

Lecture Notes for **Neural Networks** **and Machine Learning**



More CNN Circuits



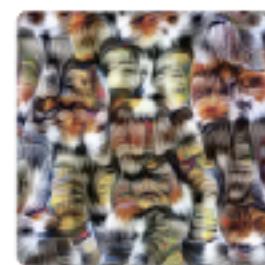
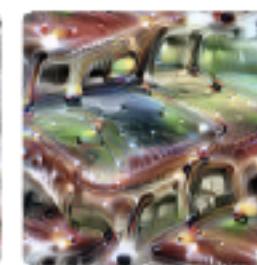
Logistics and Agenda

- Logistics
 - Grading Update
- Agenda
 - Last Time: Visualization + Circuits in CNNs
 - More Circuits
 - Continued Town Hall
 - Branch Specialization and Universality
 - Next Time:
 - ◆ Student Paper Presentation
 - ◆ Fully Convolutional Networks



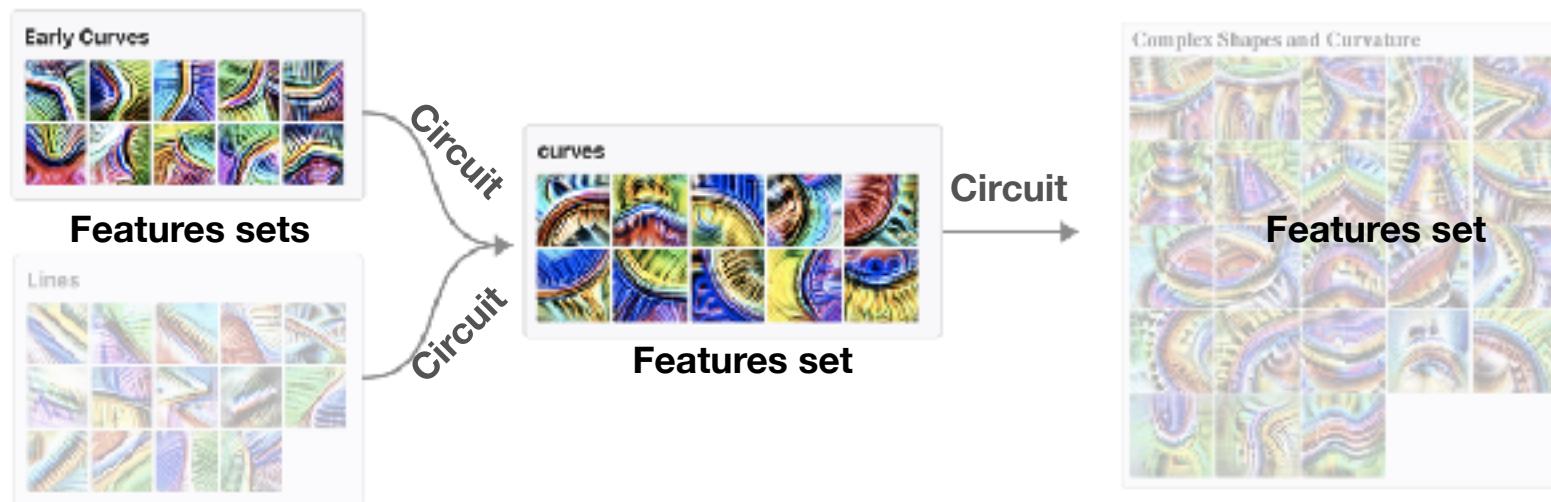
Last Time

- *Features are connected by weights, forming circuits*
- *"All neurons in our network are formed from linear combinations of neurons in the previous layer, followed by ReLU. If we can understand the features in both layers, shouldn't we also be able to understand the connections between them?"*
- *"Once you understand what features they're connecting together... You can literally read meaningful algorithms off of the weights."*



From Features to Circuits

- *Features are connected by weights, forming circuits*
- *“All neurons in our network are formed from linear combinations of neurons in the previous layer, followed by ReLU. If we can understand the features in both layers, shouldn’t we also be able to understand the connections between them?”*
- *“Once you understand what features they’re connecting together... You can literally read meaningful algorithms off of the weights.”*



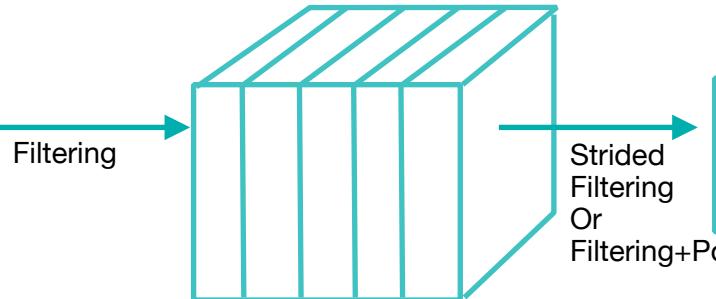
<https://microscope.openai.com/models/inceptionv1/>

51



Review of CNN Structure

Input Activations



Output Activations



```
keras_layer = model.get_layer('block3_conv1')
layer_output = keras_layer.output
weights_list = keras_layer.get_weights() # list
filters = weights_list[0]
biases = weights_list[1]
```

block4_conv1 activation size is (None, None, None, 256) (batch x H x W x filter)
block4_conv1 filters is of shape (3, 3, 128, 256) ... (k x k x channels x filters)
block4_conv1 biases is of shape (256,)
one filter in block4_conv1 is (3, 3, 128)



```
# lets look at the shapes of some of the filters above
keras_layer = model.get_layer('block3_conv1')
layer_output = keras_layer.output
weights_list = keras_layer.get_weights() # list of filter, the biases
filters = weights_list[0]
biases = weights_list[1]

# print out some specifics of how the filter is saved
print('block4_conv1 activation size is ', layer_output.get_shape(), '(batch x H x W x ')
print('block4_conv1 filters is of shape',filters.shape, '...(k x k x channels x filters)')
print('block4_conv1 biases is of shape',biases.shape)

idx = 32
print('one filter in block4_conv1 is ', filters[:,:,:,:,idx].shape )
channel = 2
print('one channel in the the filter is', filters[:,:,:,:,channel,idx].shape)
print('The weights of that channel in the filter are:\n', filters[:,:,:,:,channel,idx])
print('The bias of the filter is:',biases[idx])

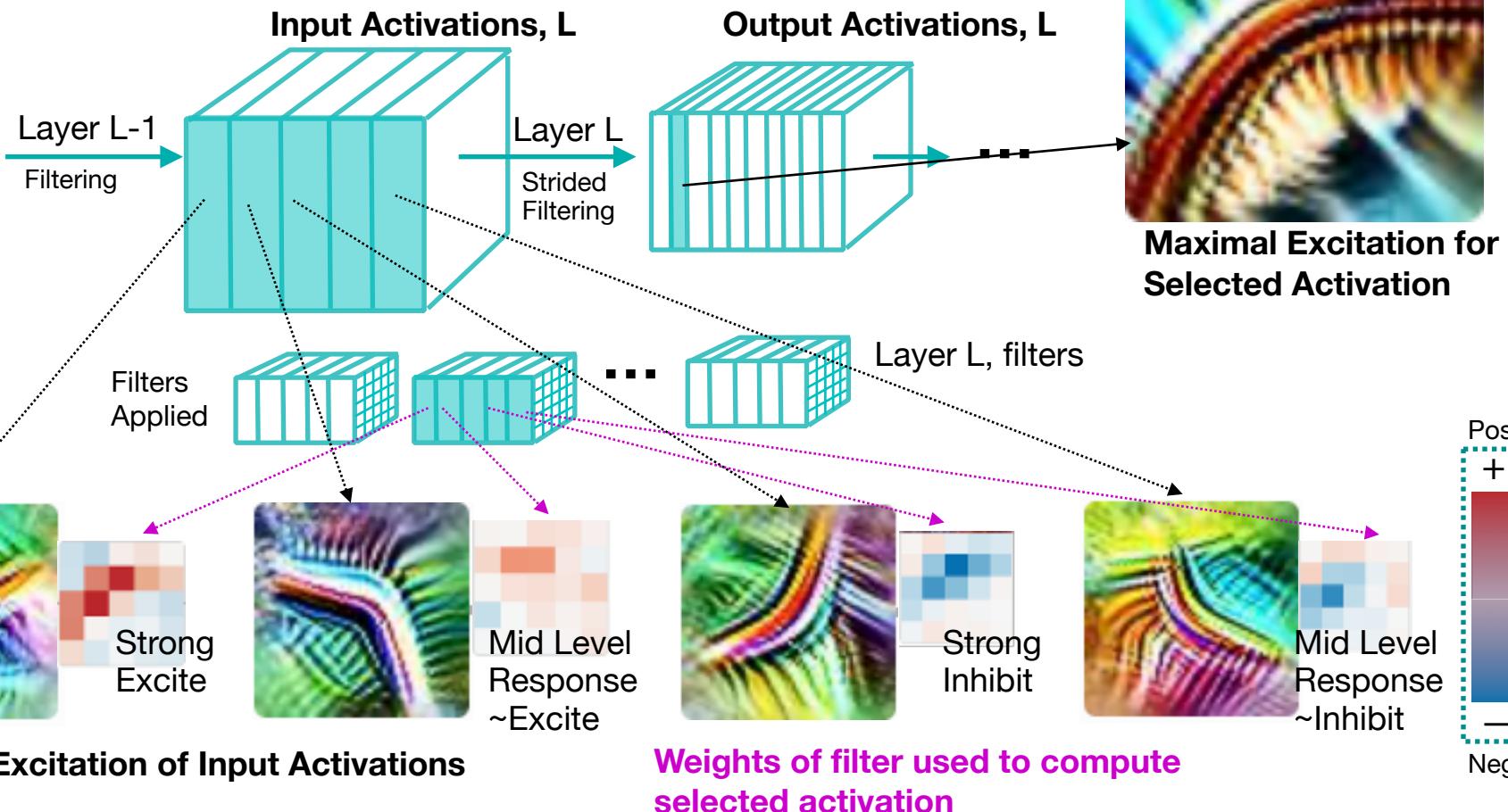
block4_conv1 activation size is (None, None, None, 256) (batch x H x W x filter)
block4_conv1 filters is of shape (3, 3, 128, 256) ... (k x k x channels x filters)
block4_conv1 biases is of shape (256,)
one filter in block4_conv1 is (3, 3, 128)
one channel in the the filter is (3, 3)
The weights of that channel in the filter are:
 [[-0.03330493  0.01174345  0.03184387]
 [-0.04050588 -0.02253938  0.02304637]
 [-0.00191393 -0.01501364  0.02783429]]
The bias of the filter is: 0.030420048
```



