

Neural Style Transfer

Content Image

(Benedict Cum-purr-scratch)



C

+

Style Image

(Georgia O'Keeffe)



S

=



X

Neural Style Transfer

Content Image

(Benedict Cum-purr-scratch)



C

Style Image

(Georgia O'Keeffe)



S

+

=



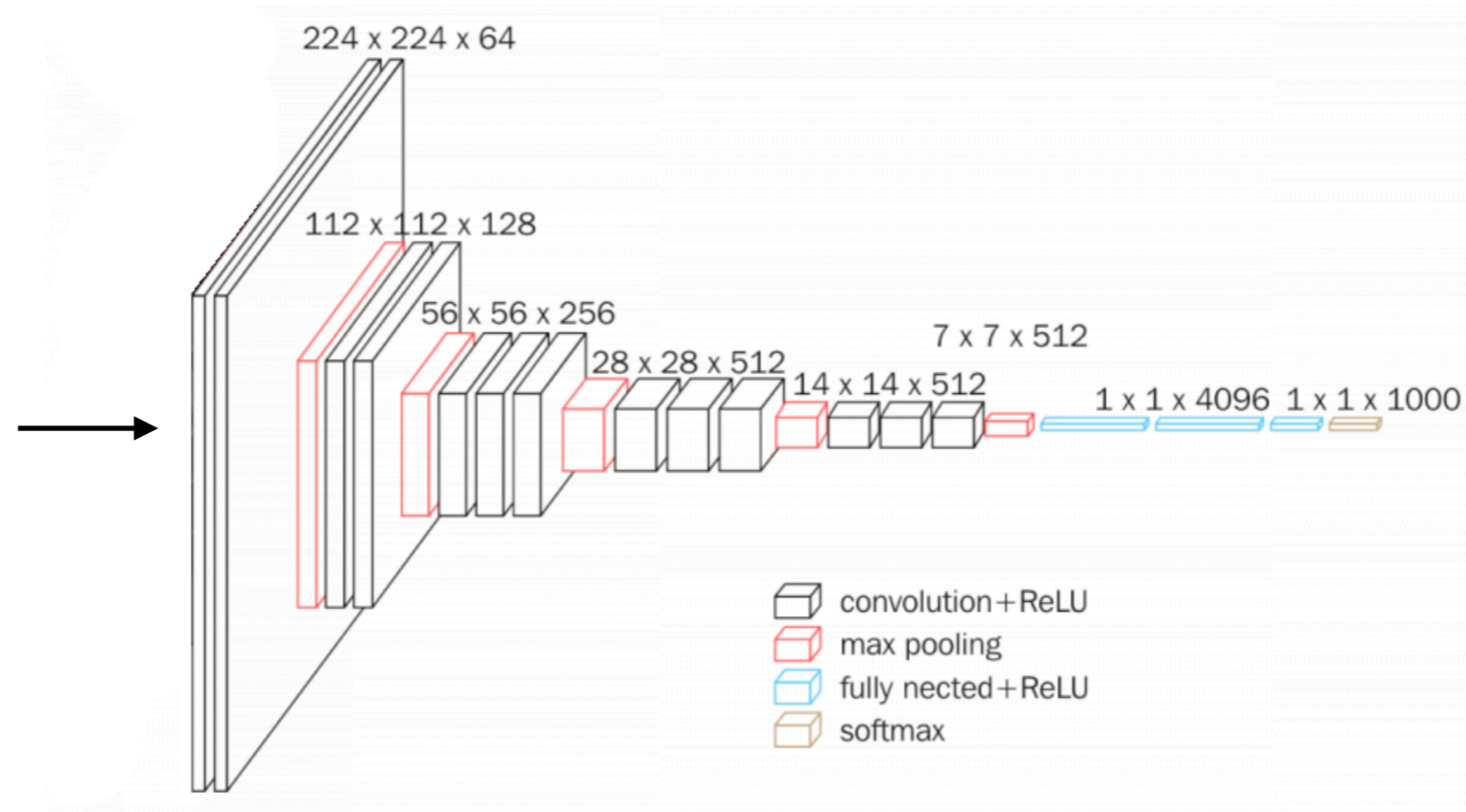
X

We *estimate* the pixel values for the image X so that:

