

6.S979

# Model Interpretability

---

Arvind Satyanarayan

# What is "Interpretability"?

**Causal:** the degree to which a person can understand the cause of the result.

**Post-Hoc:** what does the model tell me?

---

## The Mythos of Model Interpretability

---

Zachary C. Lipton<sup>1</sup>

### Abstract

Supervised machine learning models boast remarkable predictive capabilities. But can you trust your model? Will it work in deployment? What else can it tell you about the world? We want models to be not only good, but interpretable. And yet the task of *interpretation* appears underspecified. Papers provide diverse and

no one has managed to set it in writing, or (ii) the term interpretability is ill-defined, and thus claims regarding interpretability of various models may exhibit a quasi-scientific character. Our investigation of the literature suggests the latter to be the case. Both the motives for interpretability and the technical descriptions of interpretable models are diverse and occasionally discordant, suggesting that interpretability refers to more than one concept. In this paper, we seek to clarify both, suggesting that interpretability is



Session Features

Browse... New...

- Cycling Keywords
- Cycling Clubs and Organizations
- Cycling Lessons
- Bike Safety and Maintenance
- Cycling Tours and Events
- motorcycle

Inactive Features

- Shopping and Rentals

Find Items

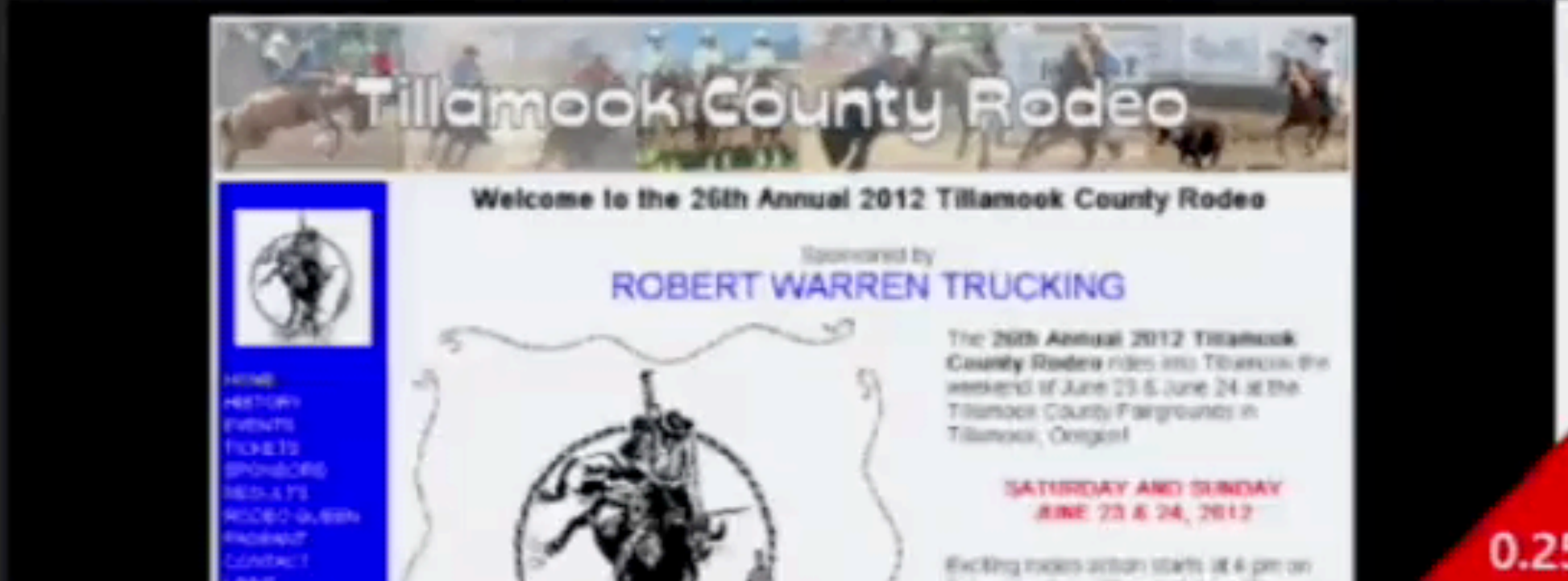
Sample 0.75

Find Diverse Find False Positives Find False Negatives

Submit

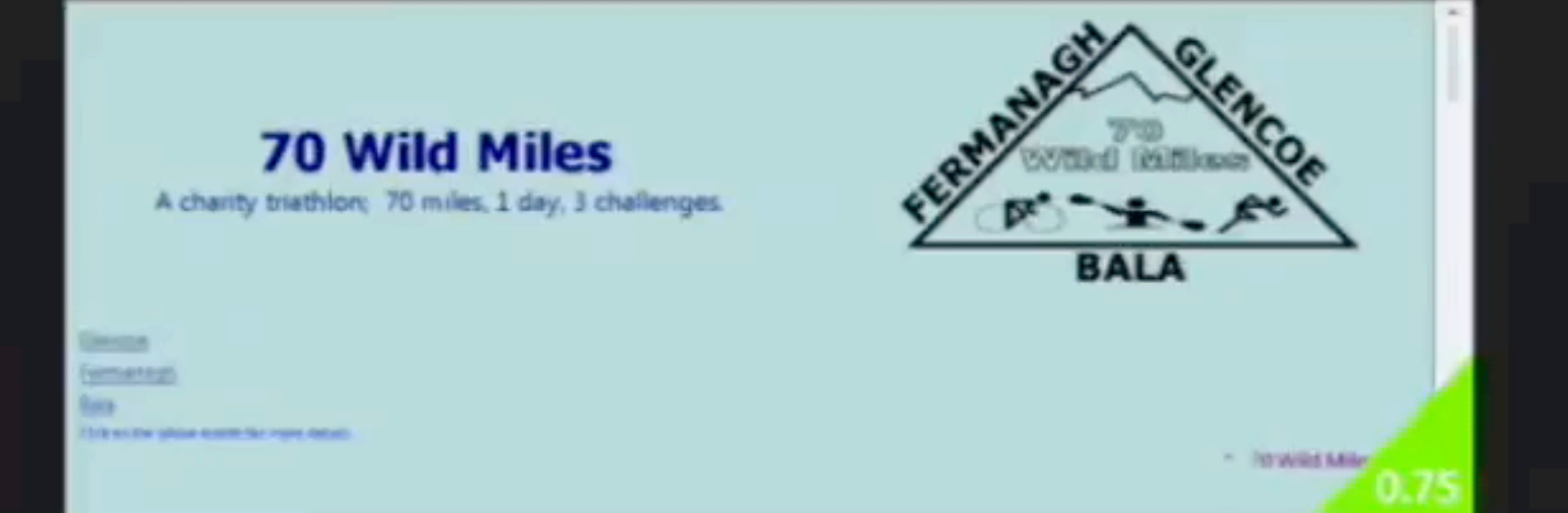
Found items around score of 0.75.

<http://www.tillamookrodeo.com>



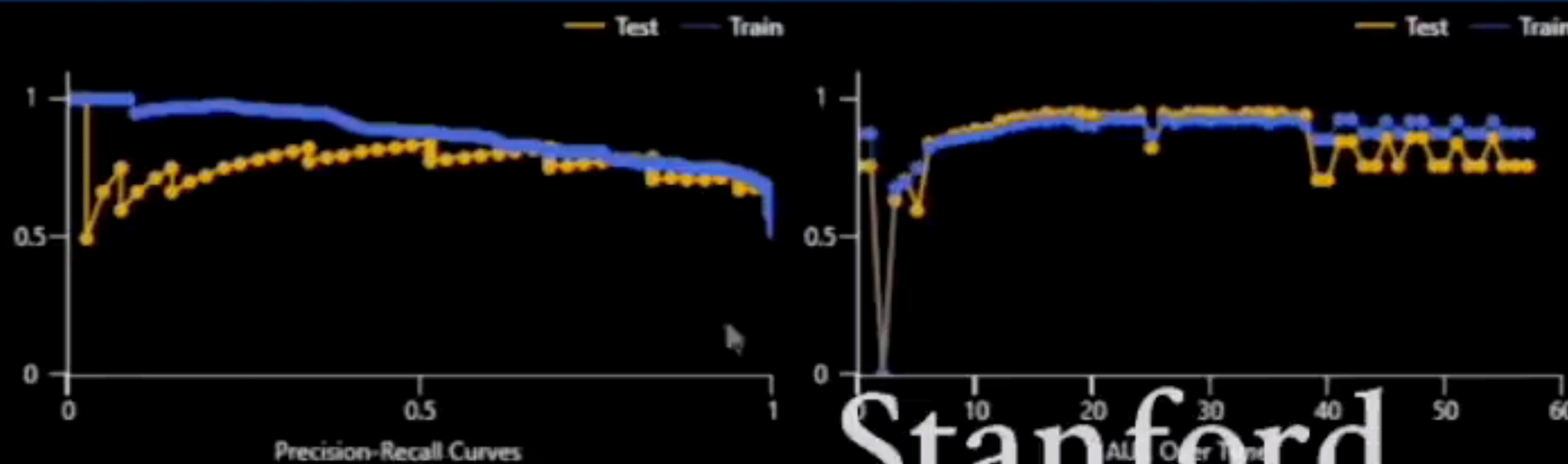
0.25

<http://www.70wildmiles.org/>



0.75

Find Items Featuring Feature Suggestions Review Labels



FEATURES: 6 ITEMS: 511

Standard Metrics ModelTracker [0.00, 1.00]

Stanford University

TEST AUC = 0.76 (+0.00) TRAIN AUC = 0.88 (+0.00)

1:15 PM 1/6/2015





Similar Items Shown on Hover

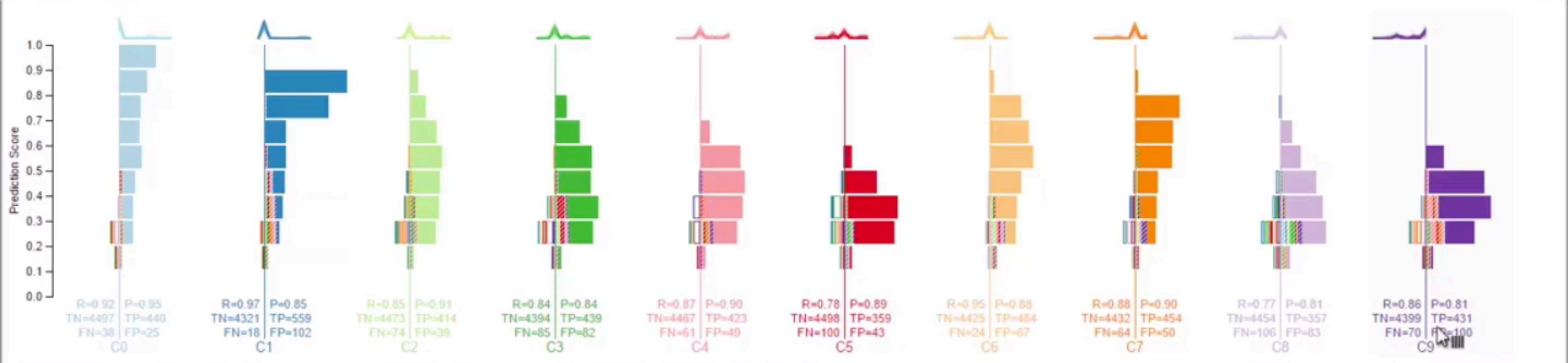
Outliers

Changes from Previous Iteration

Interactive Prediction Threshold

Tags for Tracking





Dataset: mnist\_randomforest.csv | 10 Classes, 5000 Instances, 5000 Shown | Acc.: 0.87, Prec.: 0.87, Recall: 0.87, TP: 436.0, FP: 64.0, TN: 4436.0, FN: 64.0 | 0 Hovered, 0 Selected