What weights comprise a circuit?

Structure of
Each Tensor:

Channels x Rows x Columns



Maximal Excitation of Input Activations

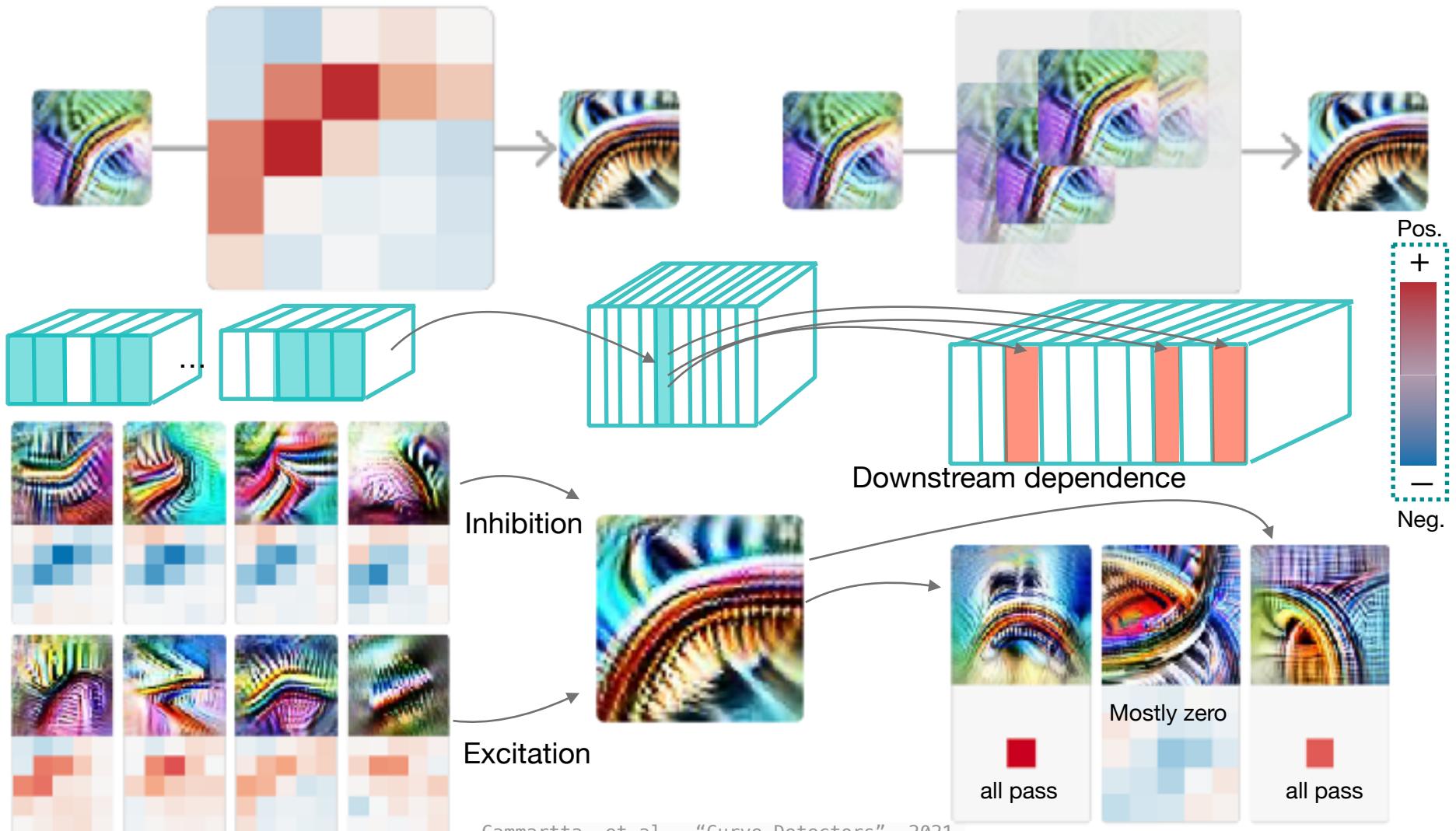
Weights of filter used to compute selected activation