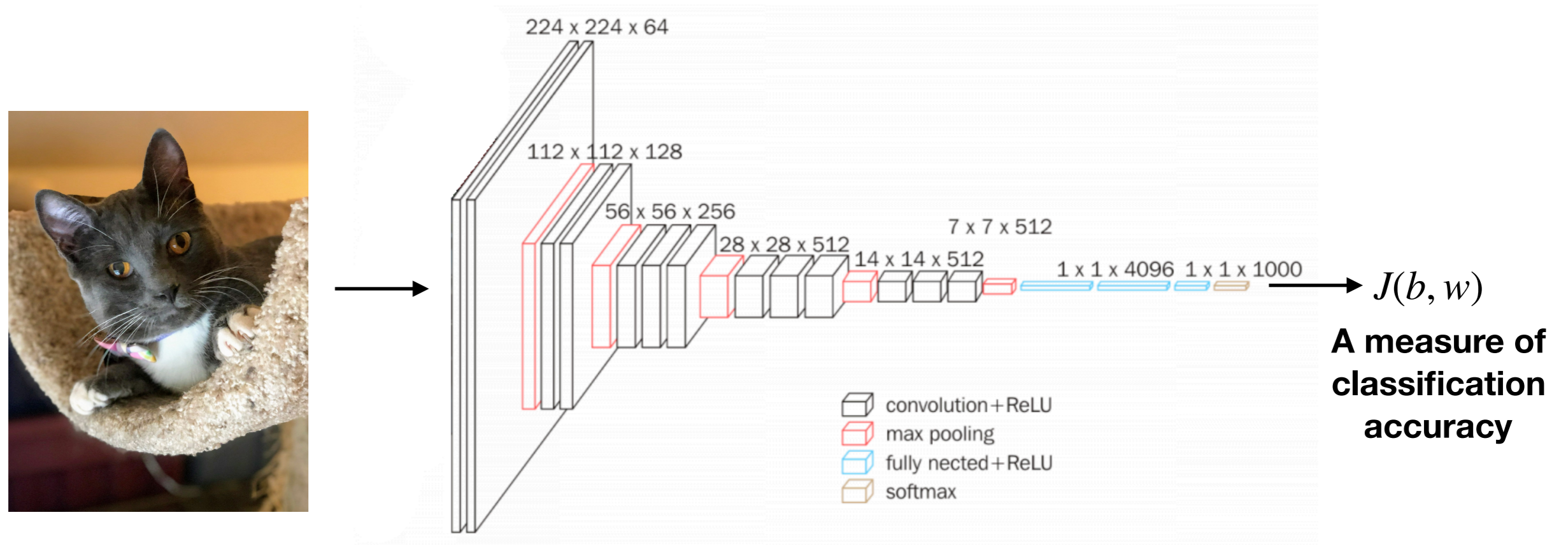
$$\text{content}(C) \approx \text{content}(X)$$

$$\text{style}(S) \approx \text{style}(X)$$

How We've Used Neural Networks So Far



How We've Used Neural Networks So Far

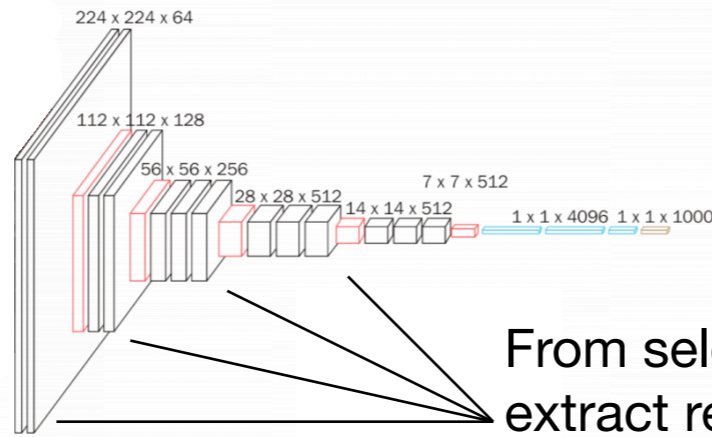


Estimation: given input data, minimize loss function by gradient descent:

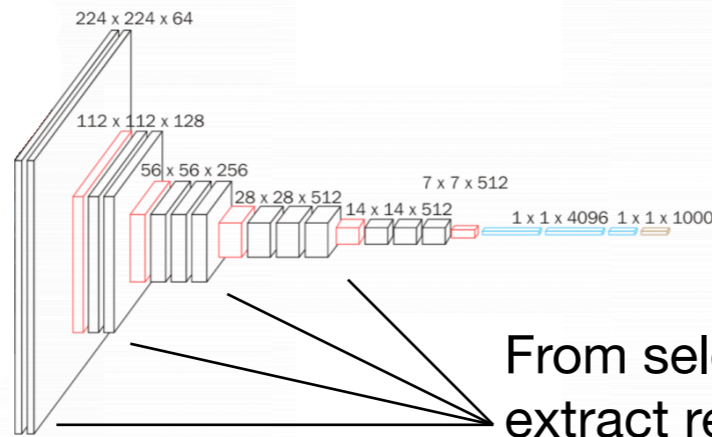
$$b^{[l]} \leftarrow b^{[l]} - \alpha \frac{\partial}{\partial b^{[l]}} J(b, w)$$
$$W^{[l]} \leftarrow W^{[l]} - \alpha \frac{\partial}{\partial W^{[l]}} J(b, w)$$

Use of Neural Networks for Style Transfer

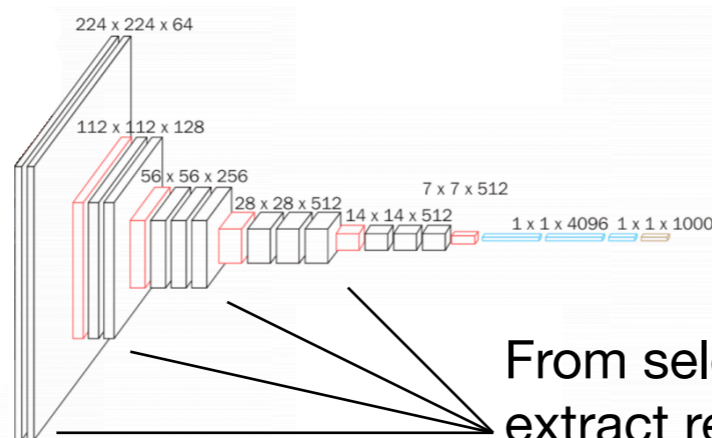
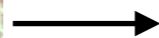
S



C



X



Loss is a measure of similarity of content and style:

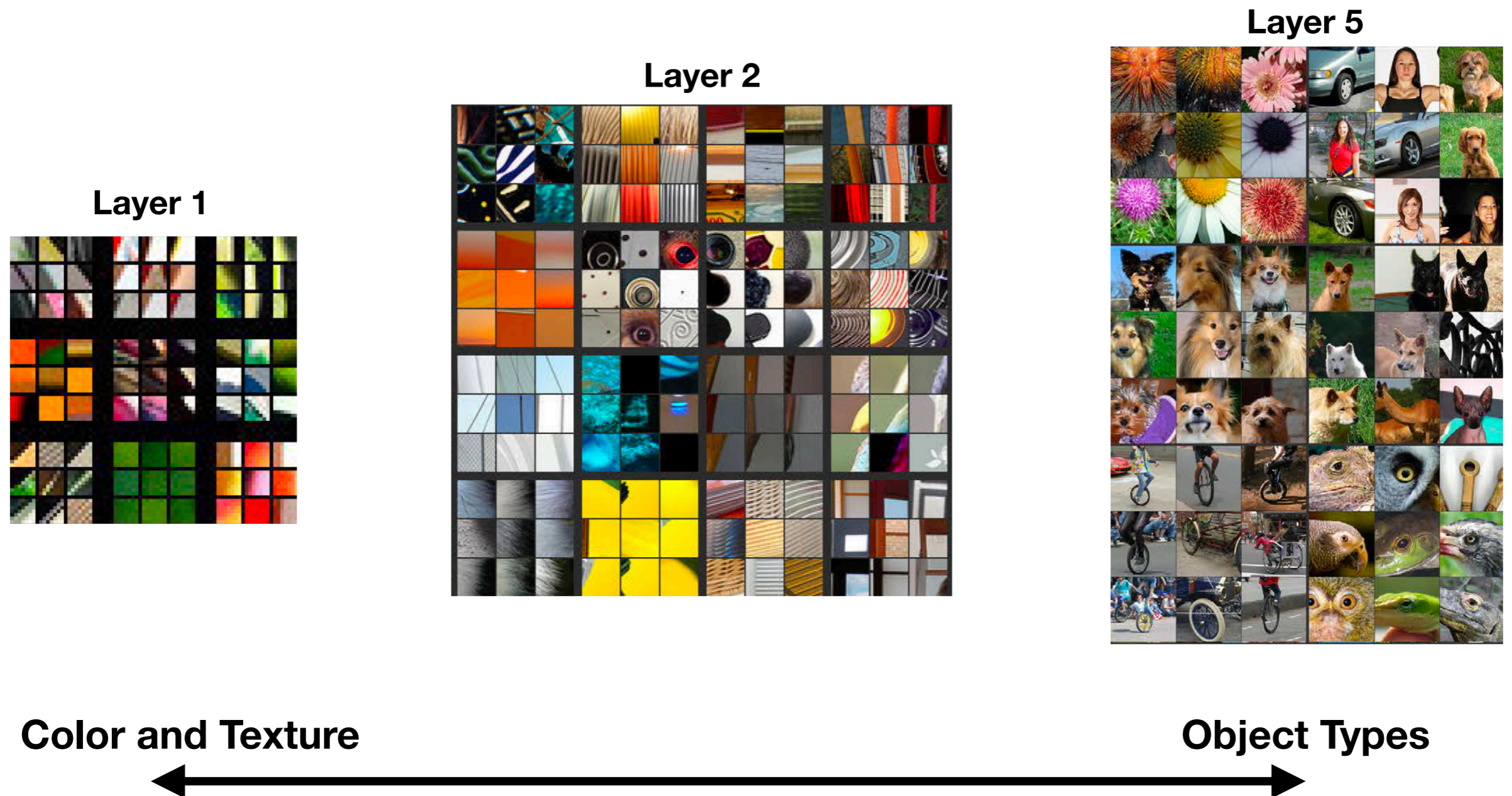
$$J(X) = [style(S) - style(X)]^2 + [content(C) - content(X)]^2$$

Gradient descent is used to estimate X:

$$X \leftarrow X - \alpha \frac{\partial}{\partial X} J(X)$$

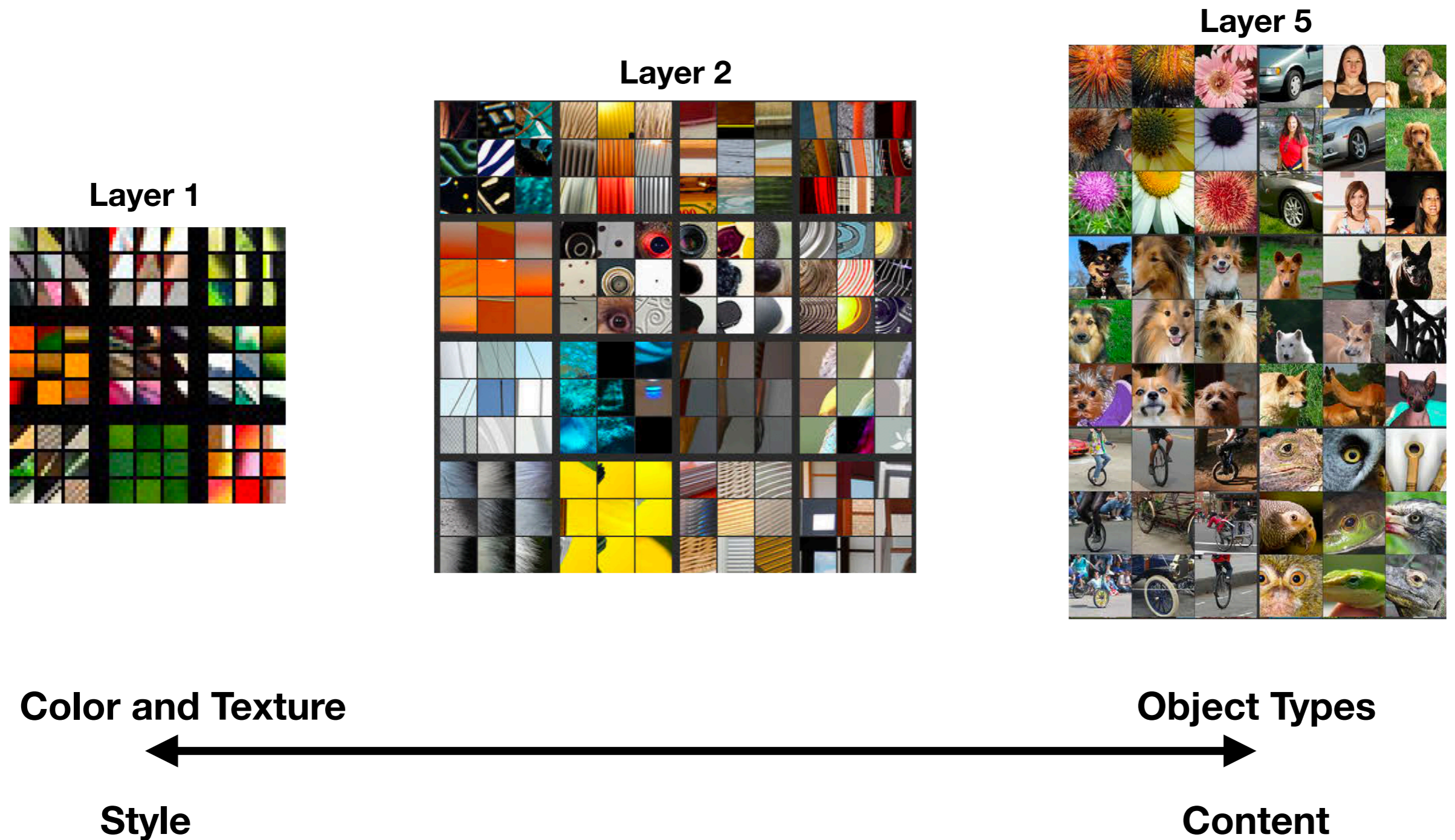
Measuring Content and Style

Figures from Zeiler and Fergus "Visualizing and Understanding Convolutional Networks" (2013)



