Lecture Notes for

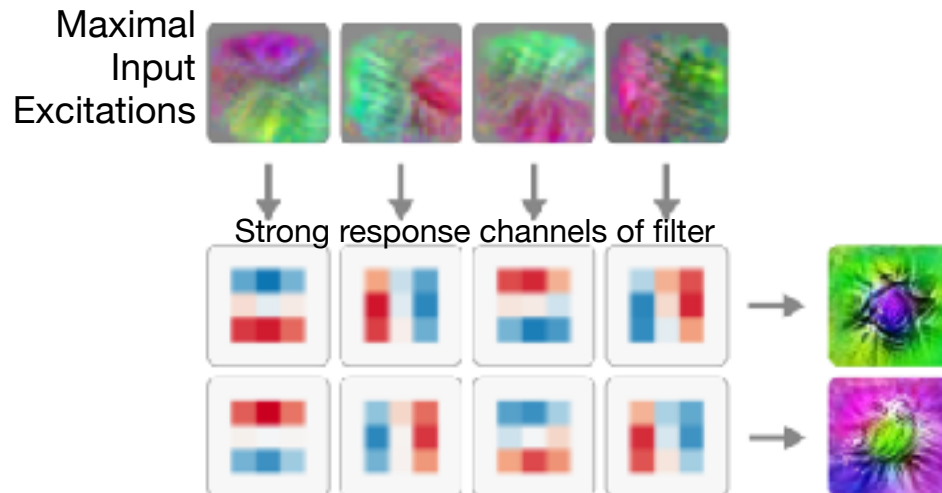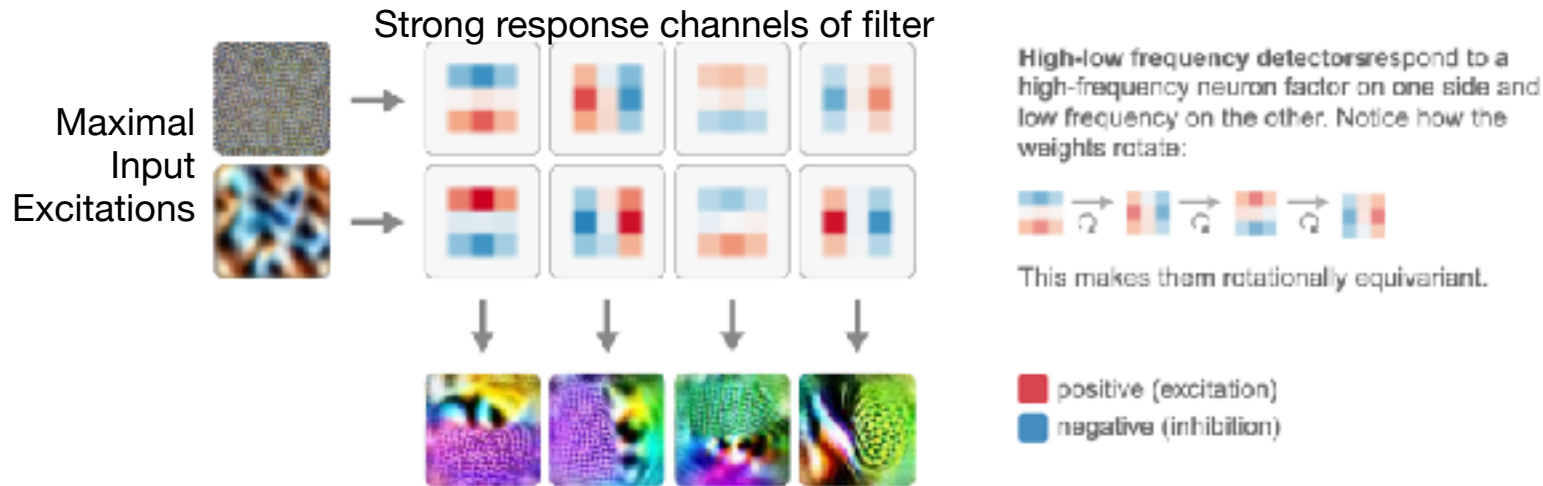# Neural Networks and Machine Learning

CNN Circuits
Continued

# Logistics and Agenda

- Logistics
  - Grading Update

- Agenda
  - Finish Circuits
  - Student Paper Presentation
  - Next Time (or today, if time):
    - Fully Convolutional Networks