★	Image	TRUE	Assigned	Correct	Prediction Score	C0	C1	C2	C3	C4	C5	C6	C7
☆	7	C7	C7	1	0.5849	0.0103	0.0226	0.0646	0.0330	0.0621	0.0279	0.0383	0.5849
☆	1	C1	C1	1	0.3576	0.0111	0.3576	0.1180	0.0431	0.0450	0.0447	0.0589	0.0235
☆	7	C7	C7	1	0.6777	0.0185	0.0310	0.0485	0.0324	0.0328	0.0491	0.0152	0.6777
☆	1	C9	C9	1	0.2158	0.0146	0.0634	0.0233	0.1524	0.1435	0.1895	0.0541	0.0745
☆	5	C5	C5	1	0.2964	0.0346	0.2051	0.1039	0.1367	0.0286	0.2964	0.0606	0.0314
☆	3	C3	C3	1	0.2365	0.0121	0.0946	0.0367	0.2365	0.0617	0.1952	0.0599	0.1379
☆	2	C2	C2	1	0.3628	0.0069	0.1414	0.3628	0.1552	0.0153	0.0526	0.1246	0.0190
☆	1	C1	C1	1	0.8319	0.0017	0.8319	0.0342	0.0166	0.0087	0.0201	0.0123	0.0258
☆	2	C2	C2	1	0.7116	0.0192	0.0188	0.7116	0.0230	0.0390	0.0157	0.0597	0.0198
☆	1	C1	C1	1	0.7070	0.0014	0.7070	0.0384	0.0223	0.0234	0.0203	0.0194	0.0637
☆	7	C7	C7	1	0.7910	0.0081	0.0135	0.0240	0.0252	0.0417	0.0272	0.0079	0.7910
☆	1	C1	C1	1	0.8477	0.0014	0.8477	0.0253	0.0172	0.0068	0.0230	0.0165	0.0215
☆	6	C6	C6	1	0.3483	0.0364	0.0256	0.0358	0.0728	0.1373	0.1377	0.3483	0.0257
☆	1	C8	C8	1	0.3779	0.0115	0.1461	0.0314	0.0668	0.0879	0.0825	0.0701	0.0486
☆	2	C2	C2	1	0.3196	0.0638	0.0291	0.3196	0.1187	0.0546	0.0869	0.1316	0.0617
☆	0	C0	C0	1	0.9639	0.9639	0.0003	0.0020	0.0026	0.0008	0.0198	0.0053	0.0016
☆	4	C4	C4	1	0.5391	0.0044	0.0179	0.0161	0.0156	0.5391	0.0369	0.0449	0.1140



Fit to screen

Download PNG

Run run1

(2)

Session runs (0)

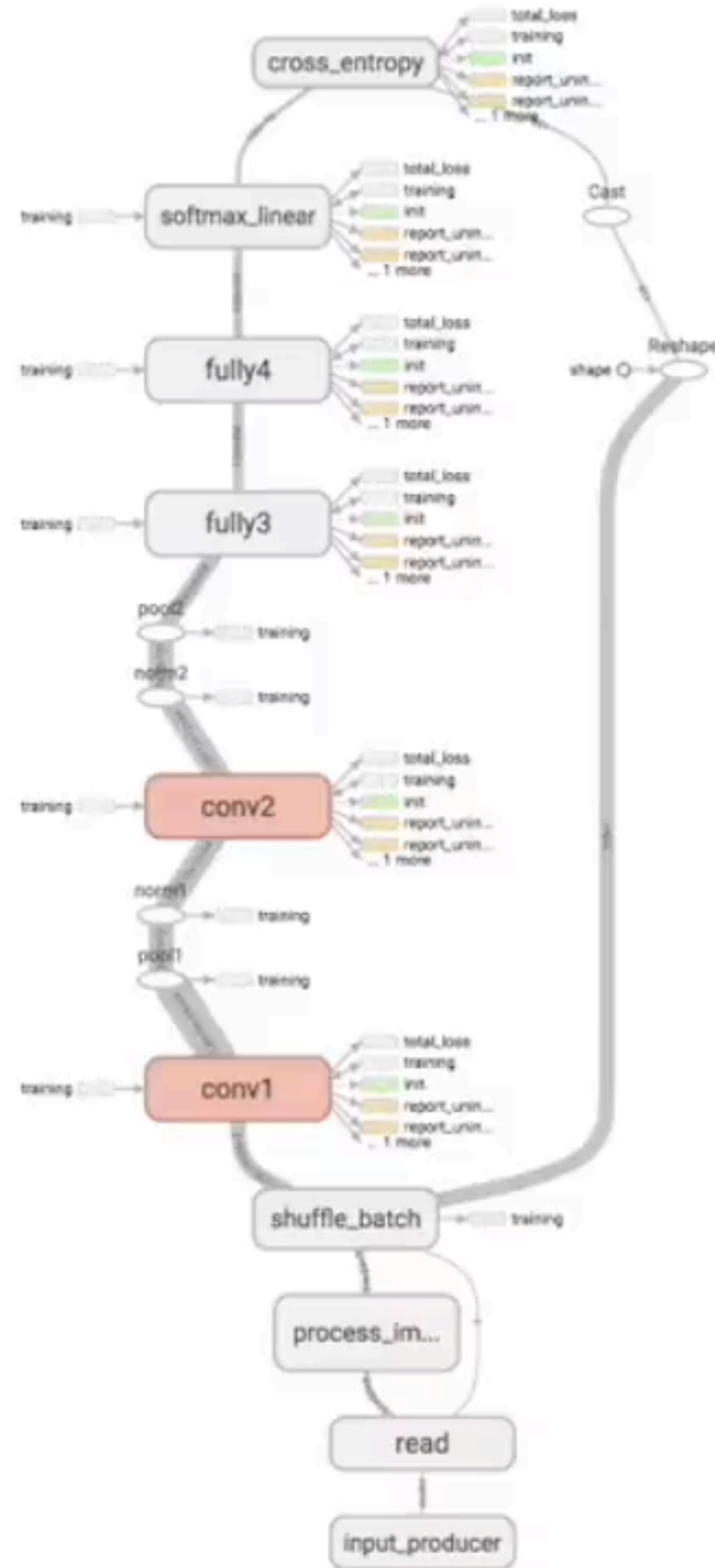
Upload

Trace inputs

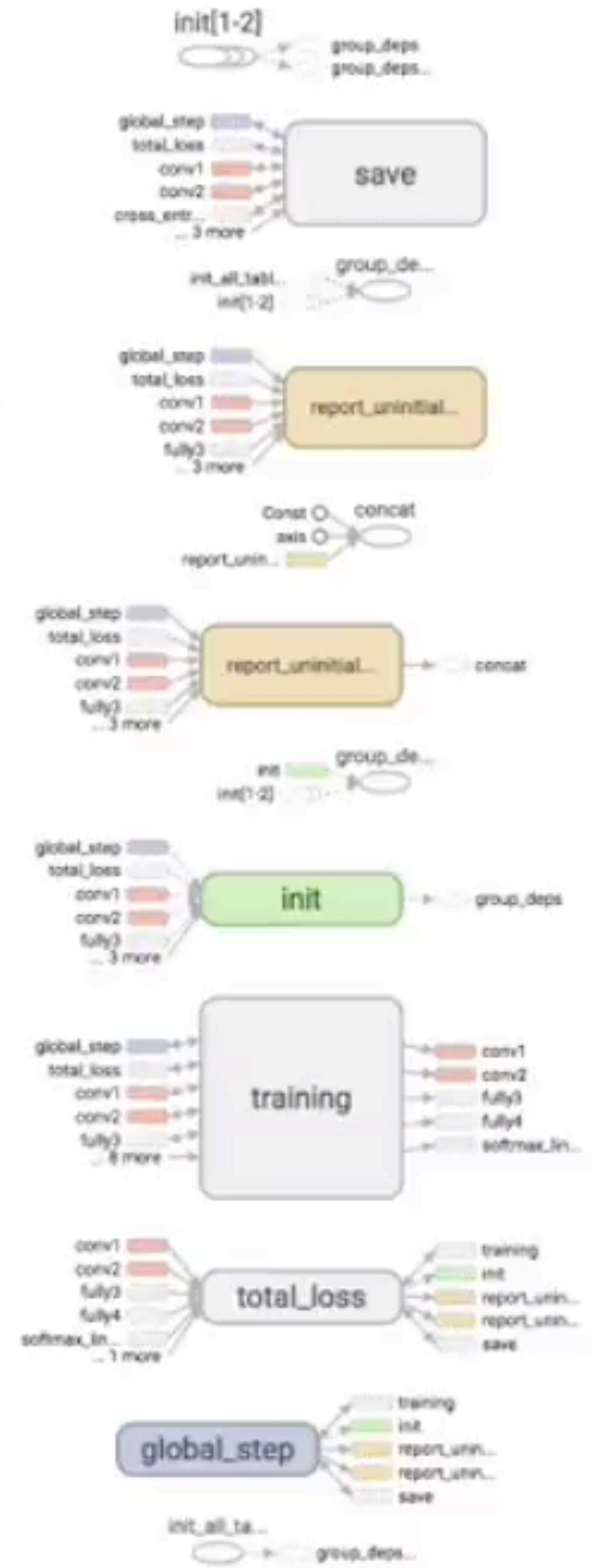
Color  Structure  Device

colors  same substructure  unique substructure

### Main Graph



### Auxiliary Nodes



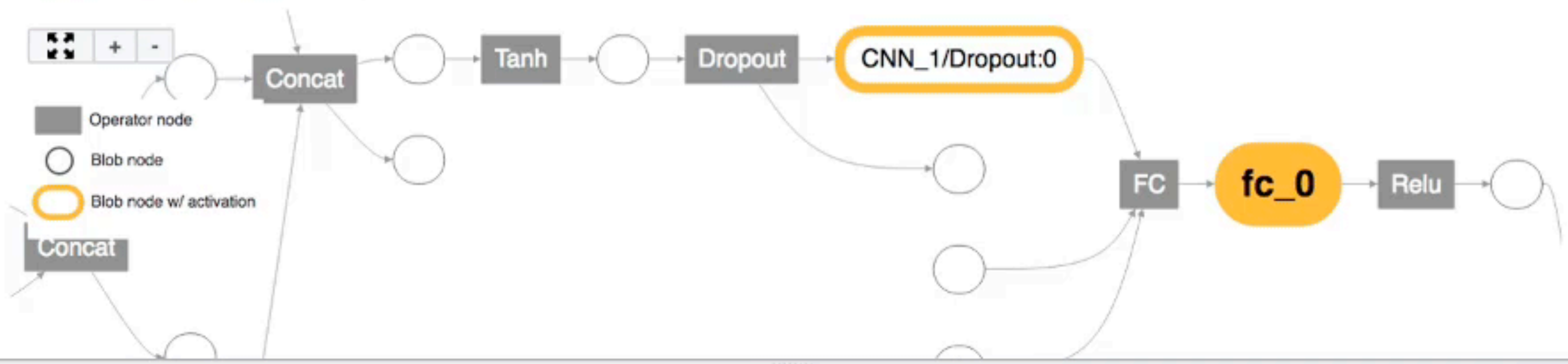
Graph (\* = expandable)

- Namespace\*
- OpNode
- Unconnected series\*
- Connected series\*
- Constant
- Summary
- Dataflow edge
- Control dependency edge
- Reference edge



# ActiVis: Visualization of Deep Neural Networks #15782570

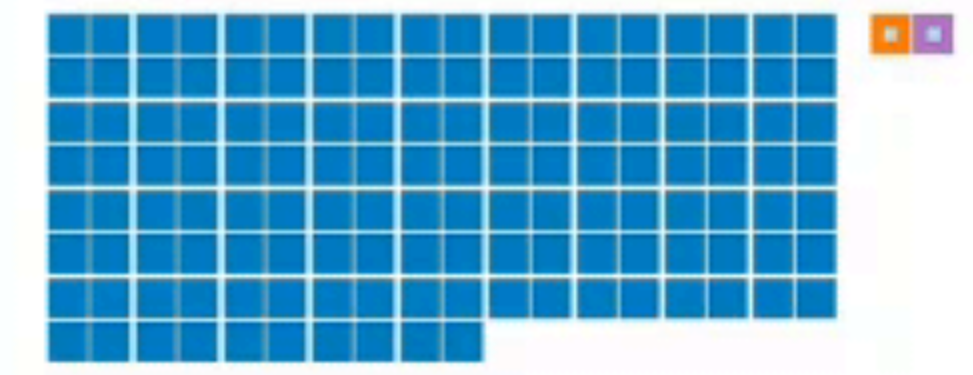
## COMPUTATION GRAPH VISUALIZATION



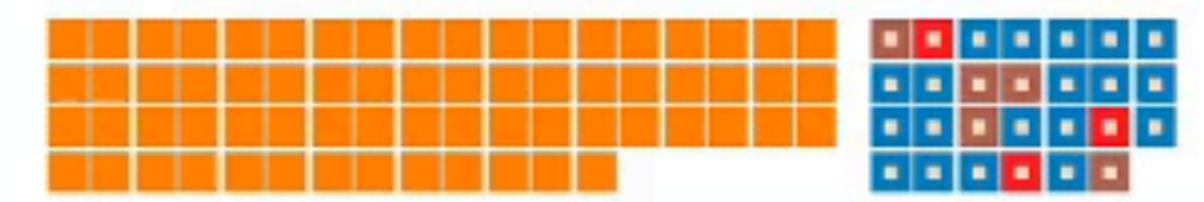
## INSTANCE SELECTION

Left column shows correctly classified instances. Right column shows misclassified instances, with border colors indicating predicted classes.