Measuring Content and Style

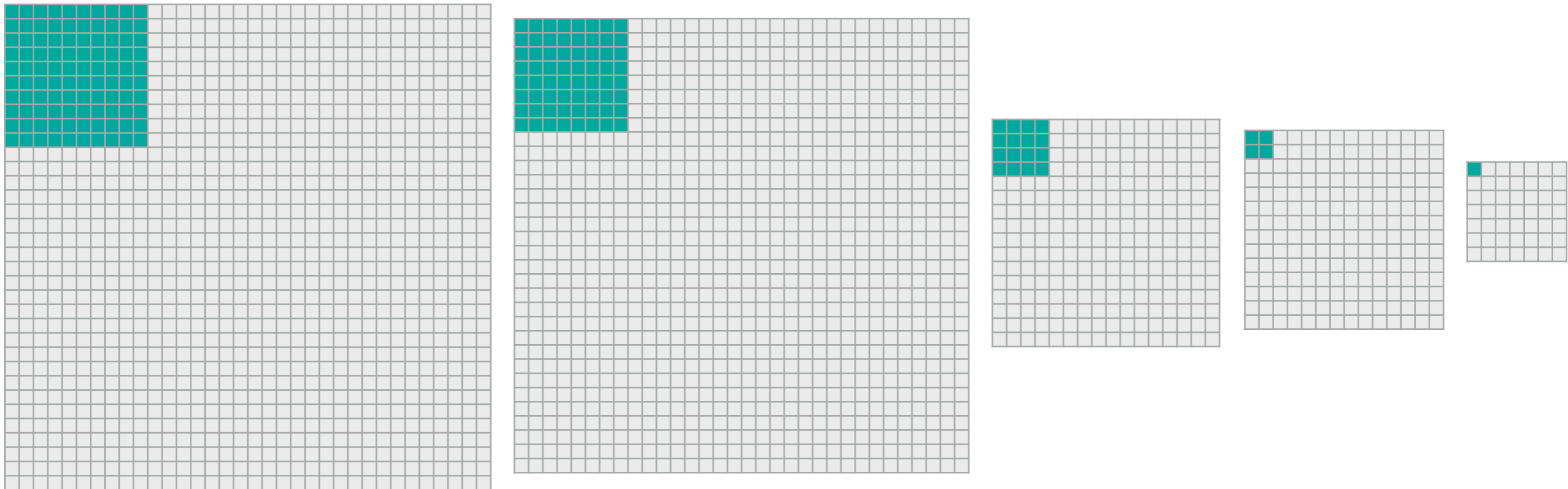
Figures from Zeiler and Fergus “Visualizing and Understanding Convolutional Networks” (2013)



Effective Receptive Field

- How many pixels of input unit are used to calculate a given activation in a later layer?

$34 \times 34 \longrightarrow 32 \times 32 \longrightarrow 16 \times 16 \longrightarrow 14 \times 14 \longrightarrow 7 \times 7$



- Activations in intermediate layers tell us about:
 - Use of color and texture at different positions in the input image
 - Types of objects at different positions in the input image

Measuring Content

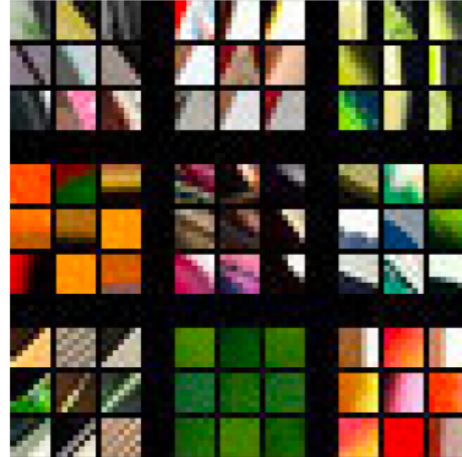
Layer 5



- Pick a layer relatively late in the network
- Its activations tell you what kind of objects are found at specific locations in the input image
- $content(image) = a^{[l]}$