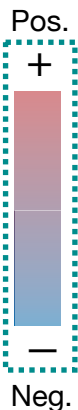# Review of Equivariant Circuits

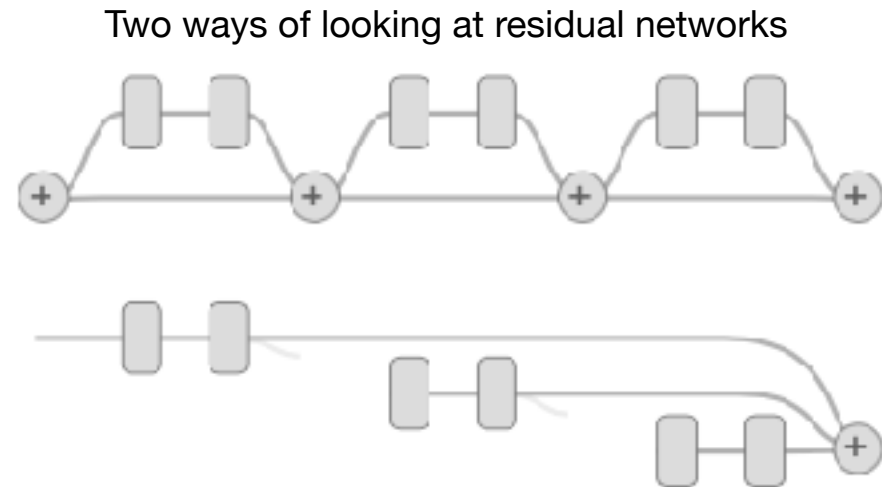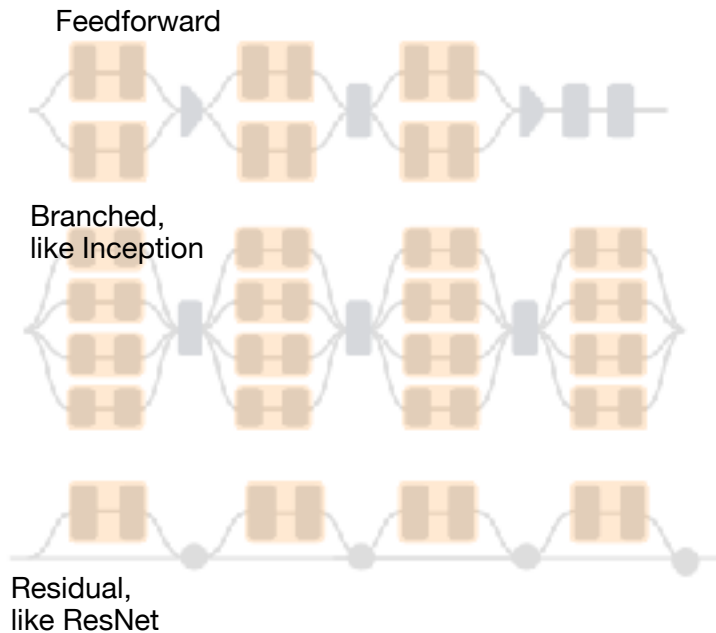- Possible to reveal patterns of circuits via sets of weights

Strong response channels of filter

Maximal Input Excitations

High-low frequency detectors respond to a high-frequency neuron factor on one side and low frequency on the other. Notice how the weights rotate:

This makes them rotationally equivariant.

positive (excitation)
negative (inhibition)

Maximal Input Excitations

Strong response channels of filter

Rotational equivariance can be turned into invariance with the transpose of an invariant -> equivariant circuit.

Here, we see color contrast units(rotationally equivariant) combine to make color center surround units(rotationally invariant). Again, notice how the weights rotate, forming the same pattern we saw above with high-low frequency detectors, but with inputs and outputs swapped.

positive (excitation)
negative (inhibition)

Pos.
+

—
Neg.

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

# Branch Specialization

Feedforward

Branched,
like Inception

Residual,
like ResNet

Two ways of looking at residual networks

Primitives

More Complex
Downstream Circuits

| | | | |
|---|---|---|---|
| A1 | | A2 | |
| B1 | | B2 | |
| C1 | | D2 | |
| D1 | | D2 | |

- Specialized branches are consistent across many architectures, support the idea of an interconnected graph of operations
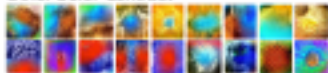
Voss, et al., "Branch Specialization", Distill, 2021.