<https://distill.pub/2020/circuits/curve-circuits/>



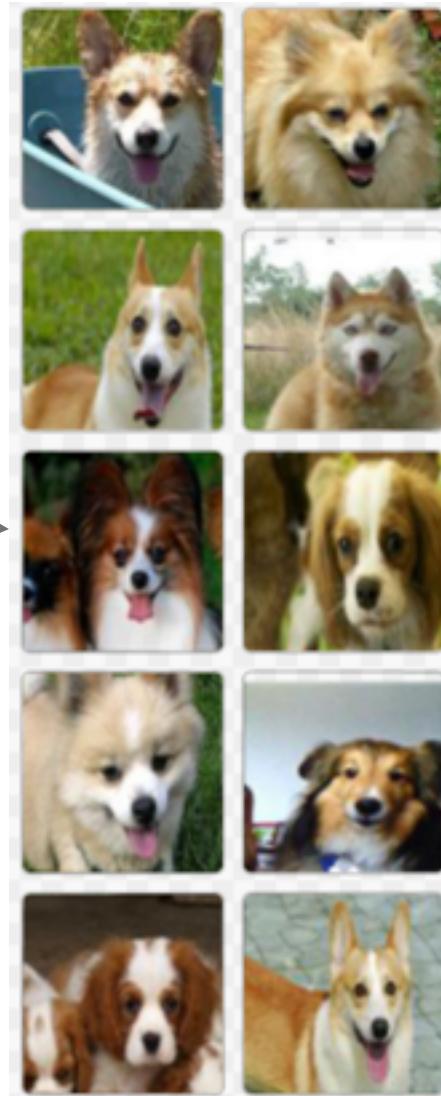
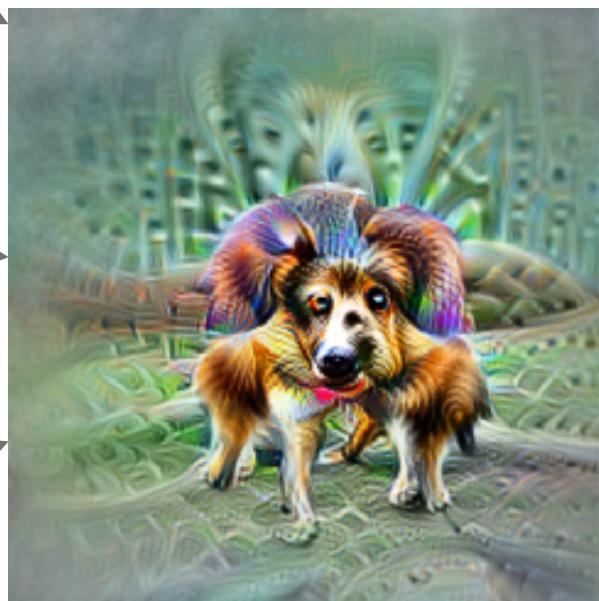
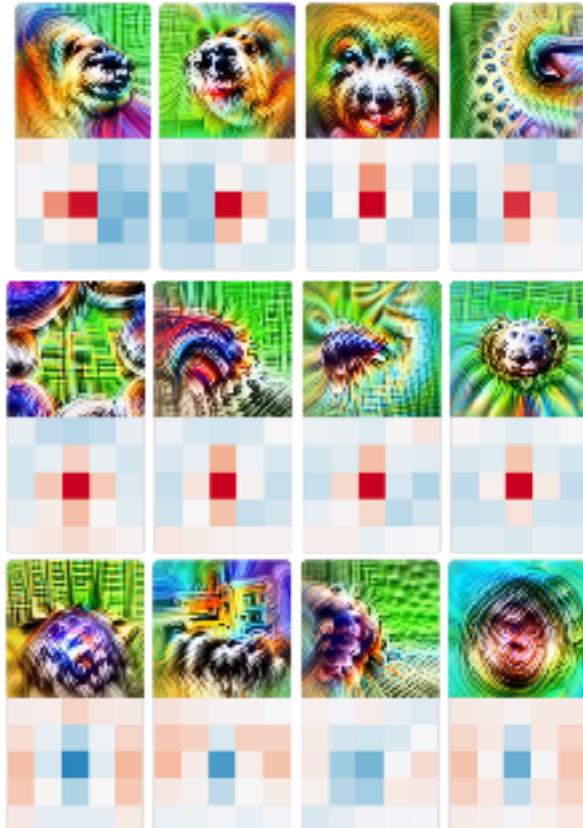
Example: Circuit for Better Curve Detection

If we visualize the 5x5 Conv Filter, we can see that this becomes a Superposition of Early Curves



Another Example: Dog head

Compact Circuit Visualization

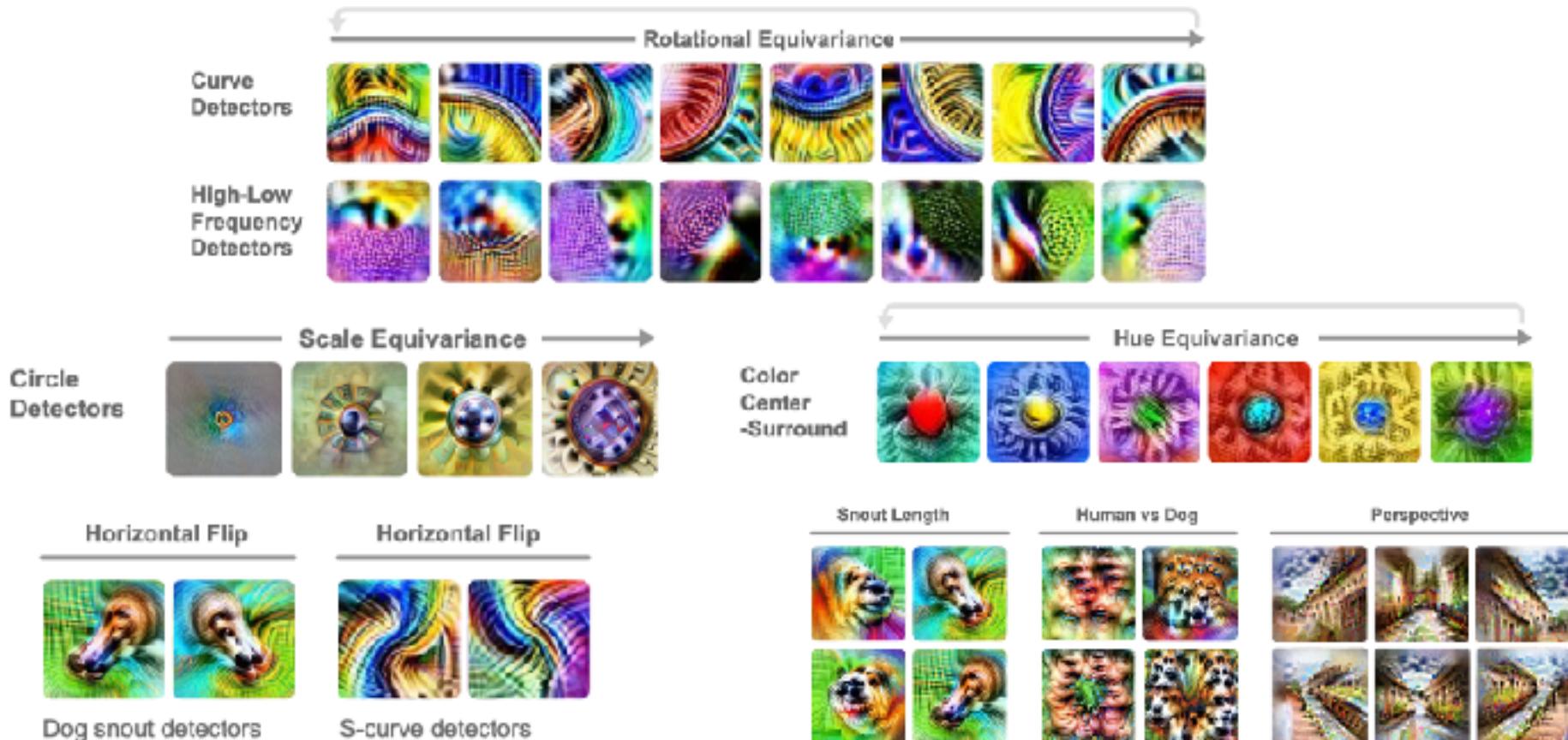


This example is also **polysemantic** due to the "**espresso maker**" class also being excited by this...