### DESC



### ENTY



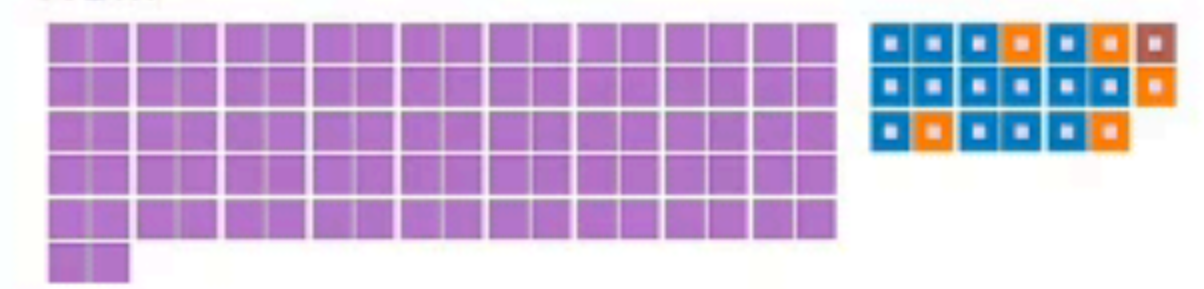
### ABBR



### HUM



### NUM



### LOC

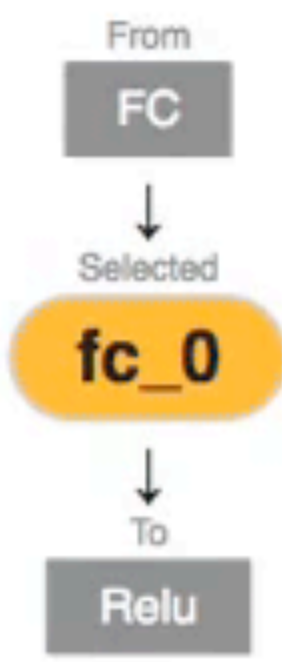


## NEURON ACTIVATION

Each row represents a group of instances. Each column is a neuron. Columns sorted by activation strength for

By class	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22
DESC		●					●	●							●						●	●	
ENTY		●					●	●							●						●	●	
ABBR																							
HUM									●	●	●	●	●								●	●	
NUM		●	●						●												●	●	
LOC							●	●													●	●	
By user-defined filters																							
Contain 'Where'							●	●													●	●	
Contain 'located'																							
Contain 'How many'		●	●																		●	●	
Contain 'How'		●	●																		●	●	
By instance ID																							

## PROJECTED



# Seq2Seq-Vis: A visual debugging tool for sequence-to-sequence models

H. Strobelt, S. Gehrmann, M. Behrisch, A. Perer, H. Pfister, A. Rush

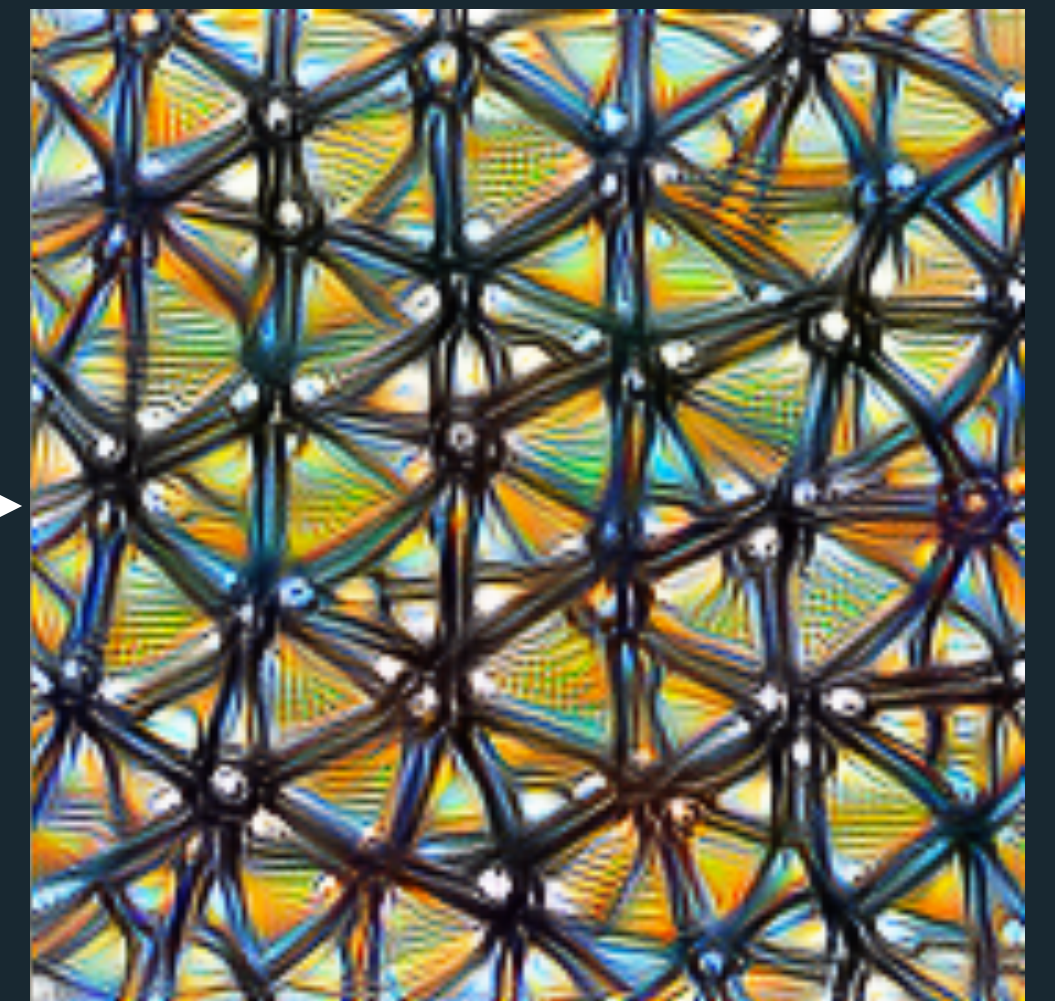
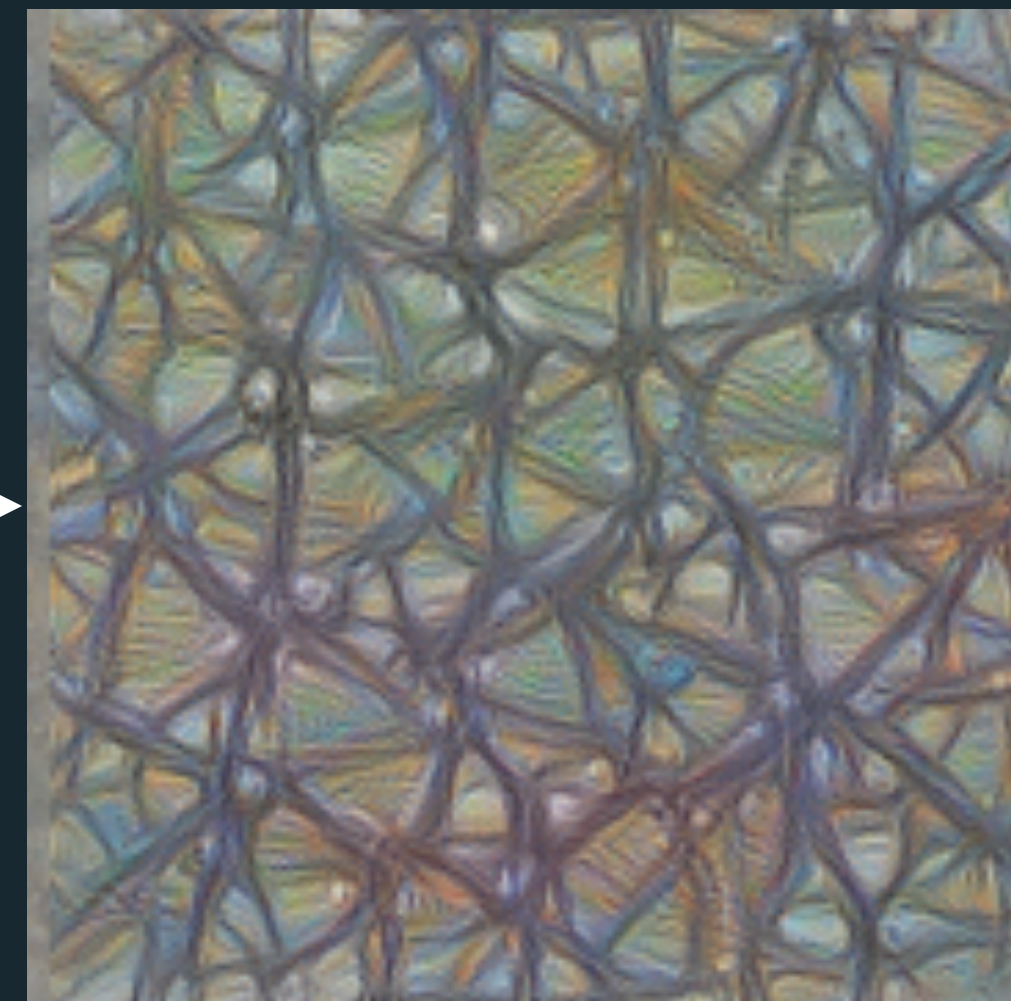
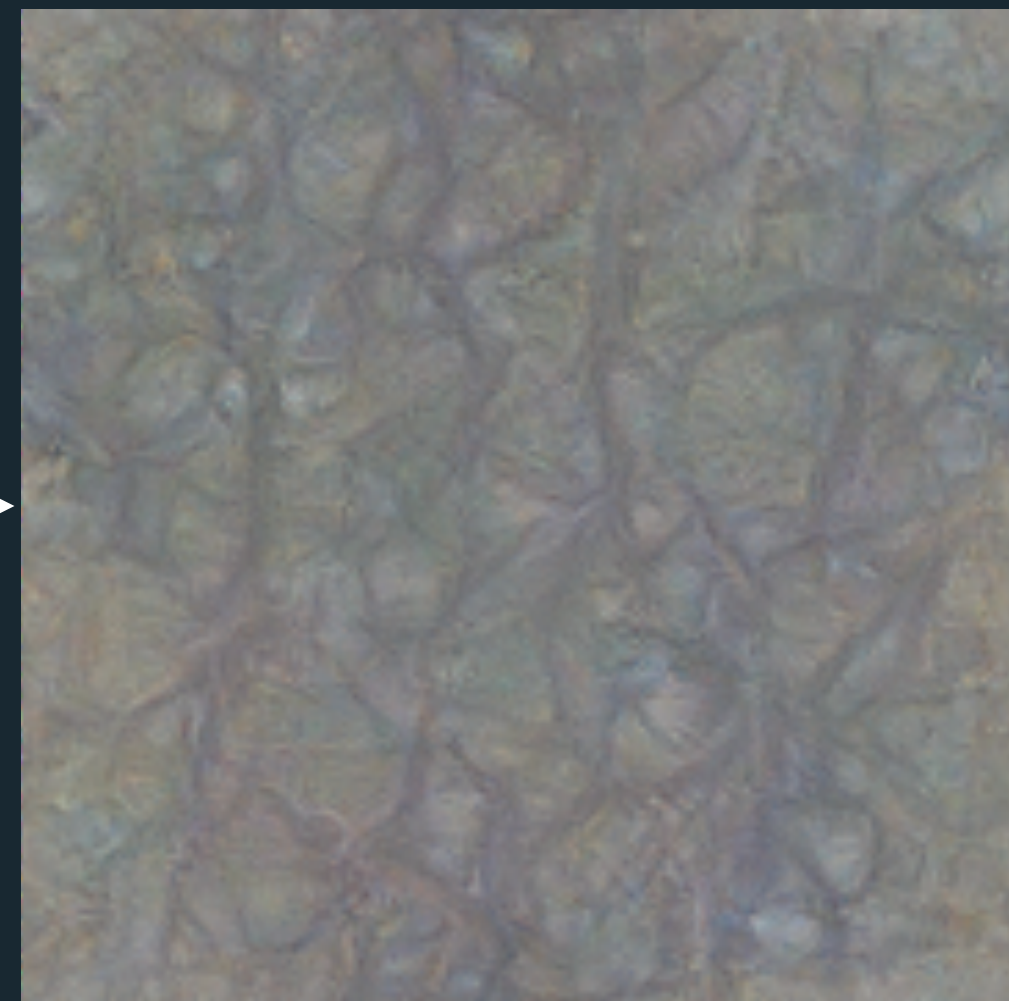
**IBM** Research





# Feature Visualization

Olah, Mordvintsev, and Schubert. Distill, 2017.  
<https://distill.pub/2017/feature-visualization/>



Step 0

Step 4

Step 48

Step 2,048

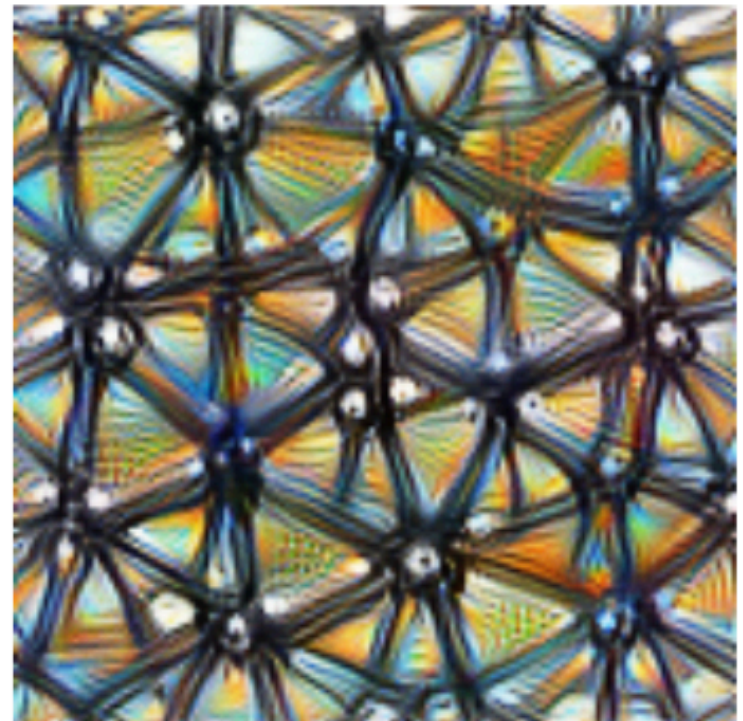


Different **optimization objectives** show what different parts of a network are looking for.