51

Lecture Notes for CS8321 Neural Networks and Machine Learning    |    Professor Eric C. Larson    |

# 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.

BW vs Color
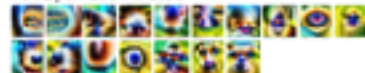
Other Color Contrast

Brightness

Other

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 detection, many of which share sub-components with curves.

Curve Related

Fur/Eye/Face Related

Boundary Detectors

Other

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.
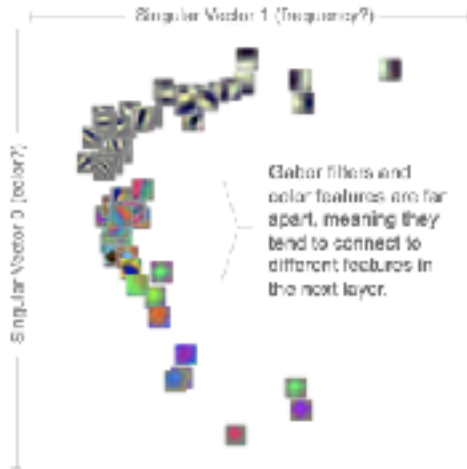
3D Geometry / Complex Shapes

Other

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

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

# Investigating Connection Clusters via SVD



Neurons in the first convolutional layer organized by the left singular vectors of |W|.
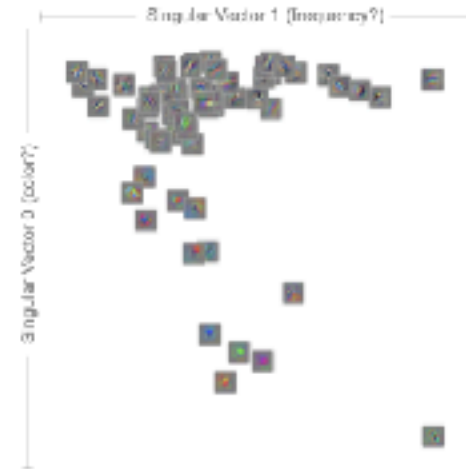
Neurons in the second convolutional layer organized by the right singular vectors of |W|.
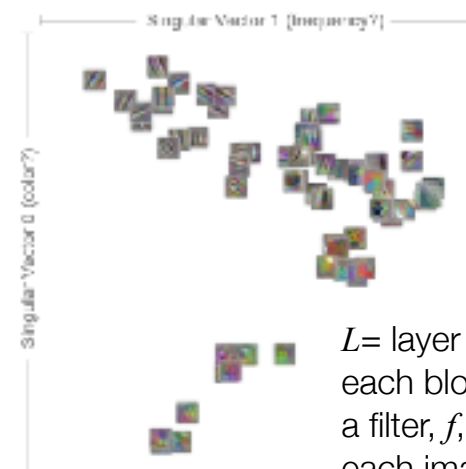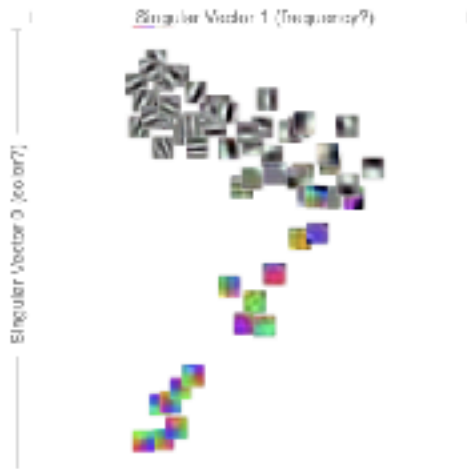
**InceptionV1 (tf-slim version) trained on ImageNet.**

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.

Gabor filters and color features are far apart, meaning they tend to connect to different features in the next layer.