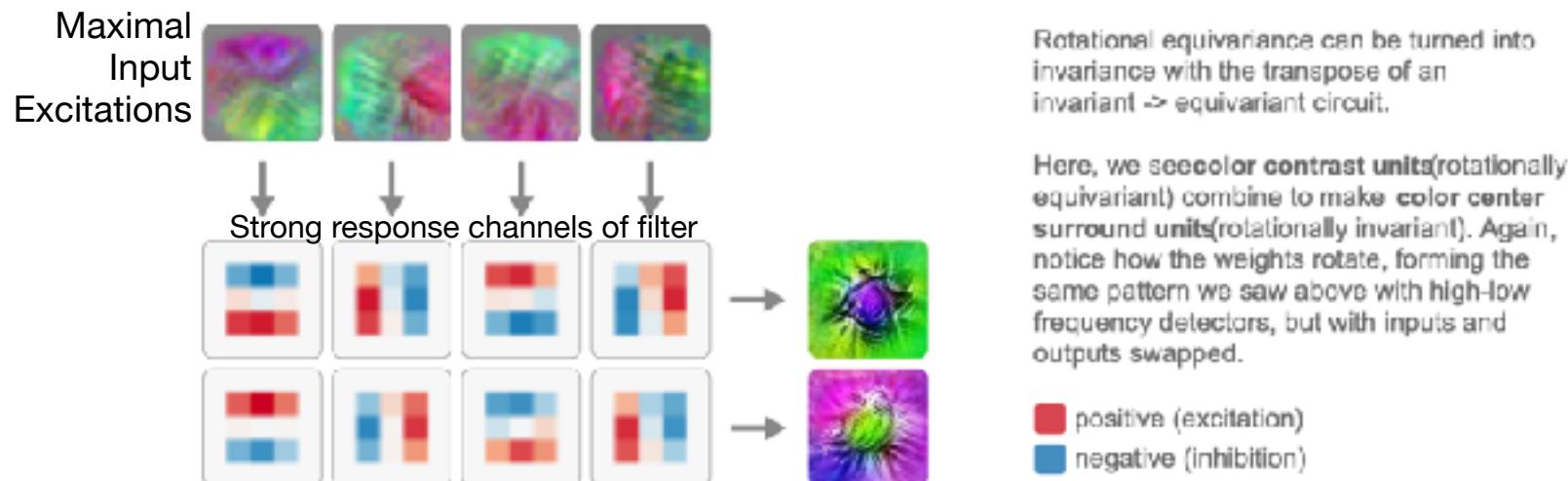
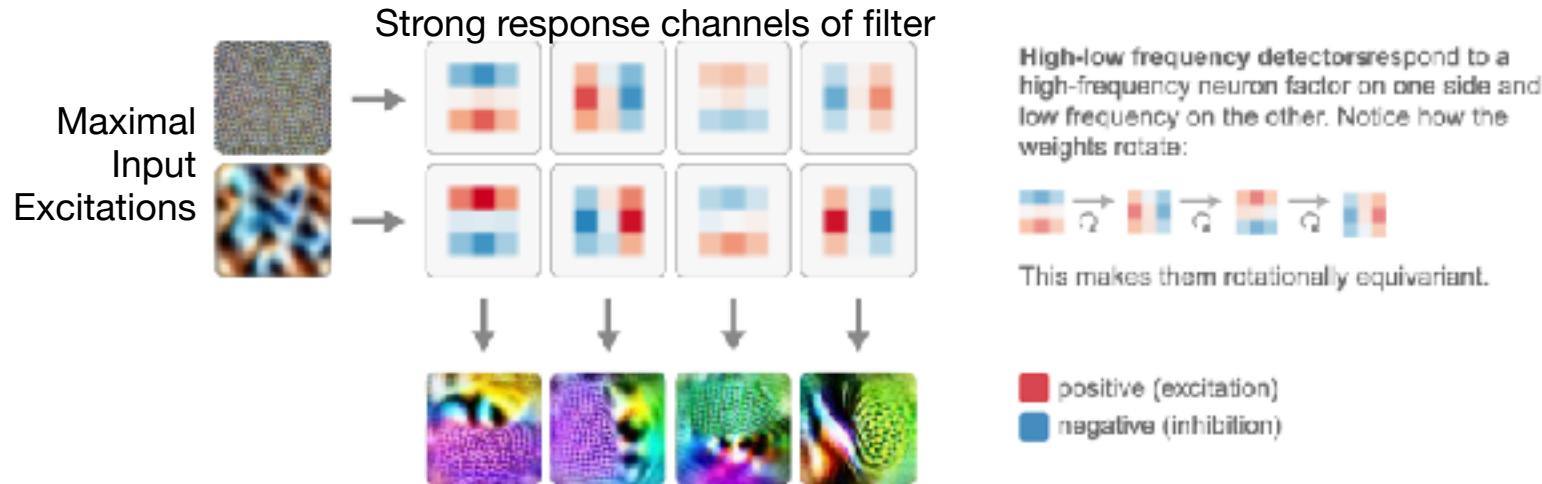
Equivariant Circuits

- Many features that are part of a circuit are clearly designed for rotation, hue, and other invariance



Equivariant circuits: a Motif

- Possible to reveal patterns of circuits via sets of weights



Olah, et al., "Naturally Occurring Equivariance in NN", 2021.



Circuits Town Hall

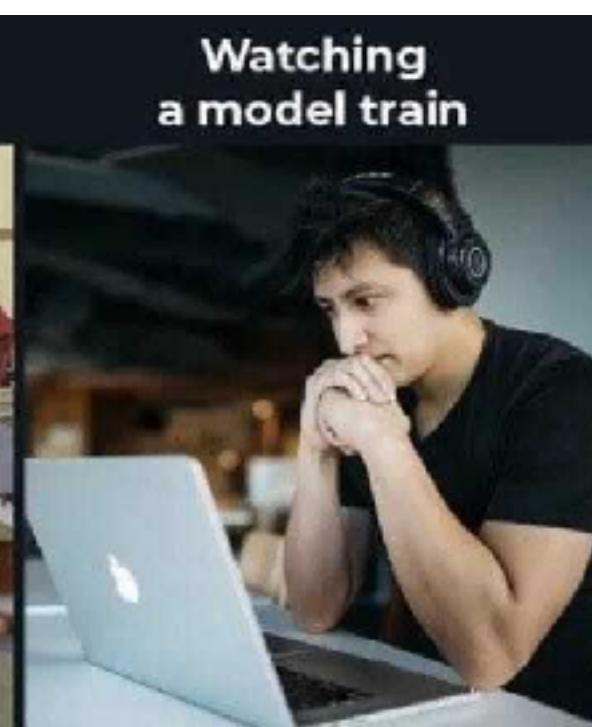
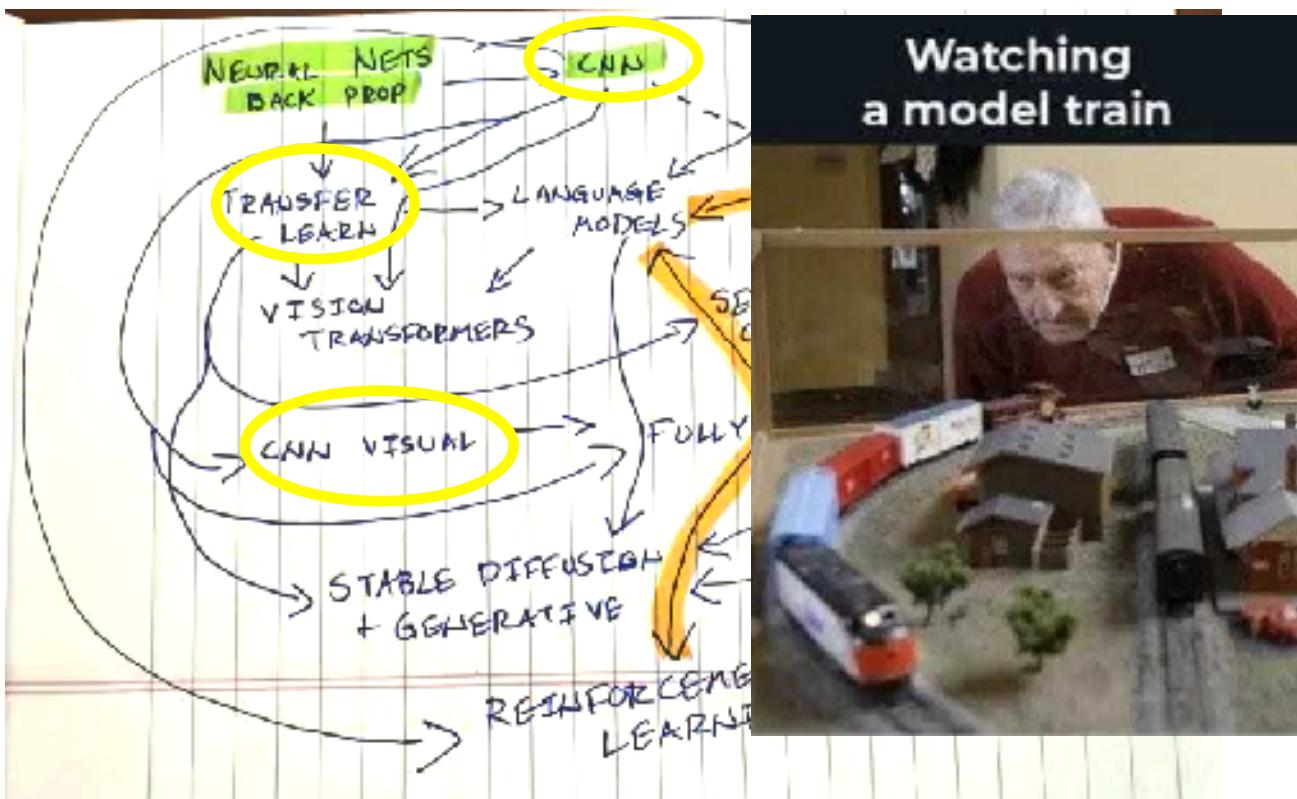


