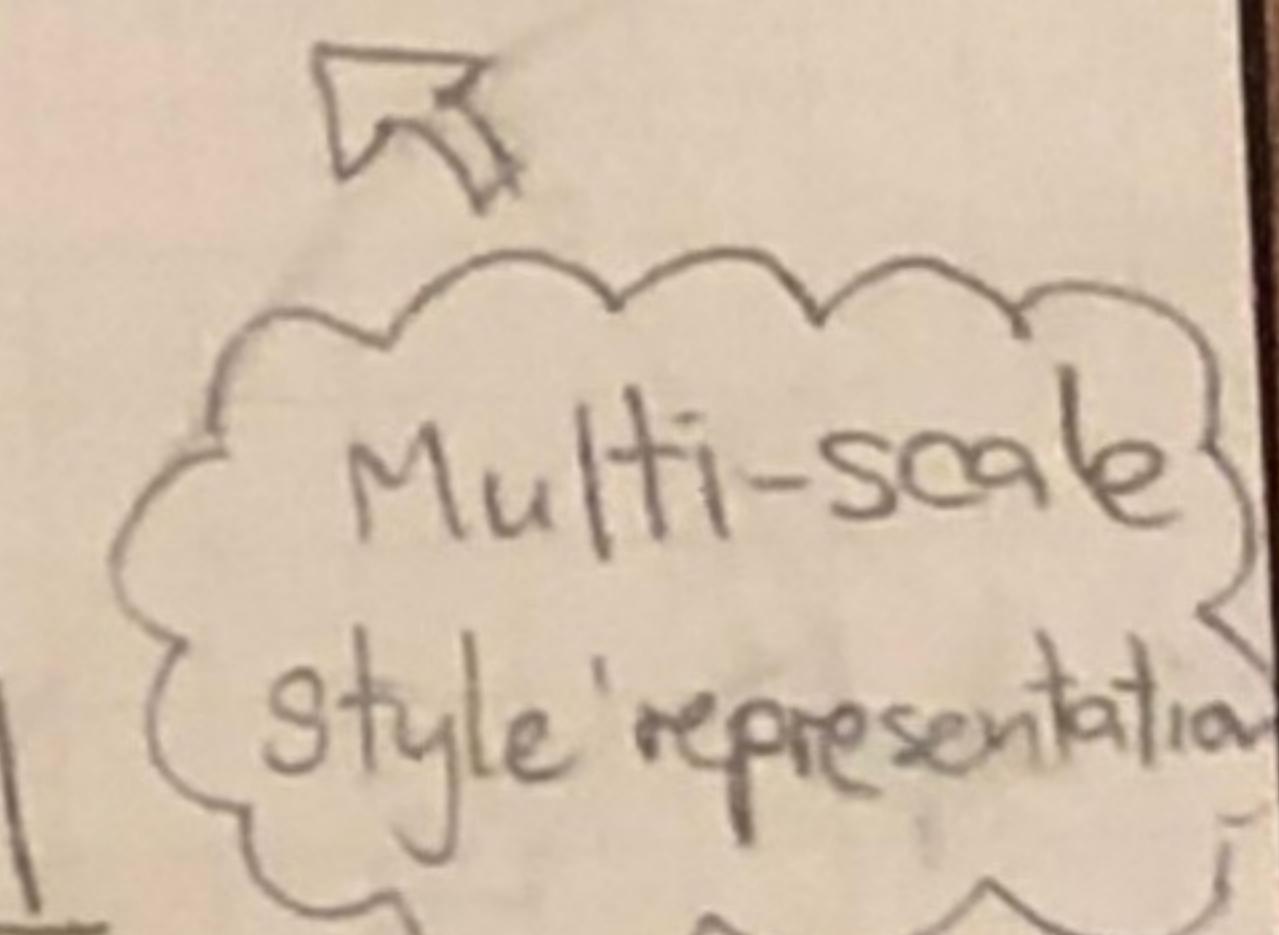
Then, iteratively change its style

Loss and Weights

Individual Layer STYLE Weights

- Suggested that you use a range between 0-1 to weight these layers
- Weight EARLIER layers (conv1-1 and conv2-1) MORE

Get LARGER style artifacts



- If choose to weight LATER layers MORE

MORE Emphasis on SMALLER feature

Content (α) and Style (β) Weight

- The ratio $\frac{\alpha}{\beta}$ will affect how Stylized your final image is
- Recommended to leave $\alpha = 1$ (content weight) set the style weight β to achieve the ratio you want