**InceptionV1 trained on Places365**

Once 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 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
    $UU^T=I$   $VV^T=I$
  - $\Sigma$ is a diagonal matrix of the singular values
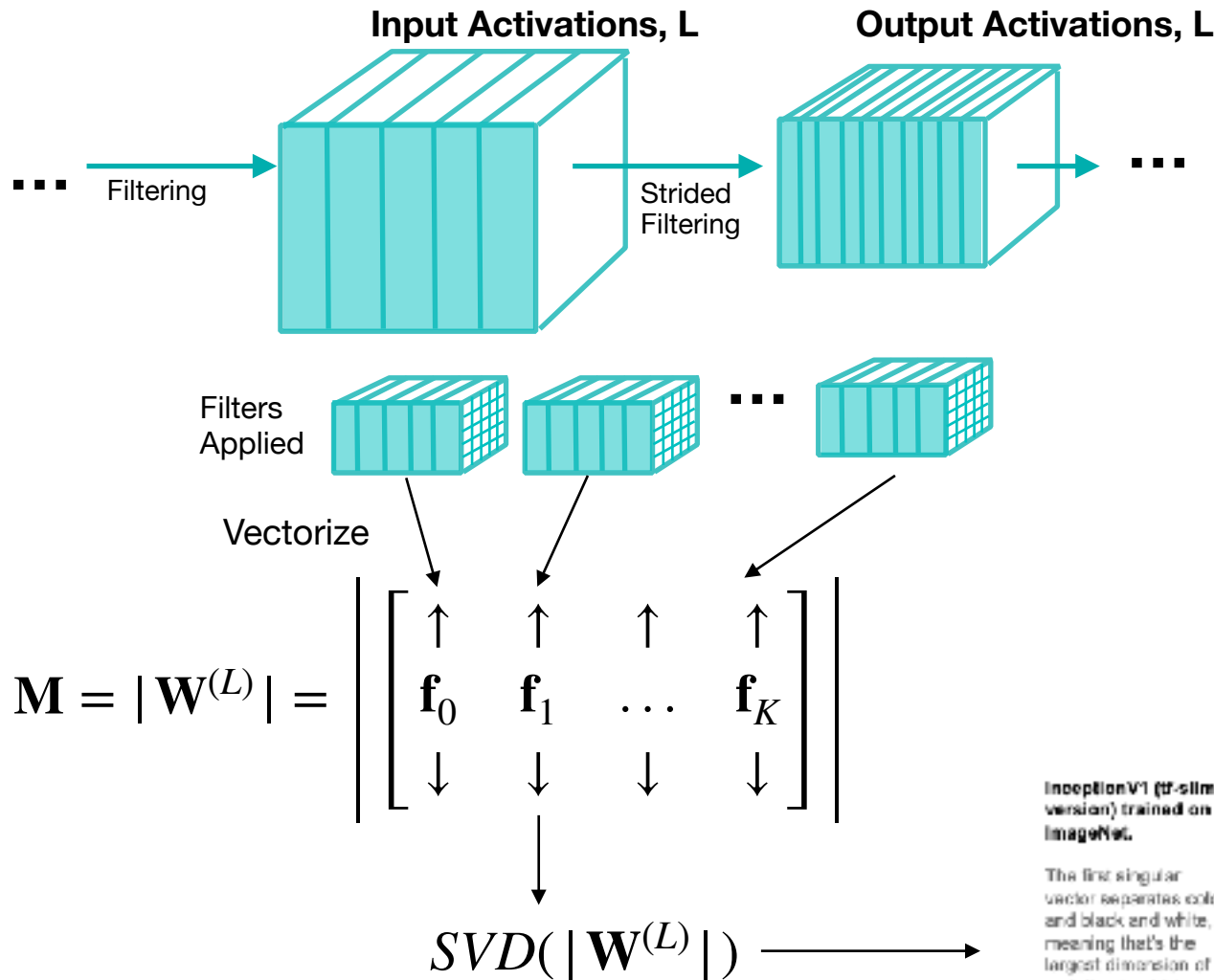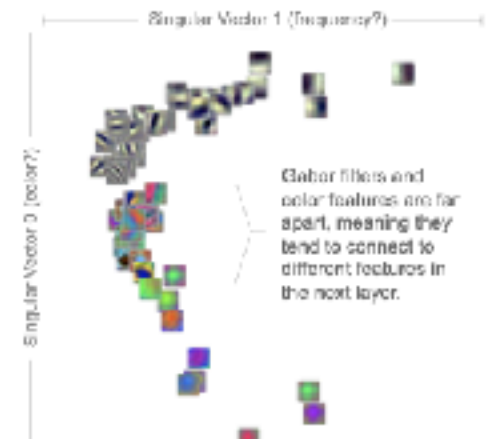  - These values characterize the variability in a matrix