- n** layer index
- x, y** spatial position
- z** channel index
- k** class index



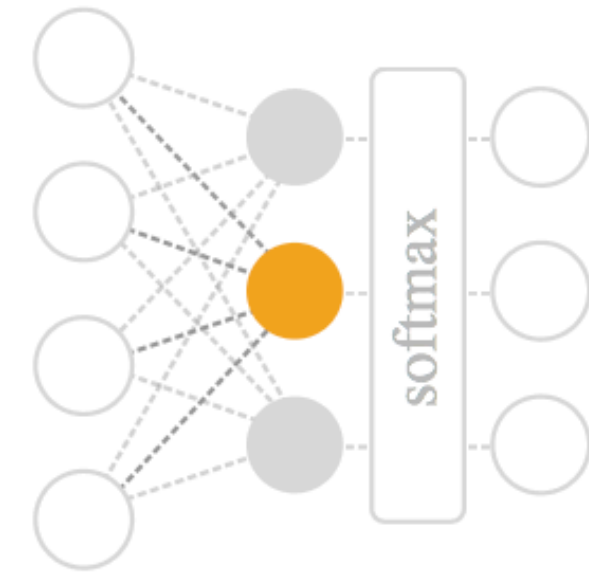
**Neuron**  
`layern[x, y, z]`



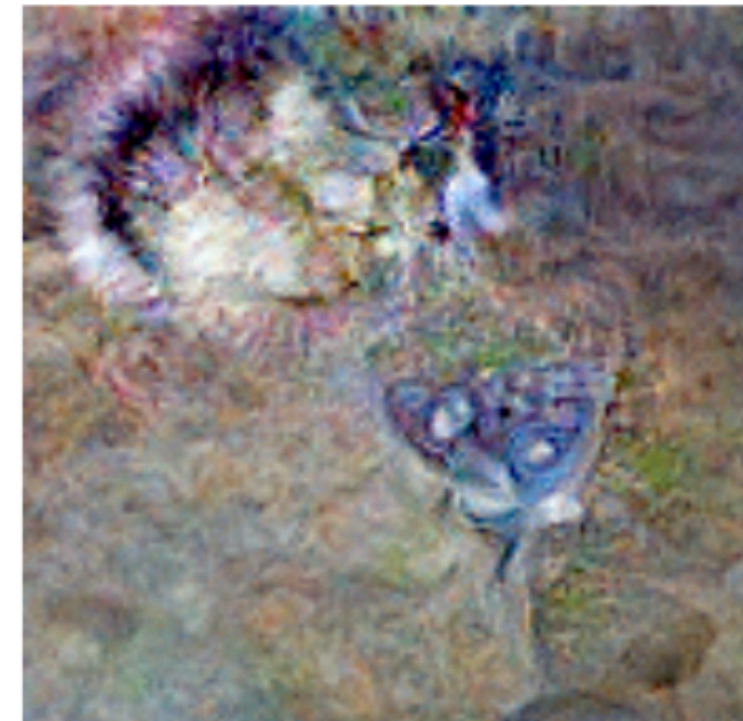
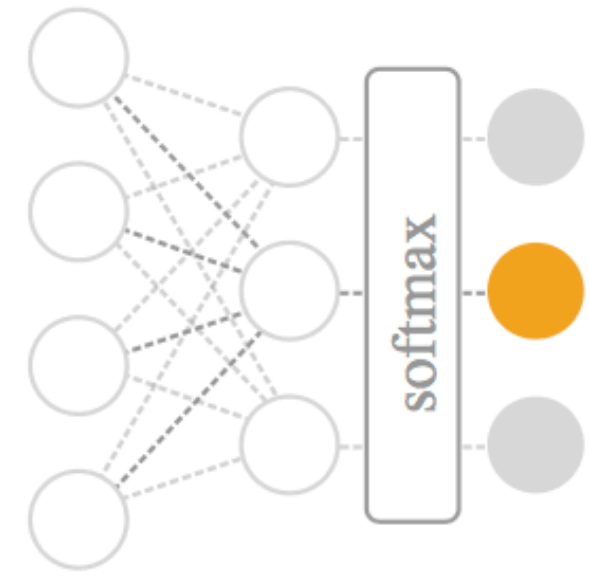
**Channel**  
`layern[:, :, z]`



**Layer/DeepDream**  
`layern[:, :, :]2`



**Class Logits**  
`pre_softmax[k]`

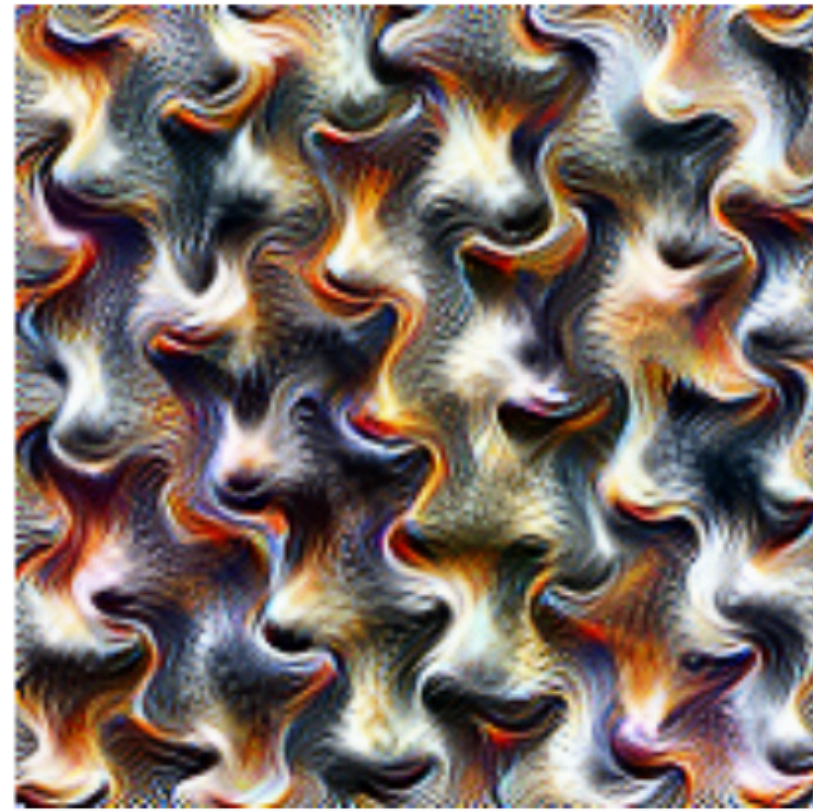
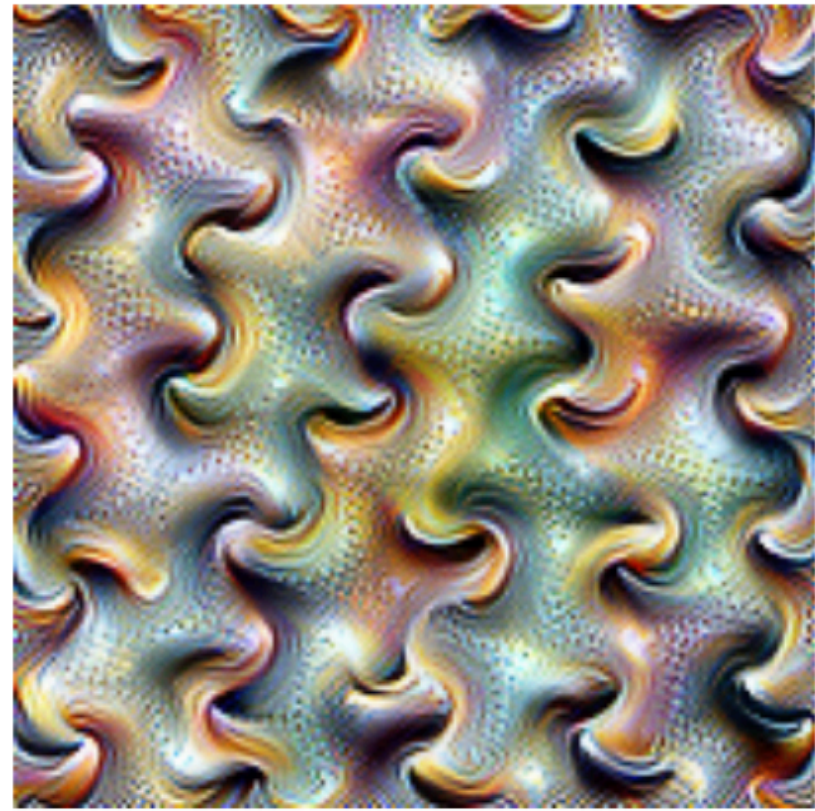


**Class Probability**  
`softmax[k]`





Simple Optimization



Optimization with diversity reveals four different, curvy facets. *Layer mixed4a, Unit 97*



Dataset examples



Simple Optimization




Optimization with diversity reveals multiple types of balls. *Layer mixed5a, Unit 9*

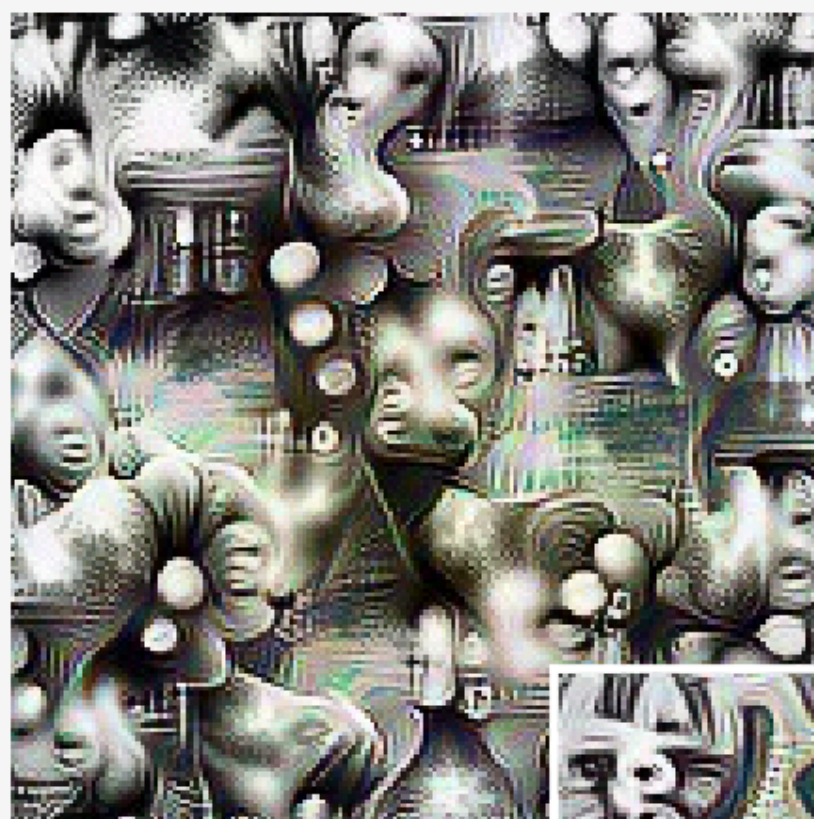
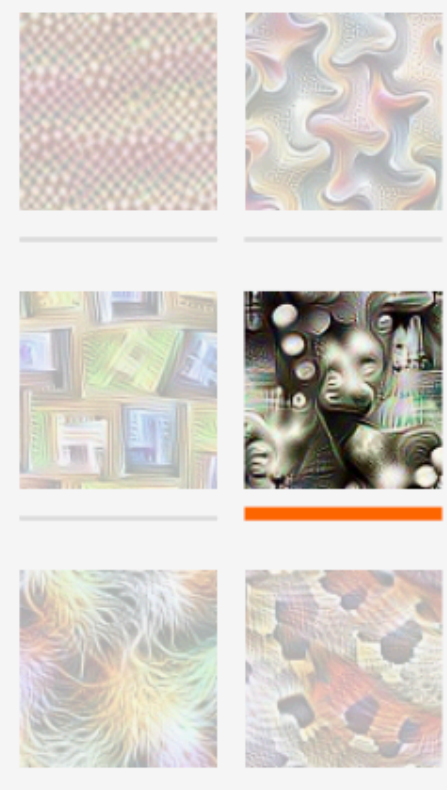


Dataset examples



By jointly optimizing two neurons we can get a sense of how they interact.

