

Huget & Wiesel: "Thought experiments" - visualized stimuli of what cat saw.

By understanding what CNNs "see", we can tap into its creativity.

Visualizing neuron response  $\rightarrow$  look for certain images that "excite" specific units  $\rightarrow$  see which pixels maximally "influence" value @ filter.

For idea 1: use gradients  $\left( \frac{d \text{unit value}}{d \text{point at image}} \right)$ . Higher = more influence. Use backprop.

"guided backprop" - a hack to drastically improve quality of "vision" by "zeroing out" negative gradients at each layer ( $\sim \text{ReLU}$ ).

Idea 2: "Optimize" image to maximally excite filter.  $x \leftarrow \arg \max_x S(x) + R(x)$

More nuanced  $R(x)$ ?  $\rightarrow$  "blur" image after updating image and zero out small values.

activations  $\uparrow$  regularizer to prevent "crazy" images (e.g.  $d||x||_2$ )

Using backprop to modify features - "transport" feature distributions to other images.

DeepDream - visualize classes thru "forced hallucination"

$\rightarrow$  literally, set  $dx = x$  (gradient to activations) to accentuate further

Style Transfer: extracting style, updating spatial positions of content

We can quantify style as a co-occurrence of features  $\rightarrow \text{Cov}_{km} = \frac{1}{N} [f_k^l f_m^l]$

$\uparrow$  average over all positions in image.

new image  $= x \leftarrow \arg \min_x \mathcal{L}_{\text{style}}(x) + \mathcal{L}_{\text{content}}(x)$

Gram matrix:  $G_{km} = \text{Cov}_{km}$

$$\mathcal{L}_{\text{style}}(x) = \sum_k \sum_m \left( G_{km}^l - A_{km}^l(x) \right)^2 w_{km}$$

$\uparrow$  source image activation  $\uparrow$  new image activation

"style"

How can we quantify content? "Match" the features.

$\uparrow$  For a specific layer (choosing said layer is also important)

$$\mathcal{L}_{\text{content}}(x) = \sum_j \sum_k \left( f_{ijk}^l(x_{\text{content}}) - f_{ijk}^l(x) \right)^2$$