Figure for Circuits Lab



Structure of Each Tensor:

Channels x Rows x Columns

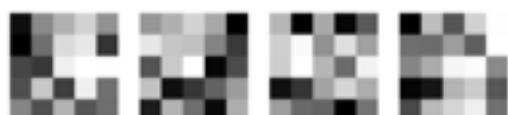
Input Activations

Filtering

Rows

5. Look at the input activations for each channel of the multi-channel filter, see which are the most influential.

Multi-channel Filters



1. Choose a network and middle layer to analyze. This example shows activations from a VGG layer with 512 channels.

2. Choose Output Activation of Interest

Output Activations

A diagram showing a stack of vertical teal-colored bars of varying heights. A pink arrow points from the top of the stack towards the right side. To the left of the stack is a green arrow pointing upwards, and to the right is a grey arrow pointing downwards.

3. Find multi-channel filter that is responsible for selected activation

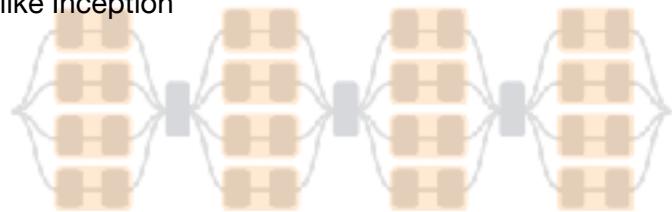
4. Look at each channel of the multi-channel filter. Each channel can be thought of as a filter applied to the activations in the previous layer.

Branch Specialization

Feedforward



Branched,
like Inception



Residual,
like ResNet

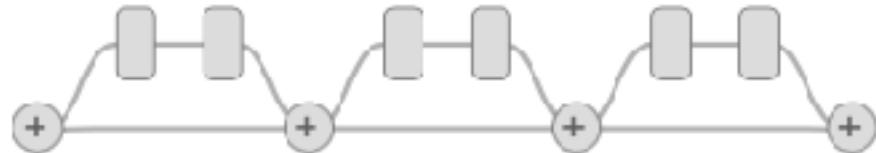


Primitives

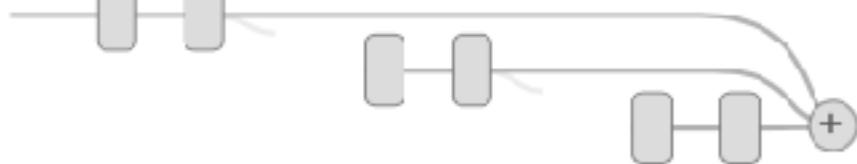


More Complex
Downstream Circuits

Two ways of looking at residual networks



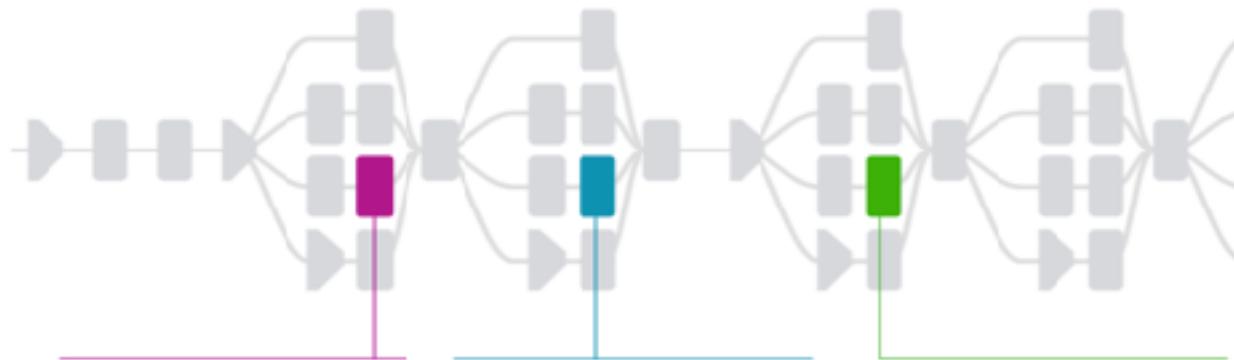
Exponential number of possible branches



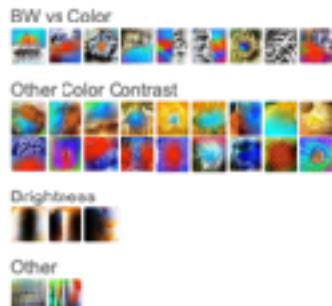
- Specialized branches are consistent across many architectures, support the idea of an interconnected graph of operations



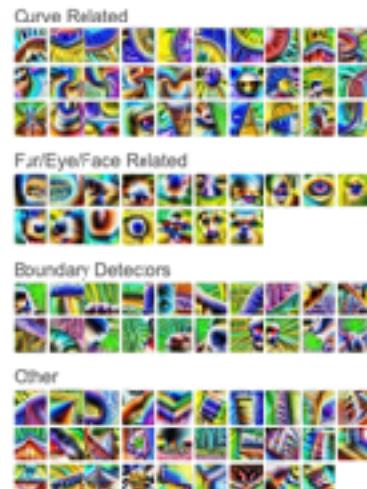
Branch Specialization



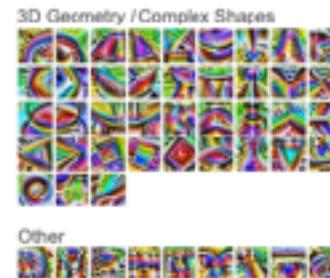
mixed3a_5x5: The 5x5 branch of mixed3a, a relatively early layer, is specialized on color detection, and especially black-and-white vs. color detection.



mixed3b_5x5: This branch contains all 30 of the curve-related features for this layer (all curves, double curves, circles, spirals, S-shape and more features, etc). It also contains a disproportionate number of boundary, eye, and fur detectors, many of which share sub-components with curves.



mixed4a_5x5: This branch appears to be specialized in complex shapes and 3D geometry detectors. We don't have a full taxonomy of this layer to allow for a quantitative assessment.