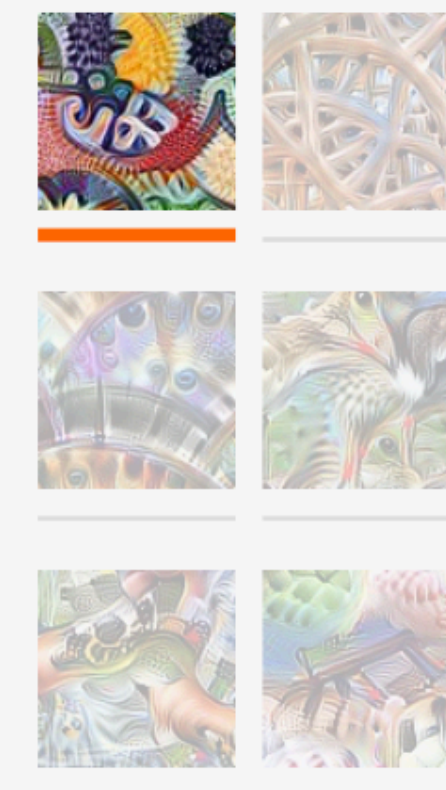
REPRODUCE IN A  
 NOTEBOOK



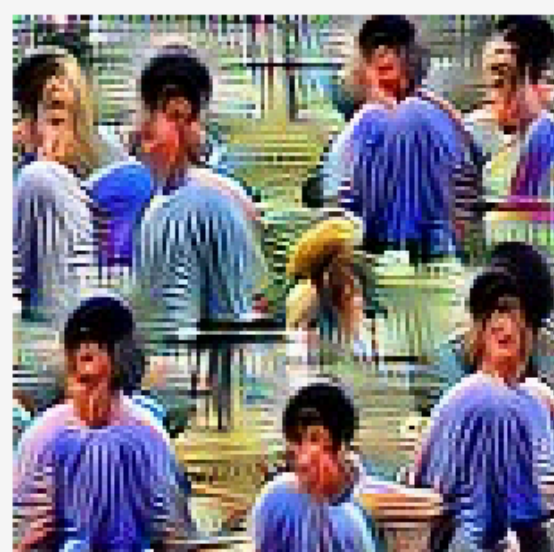
Neuron 1



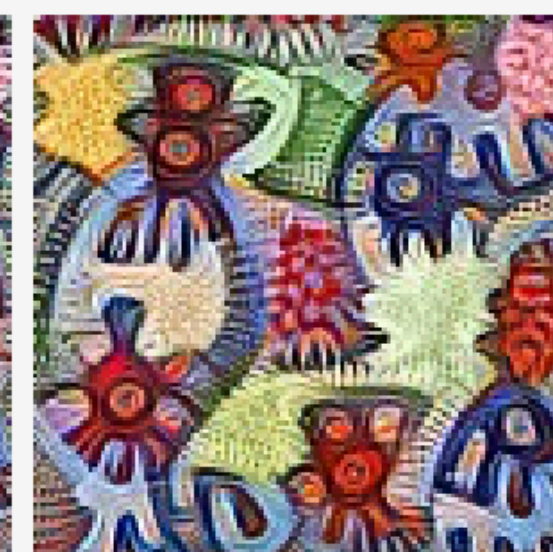
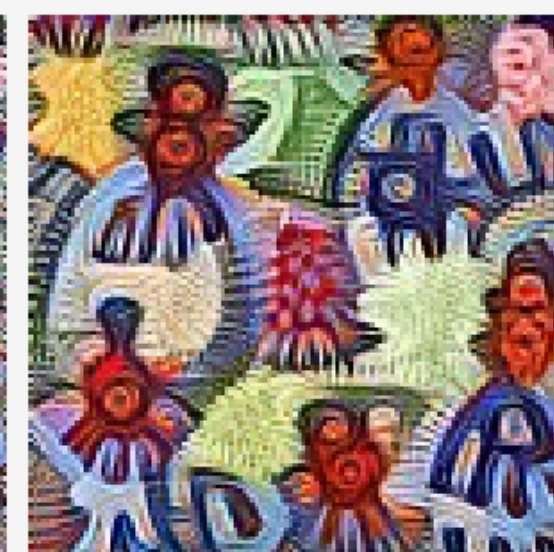
Neuron 2



Jointly optimized



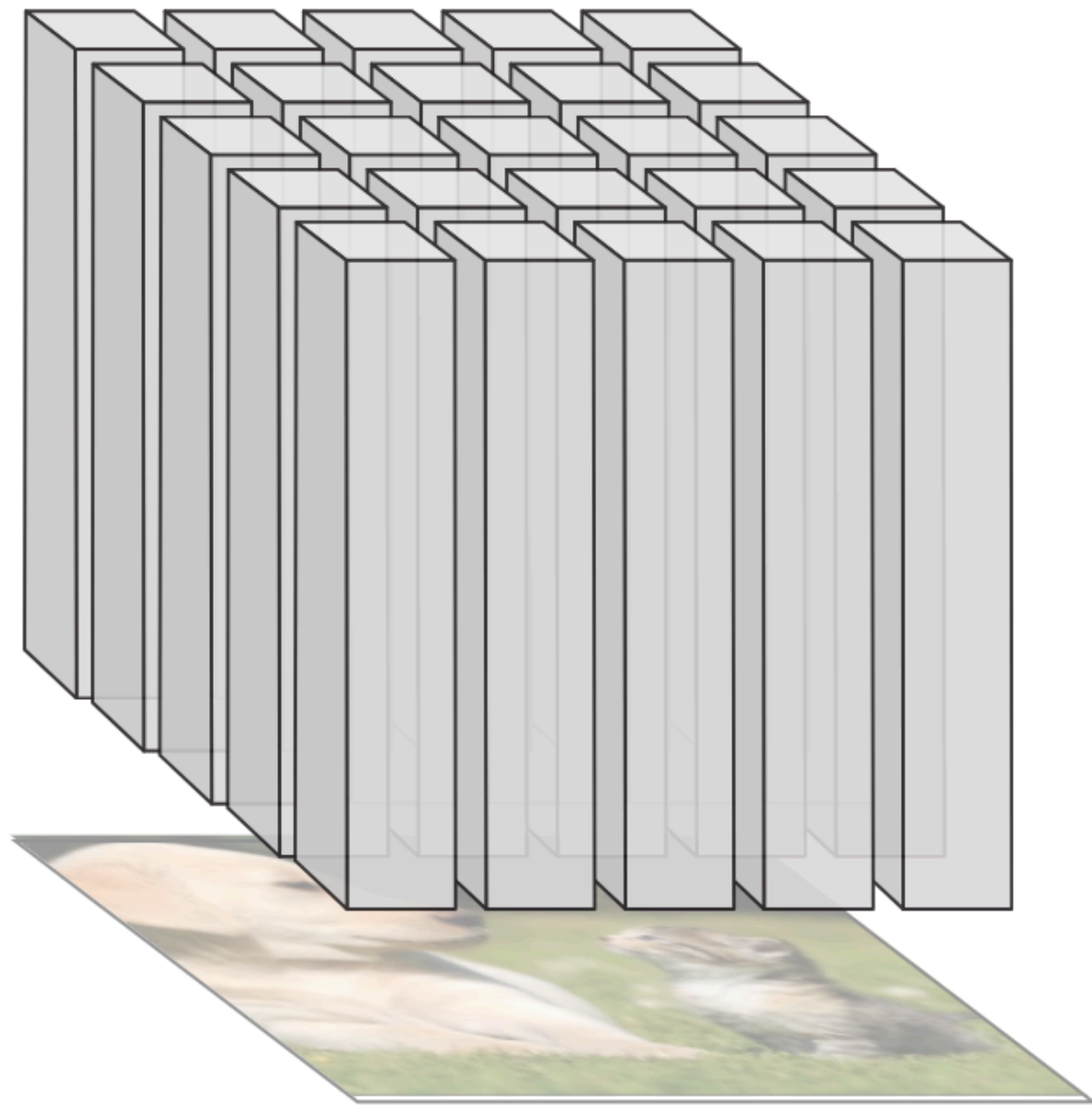
Layer 4b, Unit 475



Layer 4a, Unit 476



## Spatial Activations



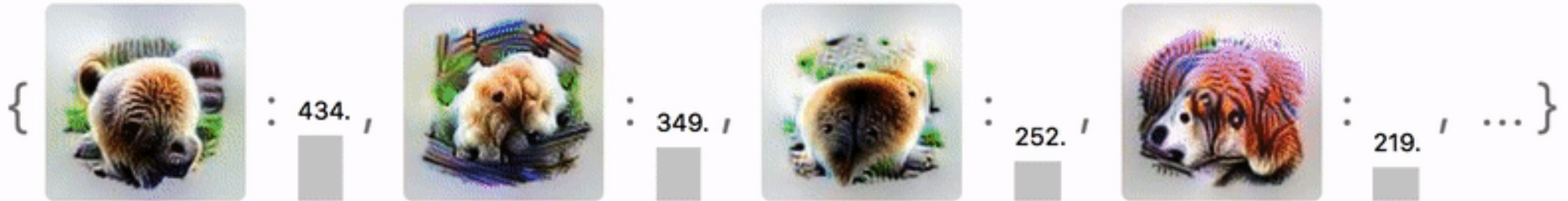
$a_{1,0} = [0, 0, 0, 0, 49.6, 0, 43.6, 30.2, 119.8, 62.7, 0, 51...$