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

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

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

53

**Structure of Each Tensor**:
Channels x Rows x Columns

**Input Activations, L**

**Output Activations, L**

Filtering

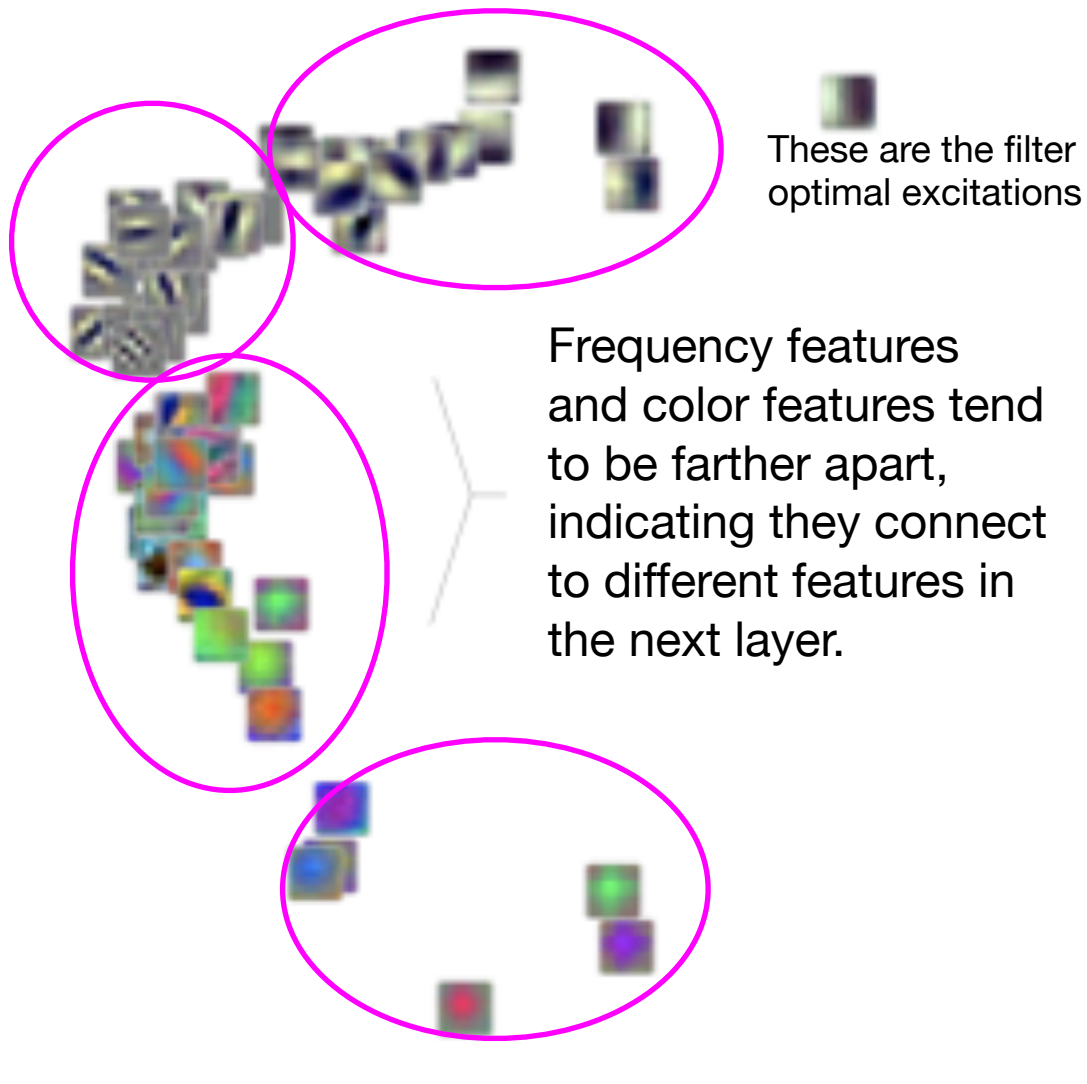Strided Filtering

$\cdots$

Filters Applied

$\cdots$

Vectorize

$$\mathbf{M} = |\mathbf{W}^{(L)}| = \left|\left|\begin{bmatrix} \uparrow & \uparrow & \uparrow & \uparrow \\ \mathbf{f}_0 & \mathbf{f}_1 & \ldots & \mathbf{f}_K \\ \downarrow & \downarrow & \downarrow & \downarrow \end{bmatrix}\right|\right|$$
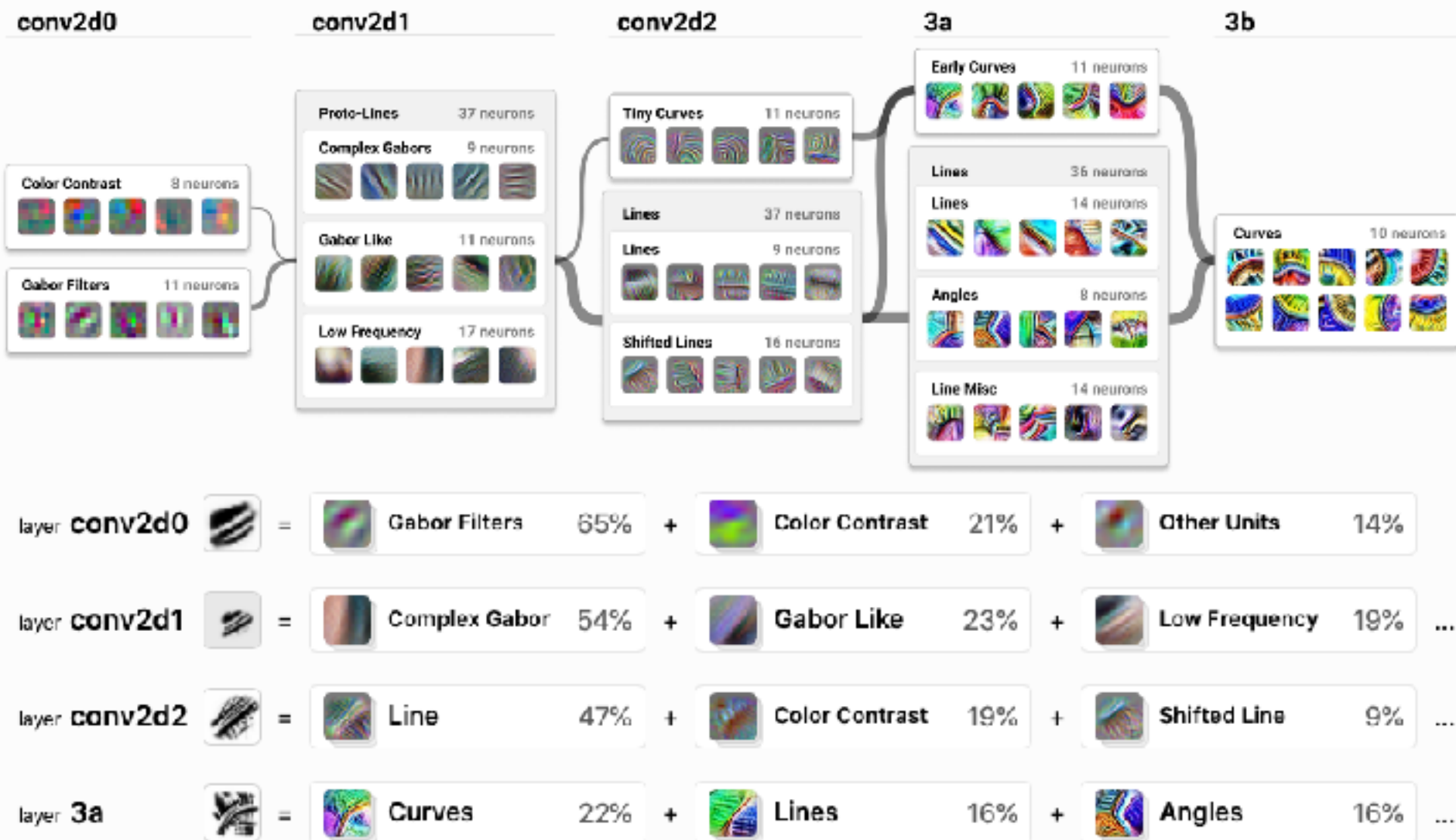
$$SVD(|\mathbf{W}^{(L)}|) \longrightarrow$$

- 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 $UU^T=I$   $VV^T=I$
  ◦ $\Sigma$ 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].

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.

Singular Vector 1 (frequency?)