Motifs appear in branches. Similar clusters of operations can be found across different architectures

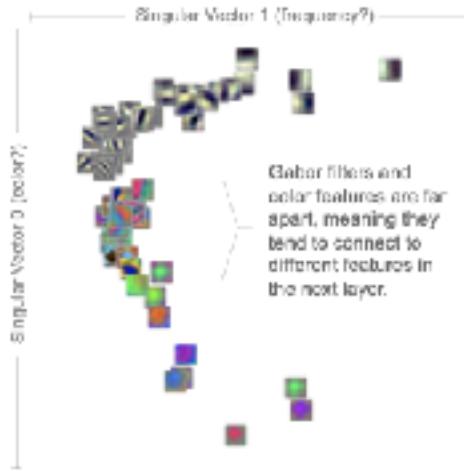


Investigating Connection Clusters via SVD

InceptionV1 (tf-slim version) trained on ImageNet

The first singular vector separates color and black and white, meaning that's the largest dimension of variation in which neurons connect to which in the next layer.

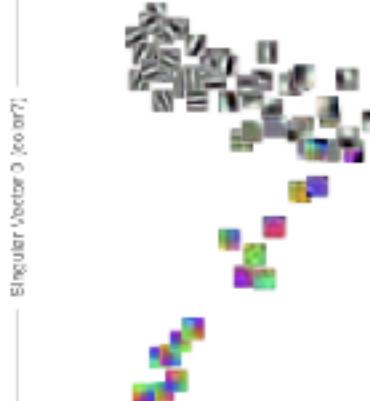
Neurons in the first convolutional layer organized by the left singular vectors of $[W]$.



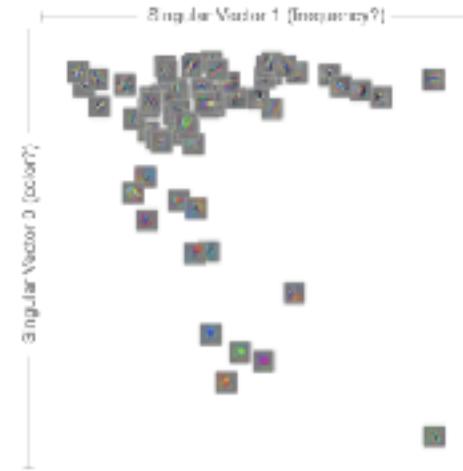
InceptionV1 trained on Places365

One more, the first singular vector separates color and black and white, meaning that's the largest dimension of variation in which neurons connect to which in the next layer.

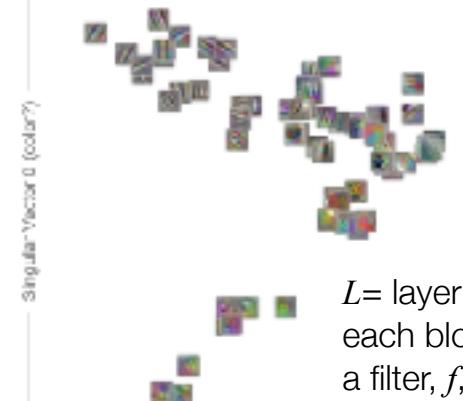
Singular Vector 1 (frequency?)



Neurons in the second convolutional layer organized by the right singular vectors of $[W]$.



Singular Vector 1 (frequency?)



- Singular Value Decomposition (SVD) decomposes a matrix into three elements

- $M = U\Sigma V^T$
- U is eig-vec of MM^T
 V is eig-vec of $M^T M$
- U and V are orthogonal such that
 $U^T U = I \quad V^T V = I$
- Σ is a diagonal matrix of the singular values
- These values characterize the variability in a matrix

$$SVD(|\mathbf{W}^{(L)}|)$$

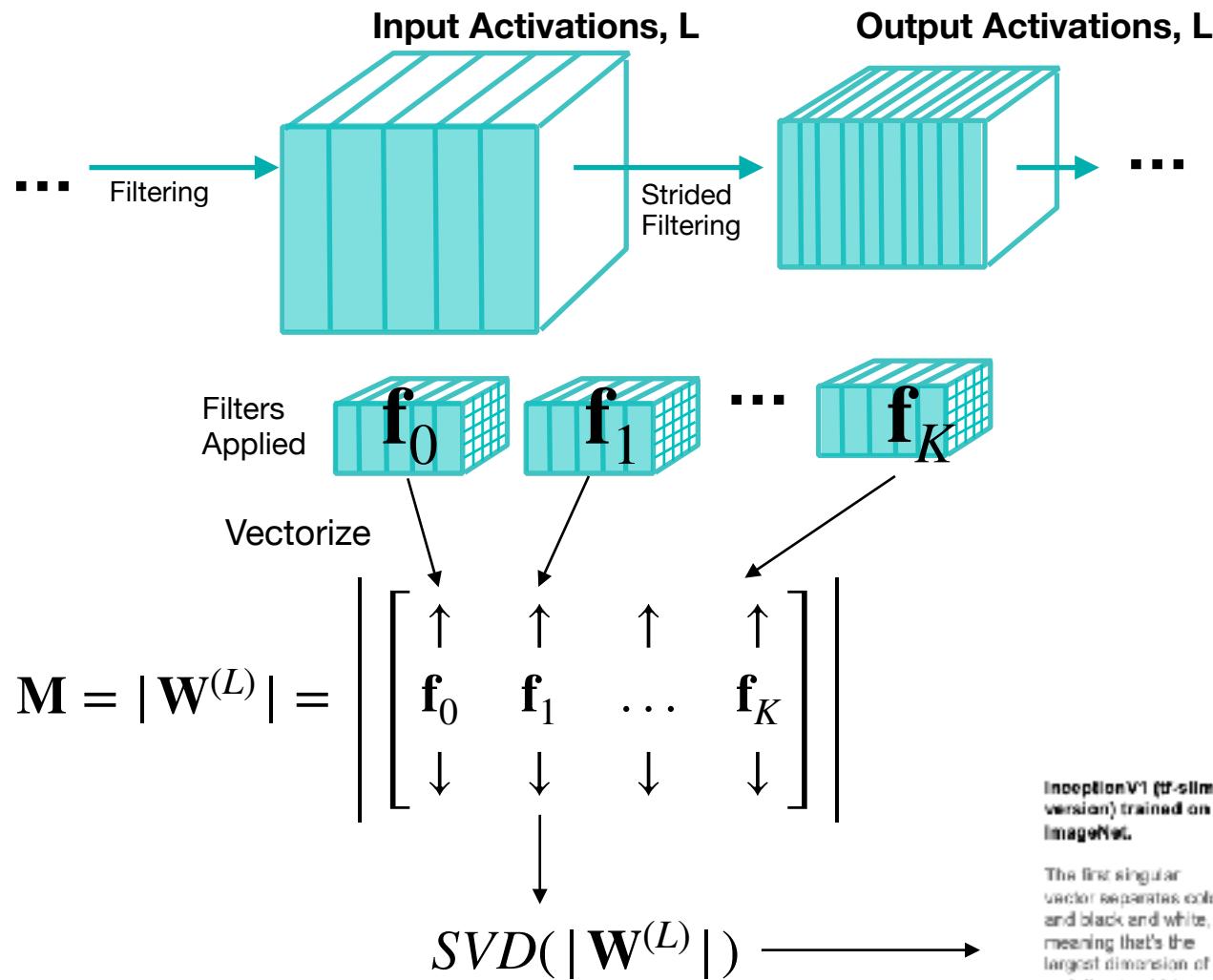
$L =$ layer
each block is a filter, f_i in layer
each image is the optimized excitation

Voss, et al., "Branch Specialization", Distill, 2021.