# Semantic Dictionaries

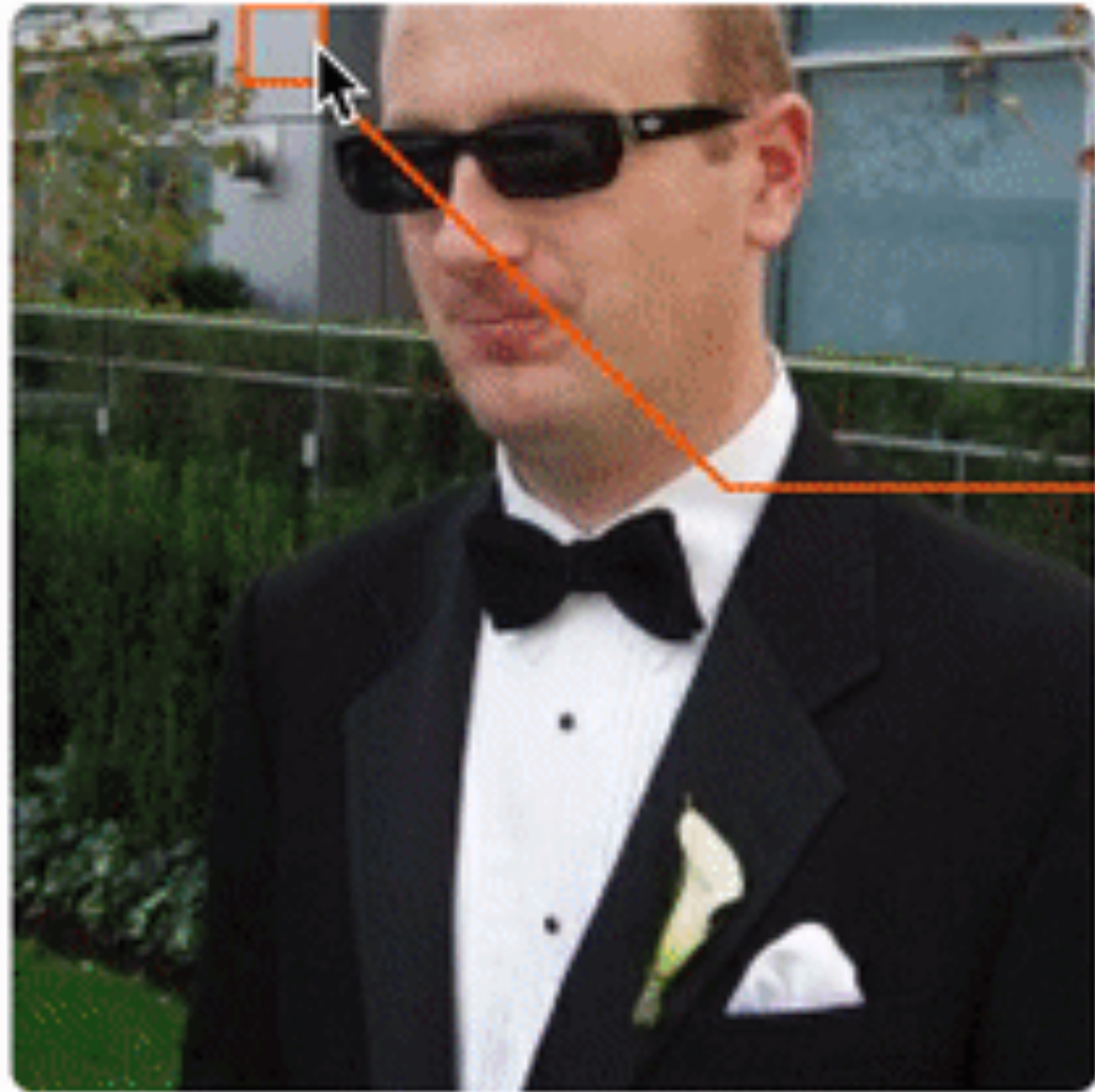


$a_{1,0} = [0, 0, 0, 0, 49.6, 0, 43.6, 30.2, 119.8, 62.7, 0, 51...$





# Semantic Dictionaries

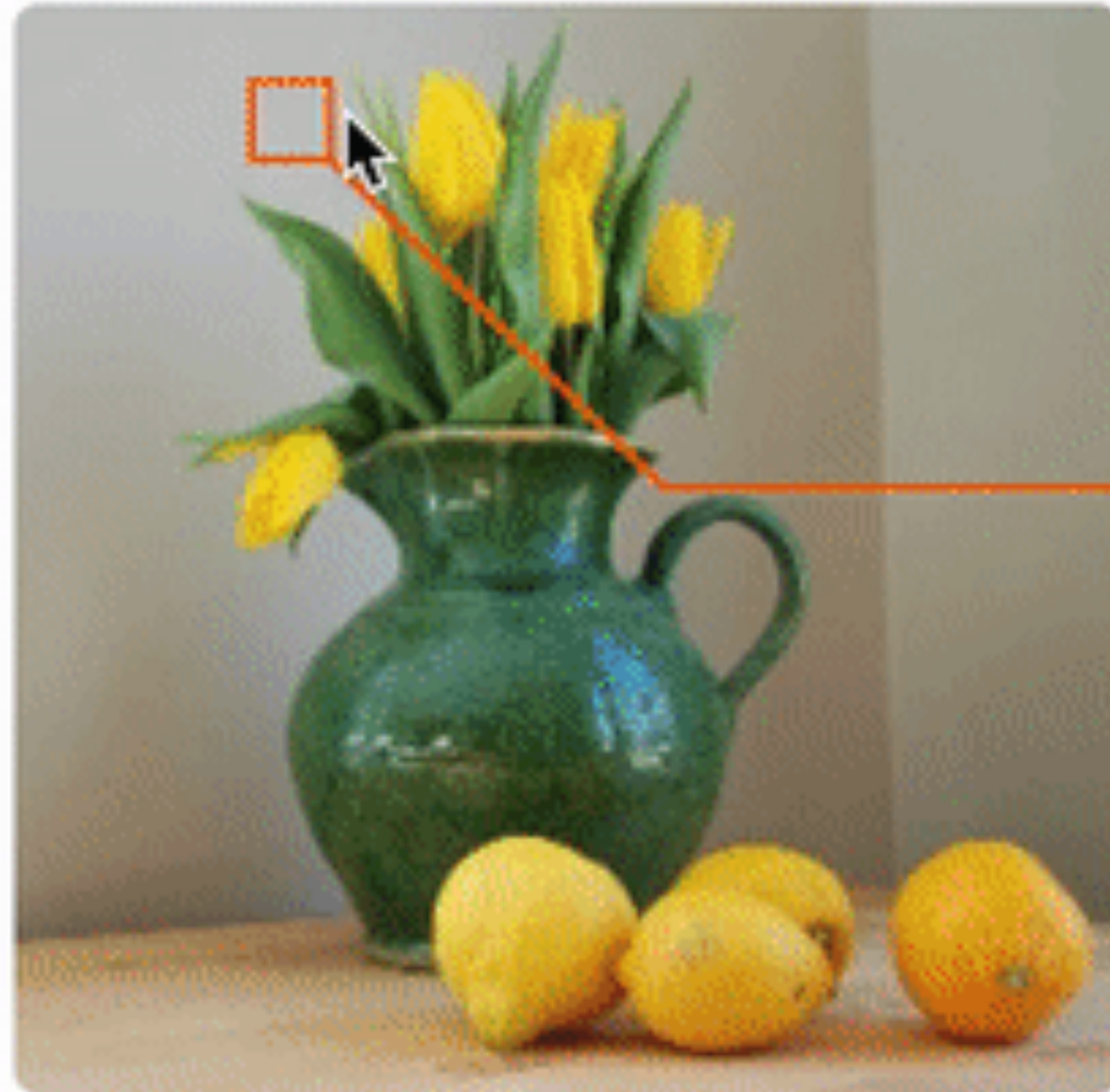


$a_{0,3} = [0, 0, 0, 76.8, 0, 38.5, 0, 0, 15.1, 0, 0, 10.4, \dots]$





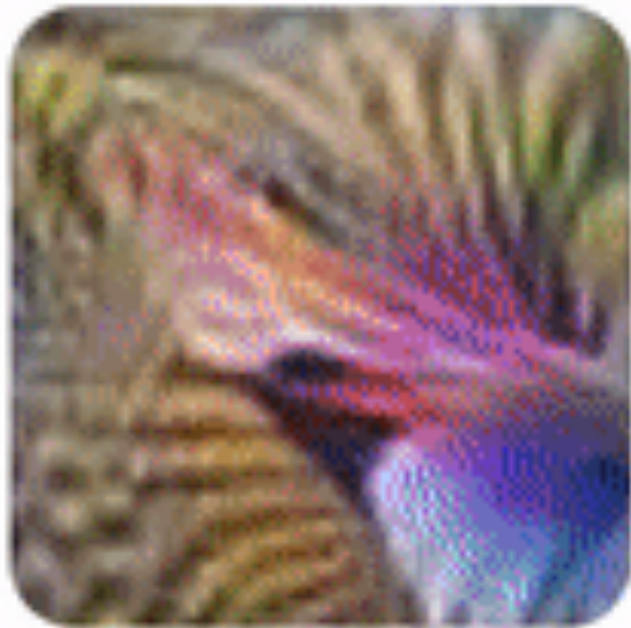
# Semantic Dictionaries



$a_{1,3} = [0, 0, 7.58, 48.4, 10.8, 0, 0, 0, 0, 0, 52.5, 0, \dots]$

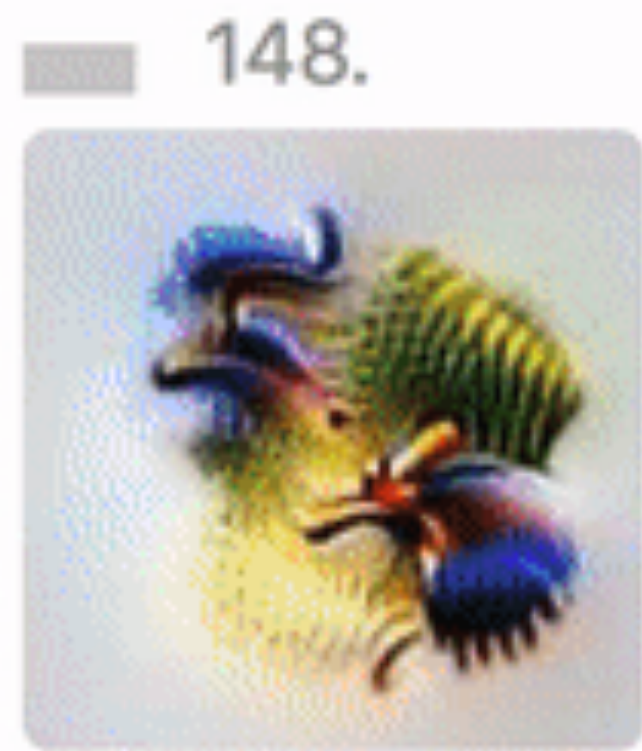




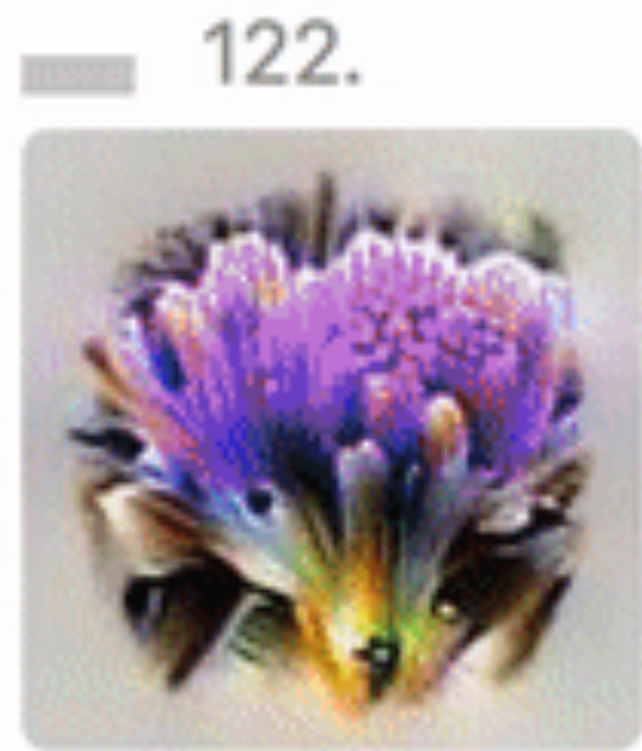


Activation Vector

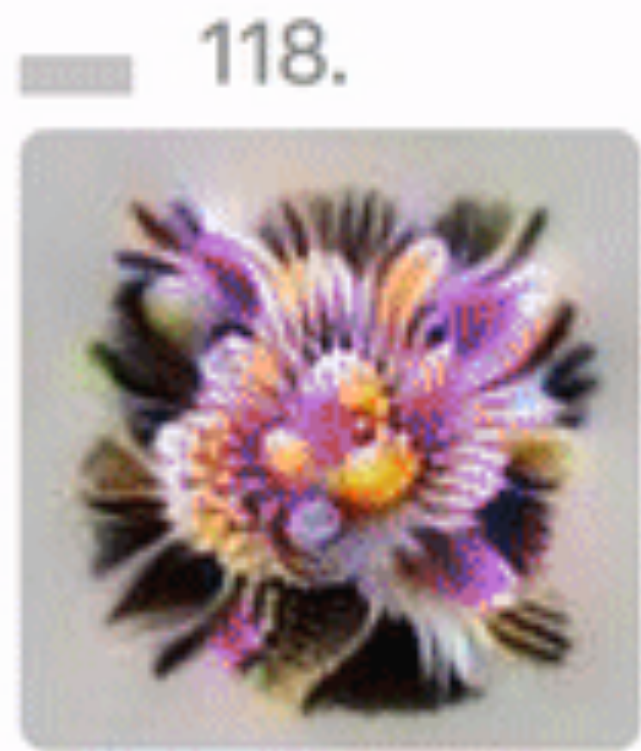
=



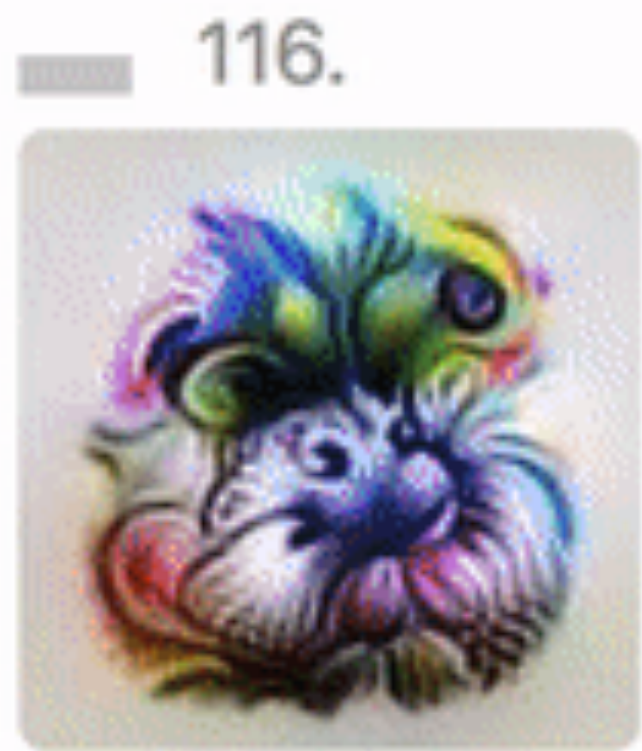
+



+



+



+

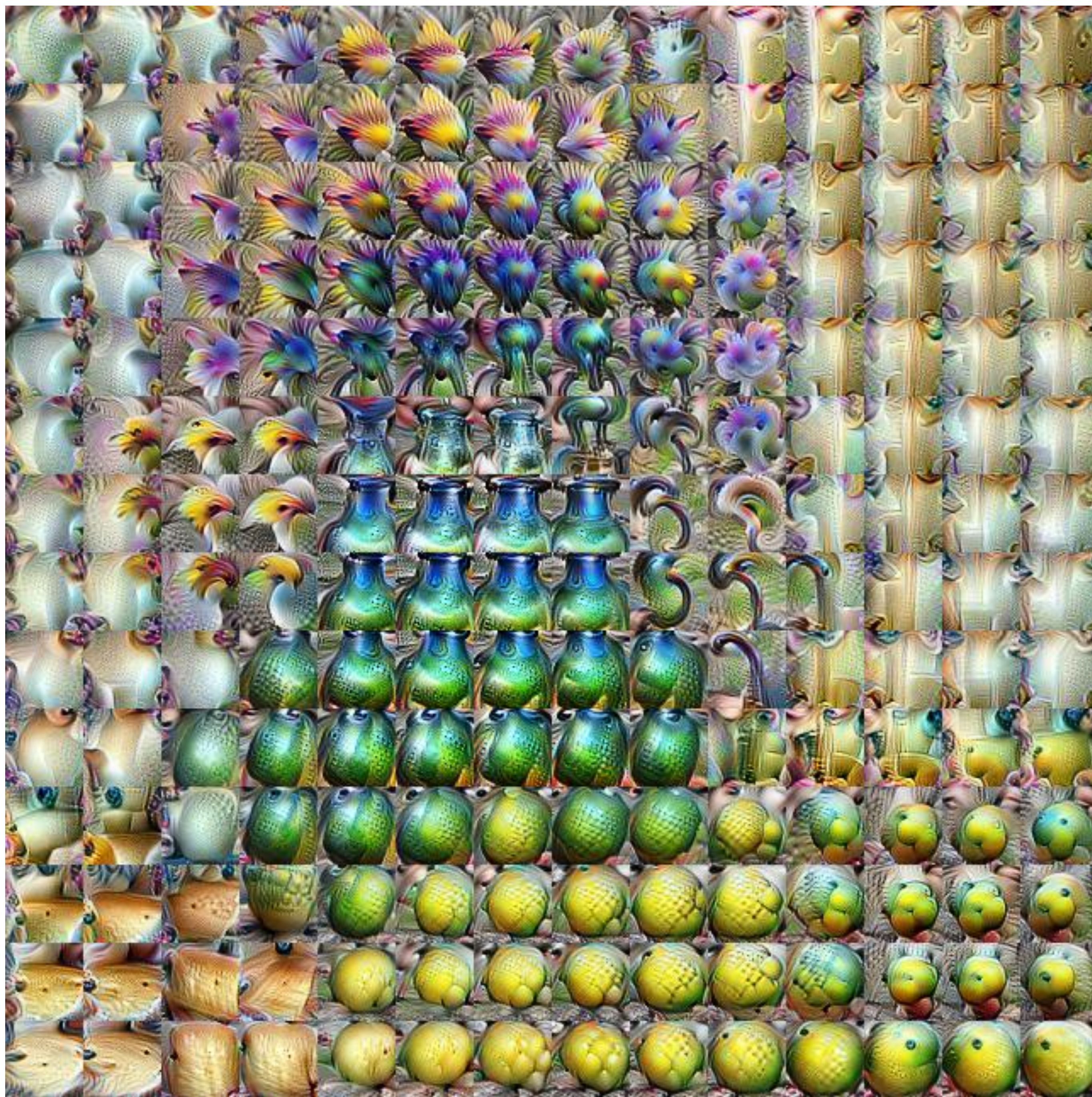
...