Gabor filters and color features are far apart, meaning they tend to connect to different features in the next layer.

Singular Vector 3 (color?)

54

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



Singular Vector 1 (Frequency?)

These are the filter optimal excitations

- $M = U\Sigma V^T$
- $U$ is eig-vec of $MM^T$

Singular Vector 0 (color?)

Frequency features and color features tend to be farther apart, indicating they connect to different features in the next layer.

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

$$\left| \begin{bmatrix} \uparrow & \uparrow & \uparrow & \uparrow \\ \mathbf{f}_0 & \mathbf{f}_1 & \dots & \mathbf{f}_K \\ \downarrow & \downarrow & \downarrow & \downarrow \end{bmatrix} \right|$$

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

# 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

# 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

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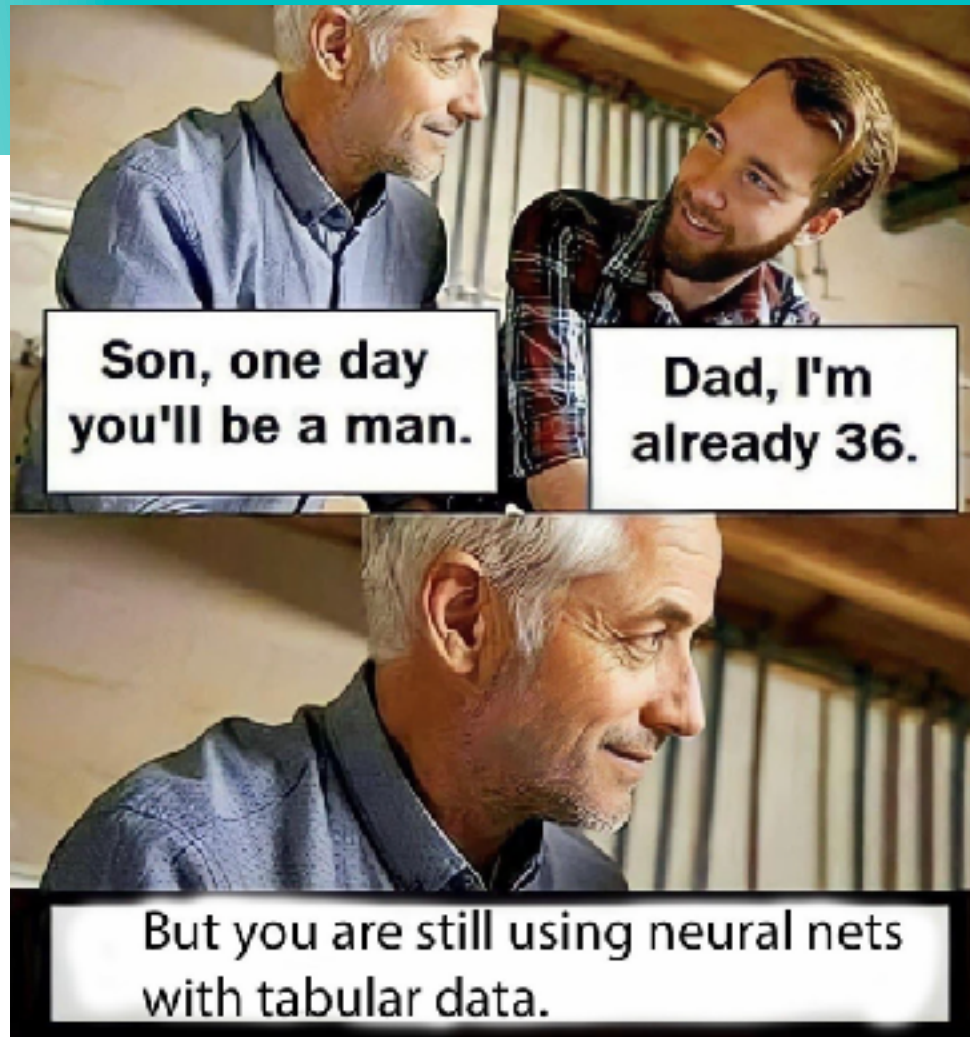
https://distill.pub/2020/circuits/zoom-in/

# Paper Presentation



Toy Models of Superposition

Lecture Notes for

# Neural Networks
# and Machine Learning

## CNN Circuits

**Next Time:**
Fully Convolutional Learning

**Reading:** Chollet 5.4