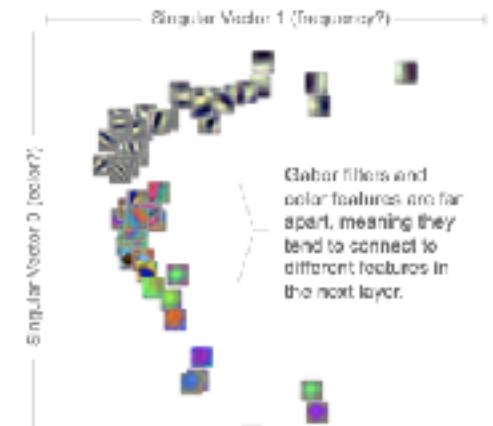
Structure of Each Tensor:

Channels x Rows x Columns



- Singular Value Decomposition (SVD) decomposes a matrix into three elements
 - $M = U\Sigma V^T$
 - U is eig-vec of MM^T
 - V is eig-vec of M^TM
 - U and V are orthogonal such that $UU^T=I \quad VV^T=I$
 - Σ is a diagonal matrix of the singular values
 - These values characterize the variability in a matrix

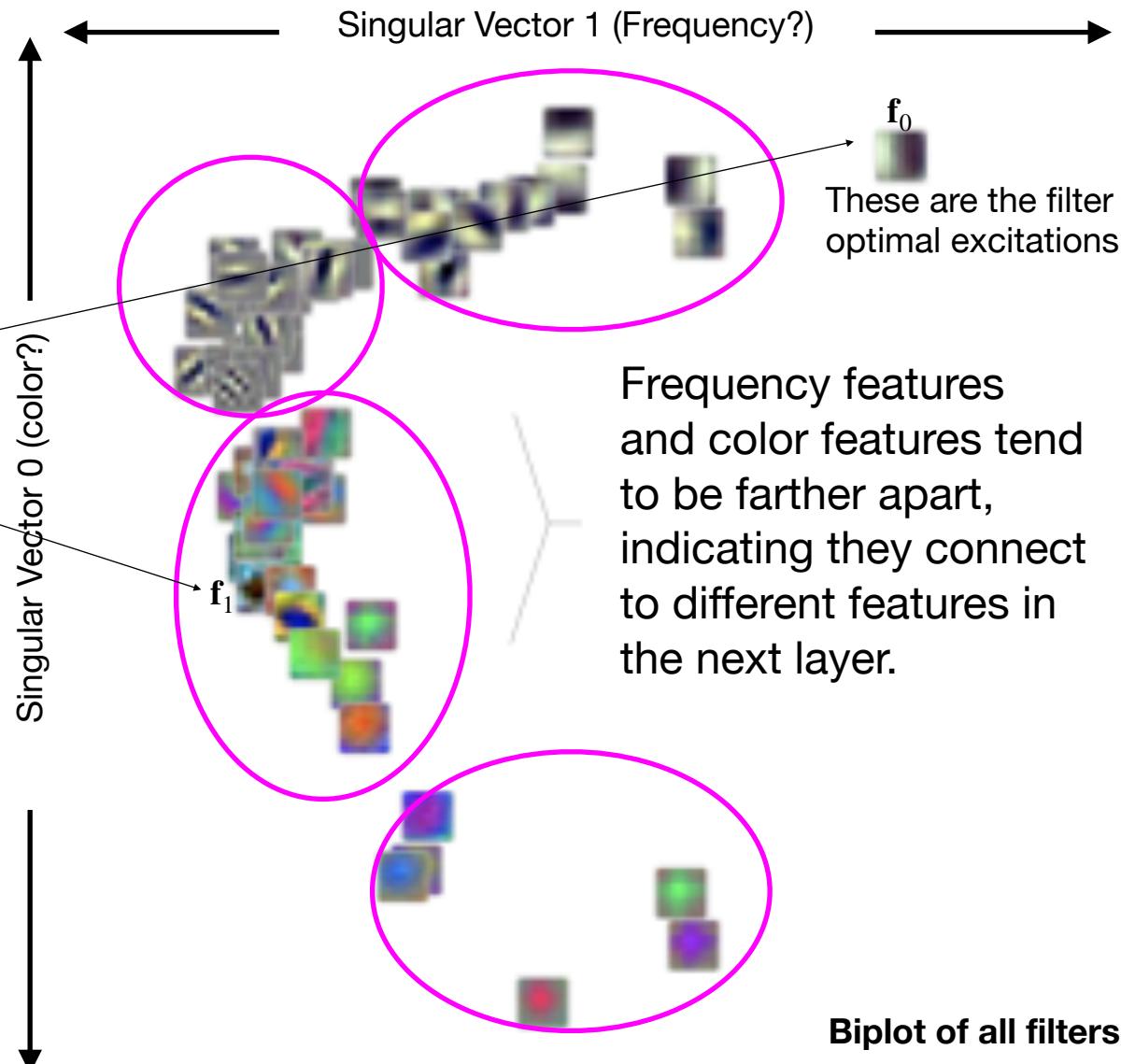
Neurons in the first convolutional layer organized by the left singular vectors of $|W|$.



$$SVD(|\mathbf{W}^{(L)}|) \rightarrow \mathbf{U}^{(L)}$$

Each \mathbf{U}
is weighted
sum of filters

$$\mathbf{U}^{(L)} = \begin{bmatrix} \mathbf{e}_0 & \mathbf{e}_1 & \dots & \mathbf{e}_K \\ \downarrow & \downarrow & \dots & \downarrow \end{bmatrix}$$

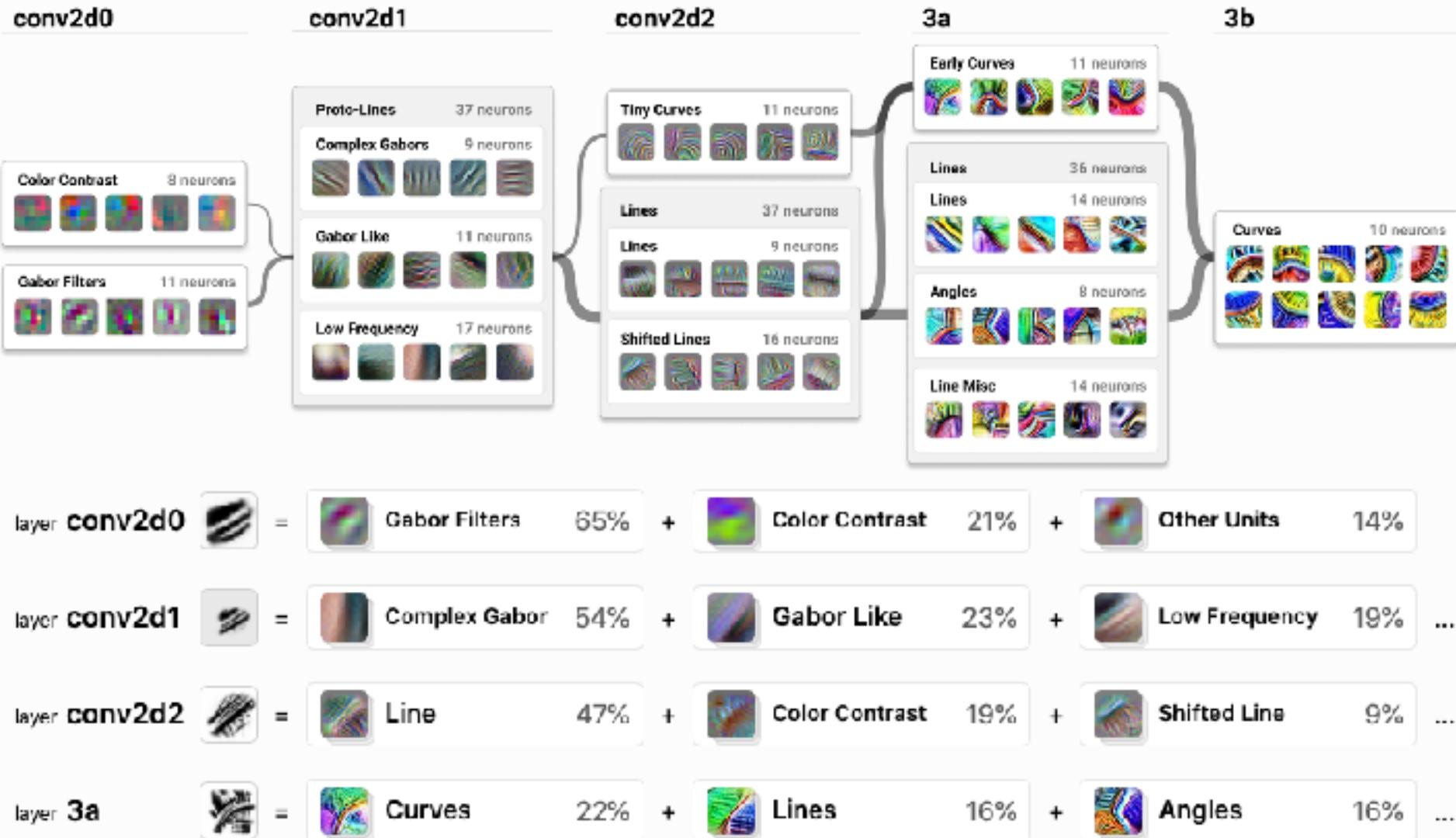


Clusters tend to be motifs, and separation of clusters reveals connections between layers (like edges in a graph)

Biplot of all filters



Neural Nets: Directed Graph of Circuits



Voss, et al., "Branch Specialization", Distill, 2021.