Channels

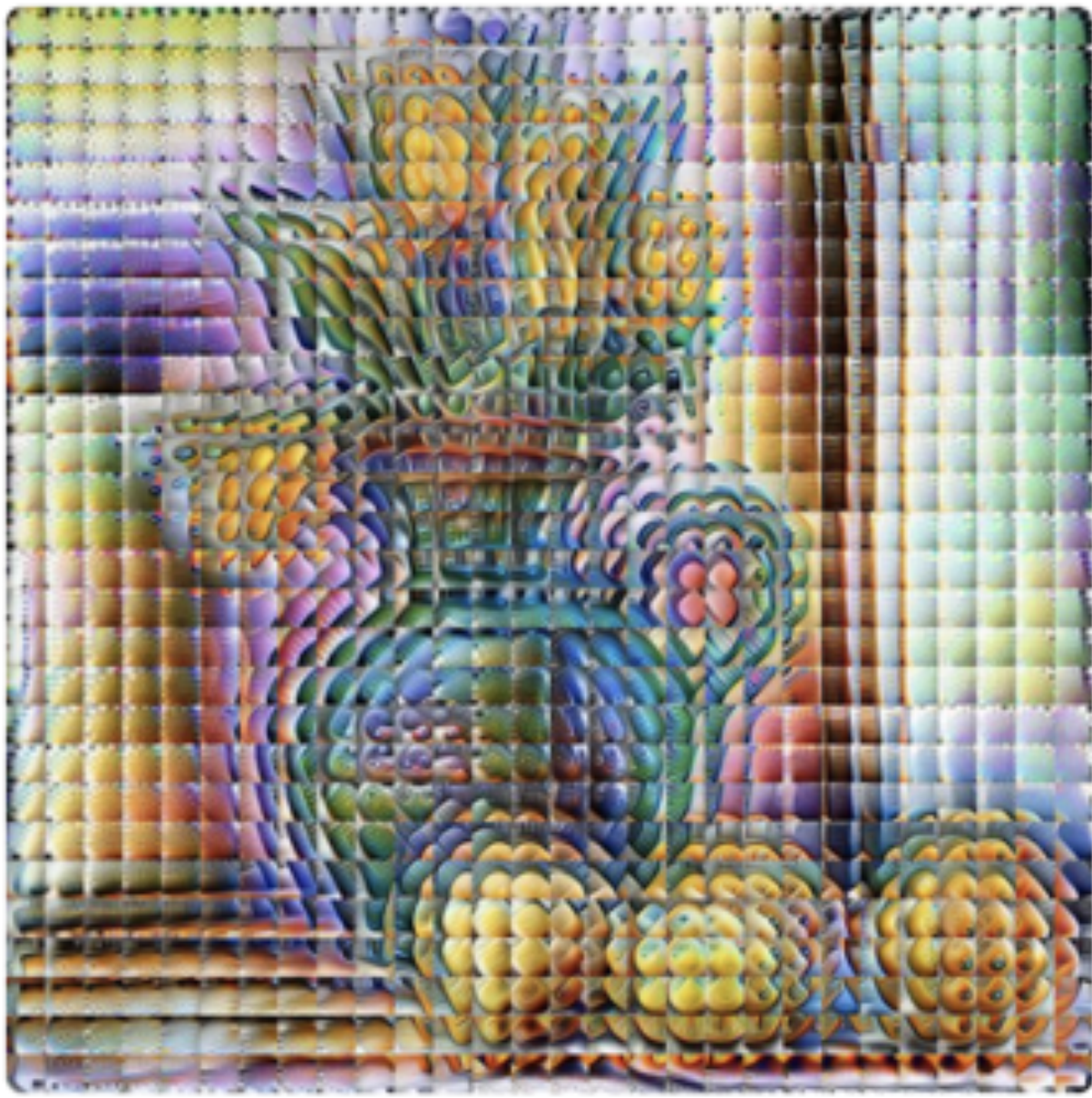




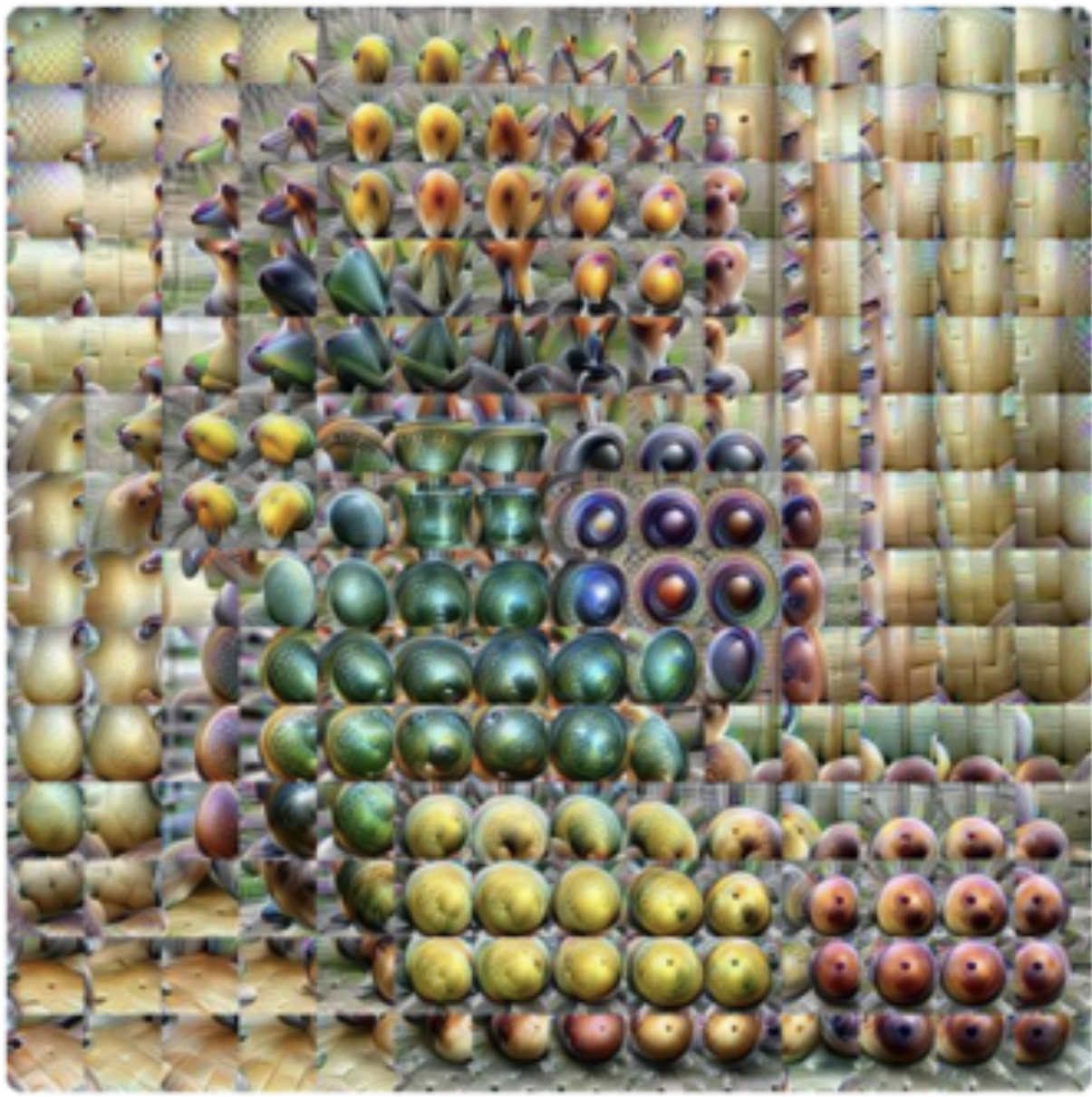
Activation Vector



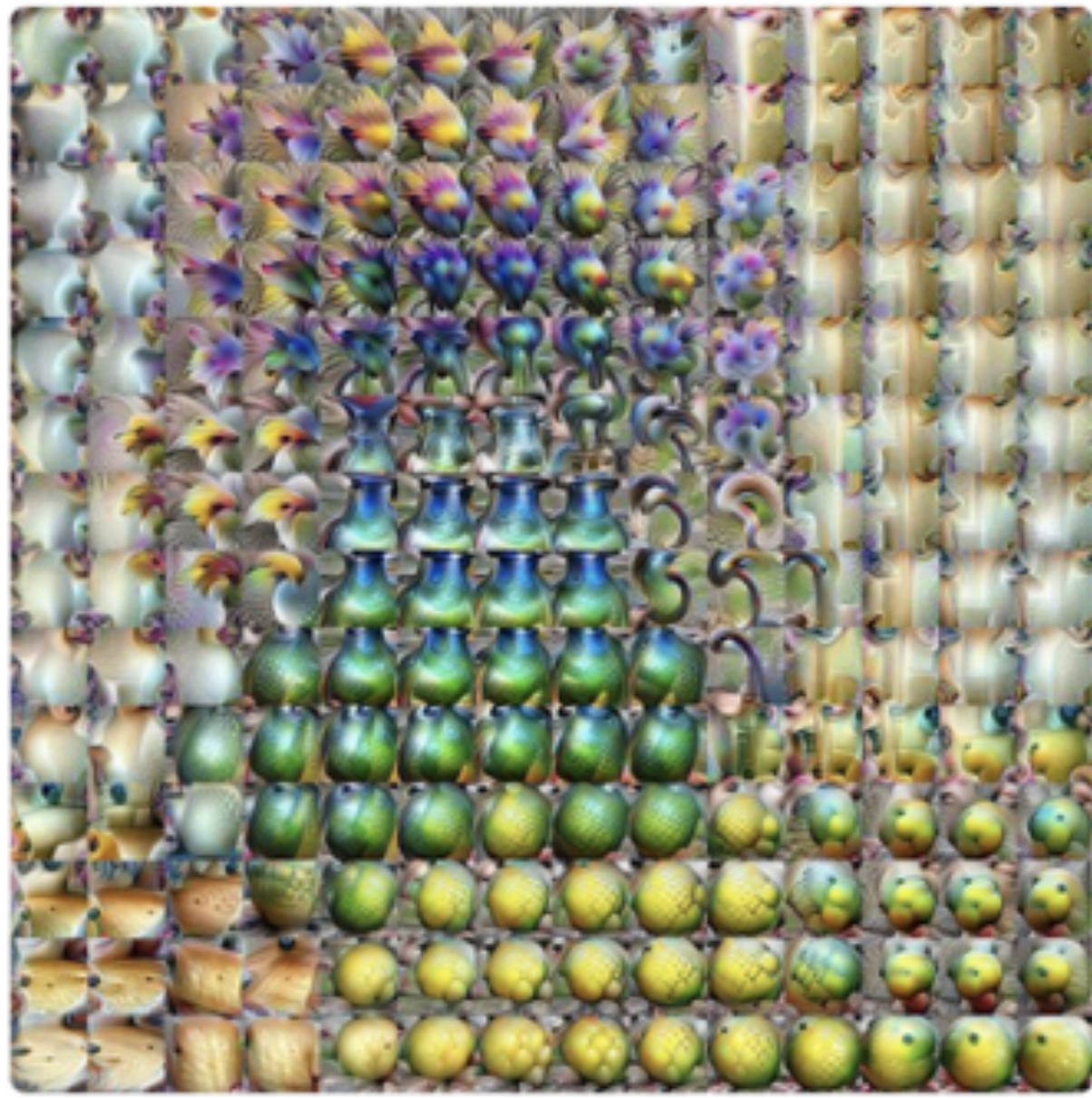




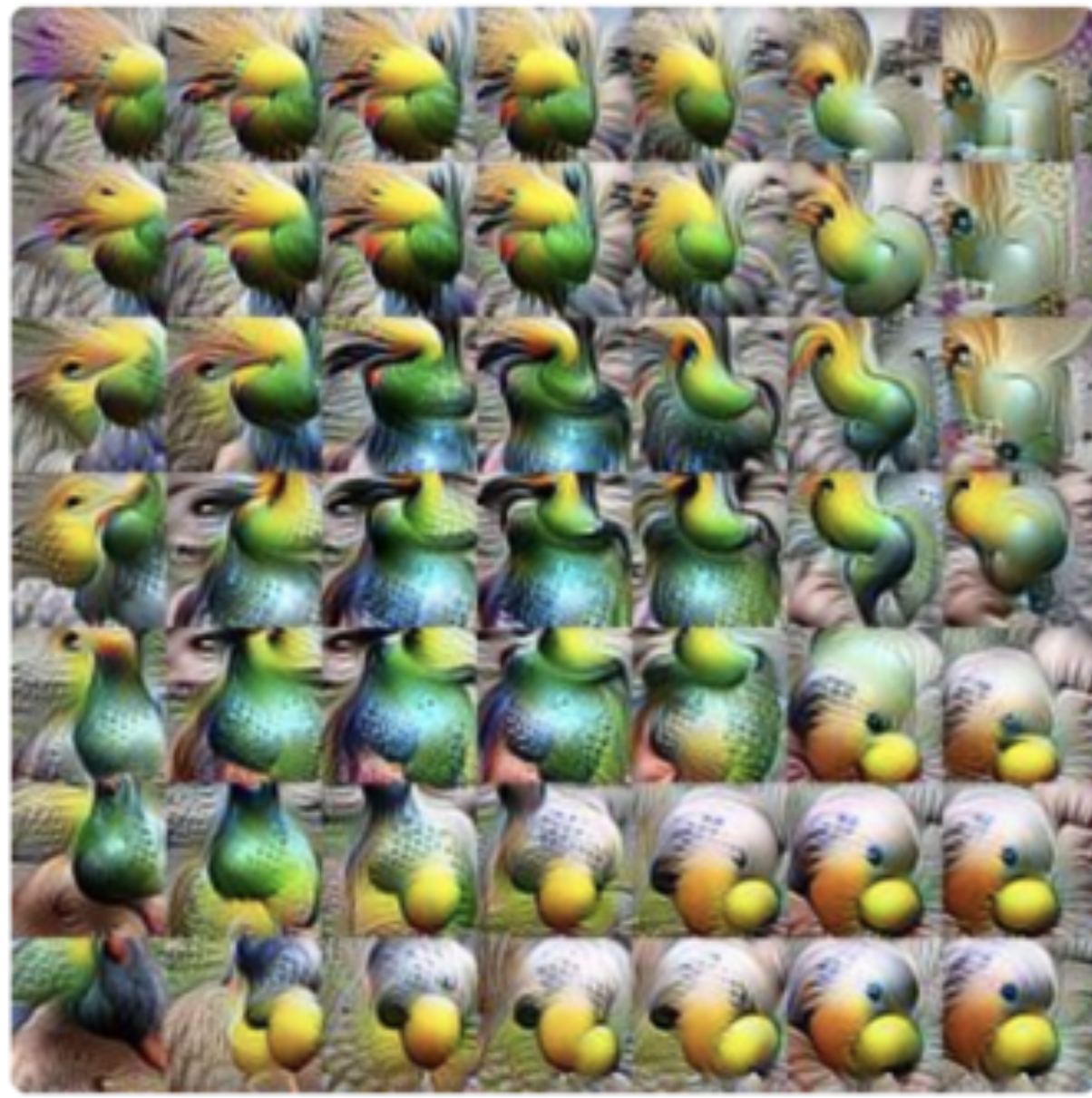
MIXED3A



MIXED4A



MIXED4D



MIXED5A




# The Building Blocks of Interpretability


Olah, Satyanarayan, et al. Distill, 2018.

<https://distill.pub/2018/building-blocks/>

Distill ABOUT PRIZE SUBMIT

CHOOSE AN INPUT IMAGE 

For instance, by combining feature visualization (what is a neuron looking for?) with attribution (how does it affect the output?), we can explore how the network decides between labels like **Labrador retriever** and **tiger cat**.



Several floppy ear detectors seem to be important when distinguishing dogs, whereas pointy ears are used to classify "tiger cat".

CHANNELS THAT MOST SUPPORT ... **LABRADOR RETRIEVER** **TIGER CAT**