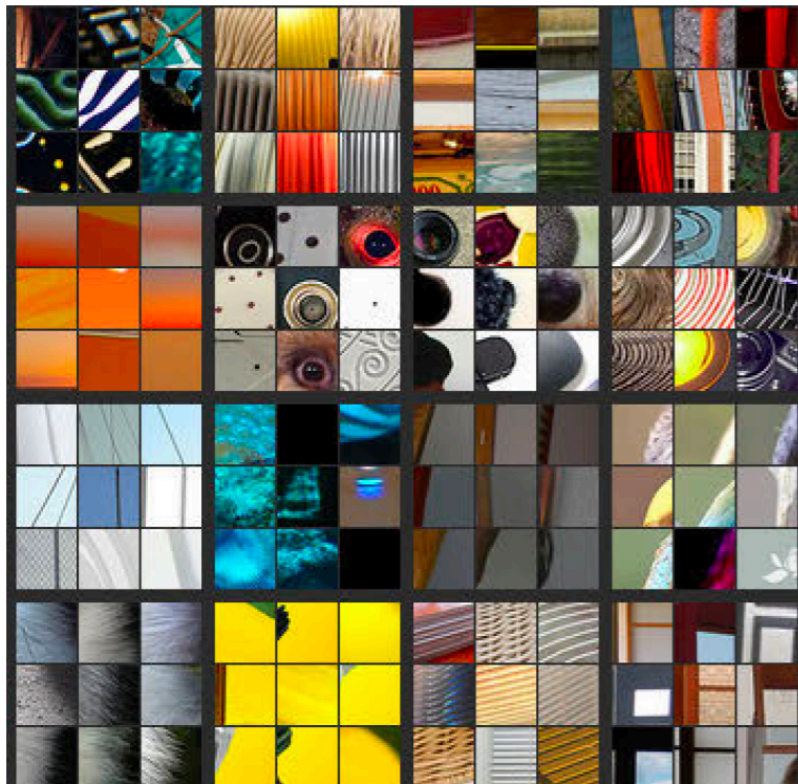
Measuring Style

Layer 1



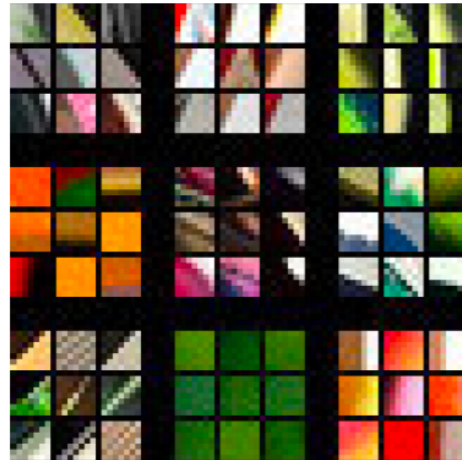
- Pick a layer (or in practice, multiple layers) relatively early in the network
- Its activations tell you what kind of colors and textures are found at specific locations in the input image

Layer 2



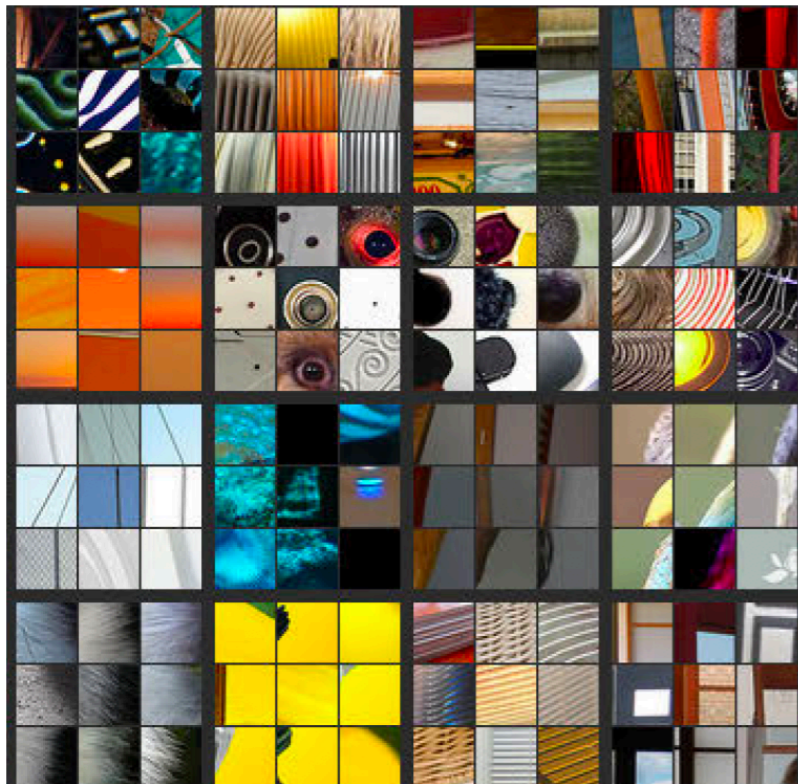
Measuring Style

Layer 1



- Pick a layer (or in practice, multiple layers) relatively early in the network
- Its activations tell you what kind of colors and textures are found at specific locations in the input image

Layer 2



Now what?

- We don't necessarily want to match use of color and texture in specific positions of the style image
- We want to match use of color and texture across the full image