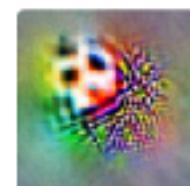
Universality of Circuits

- Analogous features and circuits form across models and tasks

Curve detectors



High-Low Frequency detectors



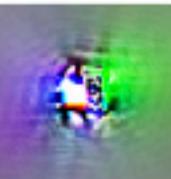
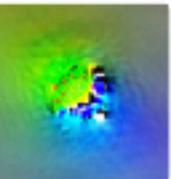
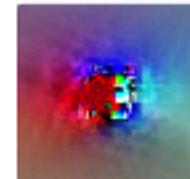
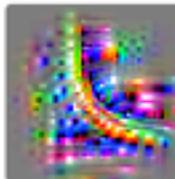
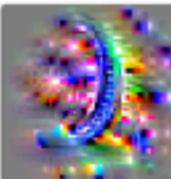
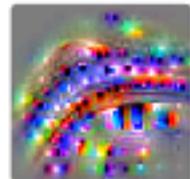
ALEXNET

Krizhevsky et al. [34]



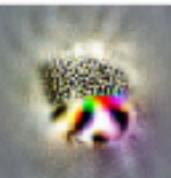
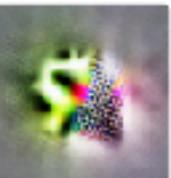
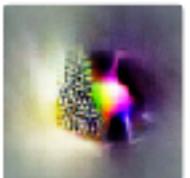
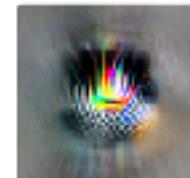
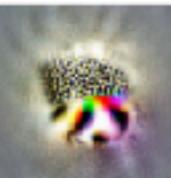
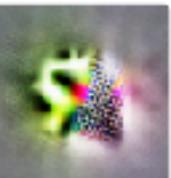
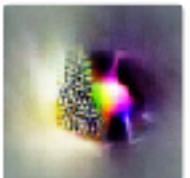
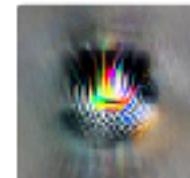
INCEPTIONV1

Szegedy et al. [26]



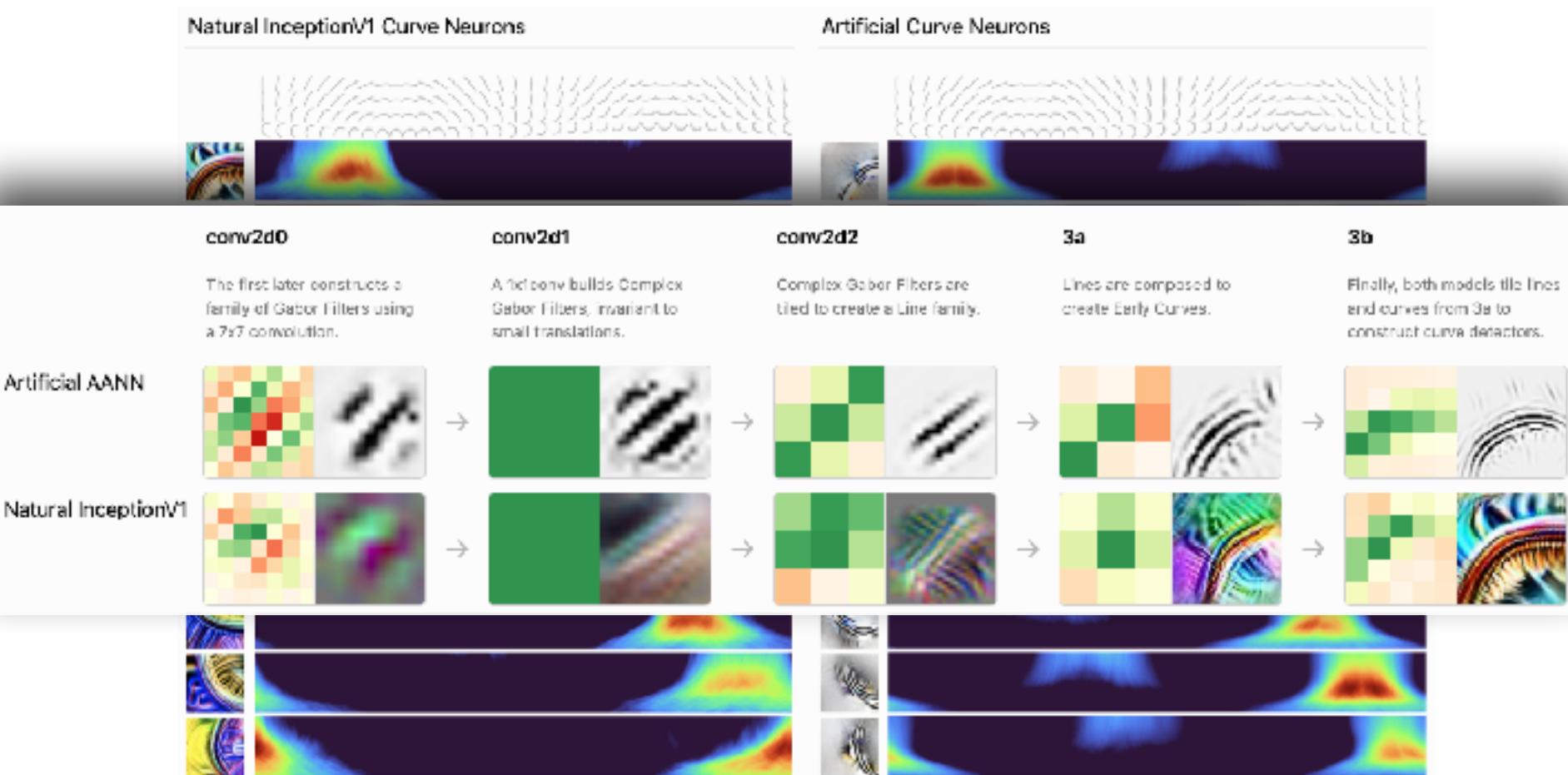
VGG19

Simonyan et al. [35]



Reverse Engineering a Circuit

- With assumption of what feature is, a circuit can be implemented by hand that nearly identically follows the assumed functionality



Closing Thoughts from OpenAI Researchers

Closing Thoughts

We take it for granted that the microscope is an important scientific instrument. It's practically a symbol of science. But this wasn't always the case, and microscopes didn't initially take off as a scientific tool. In fact, they seem to have languished for around fifty years. The turning point was when Robert Hooke published *Micrographia* [1], a collection of drawings of things he'd seen using a microscope, including the first picture of a cell.

Our impression is that there is some anxiety in the interpretability community that we aren't taken very seriously. That this research is too qualitative. That it isn't scientific. But the lesson of the microscope and cellular biology is that perhaps this is expected. The discovery of cells was a qualitative research result. That didn't stop it from changing the world.

<https://distill.pub/2020/circuits/zoom-in/>



Lecture Notes for Neural Networks and Machine Learning

CNN Circuits



Next Time:
Fully Convolutional Learning
Reading: Chollet 5.4