feature visualization of channel

hover for attribution maps →

net evidence	1.63	1.51	1.19	...	1.32	1.54	1.72
for "Labrador retriever"	1.22	1.24	1.32	...	-0.70	-1.24	-0.43
for "tiger cat"	-0.40	-0.27	0.13	...	0.62	0.30	1.29

REPRODUCE IN A CO NOTEBOOK



# SUMMIT

## Scaling Deep Learning Interpretability by Visualizing Activation and Attribution Summarizations

*Fred Hohman, Haekyu Park, Caleb Robinson, Duen Horng (Polo) Chau*

*IEEE VIS 2019*

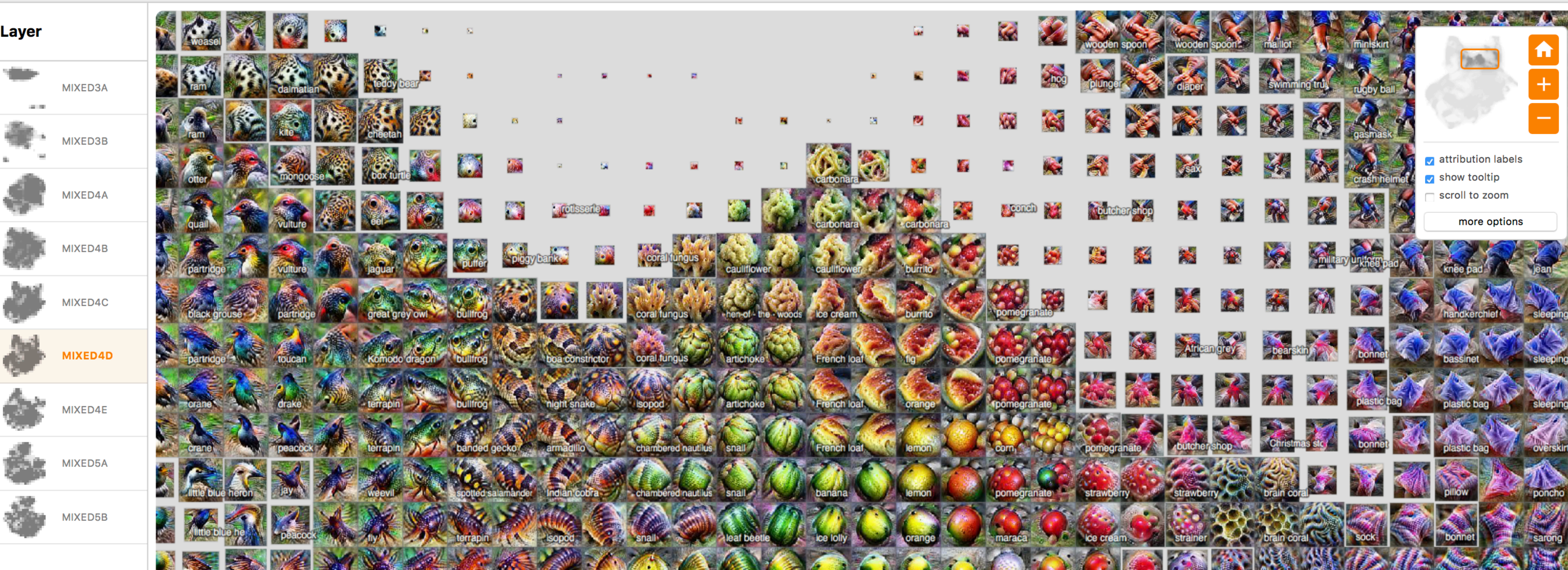
*Vancouver, Canada*



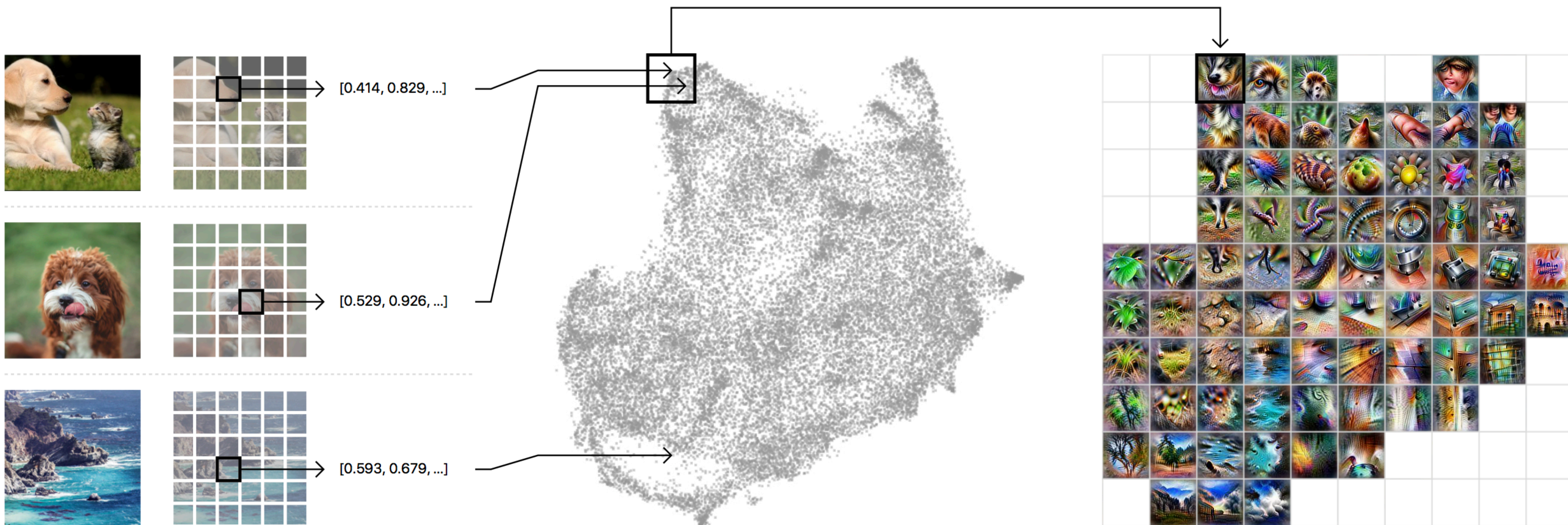


# Exploring Neural Networks with Activation Atlases

By using feature inversion to visualize millions of activations from an image classification network, we create an explorable *activation atlas* of features the network has learned which can reveal how the network typically represents some concepts.







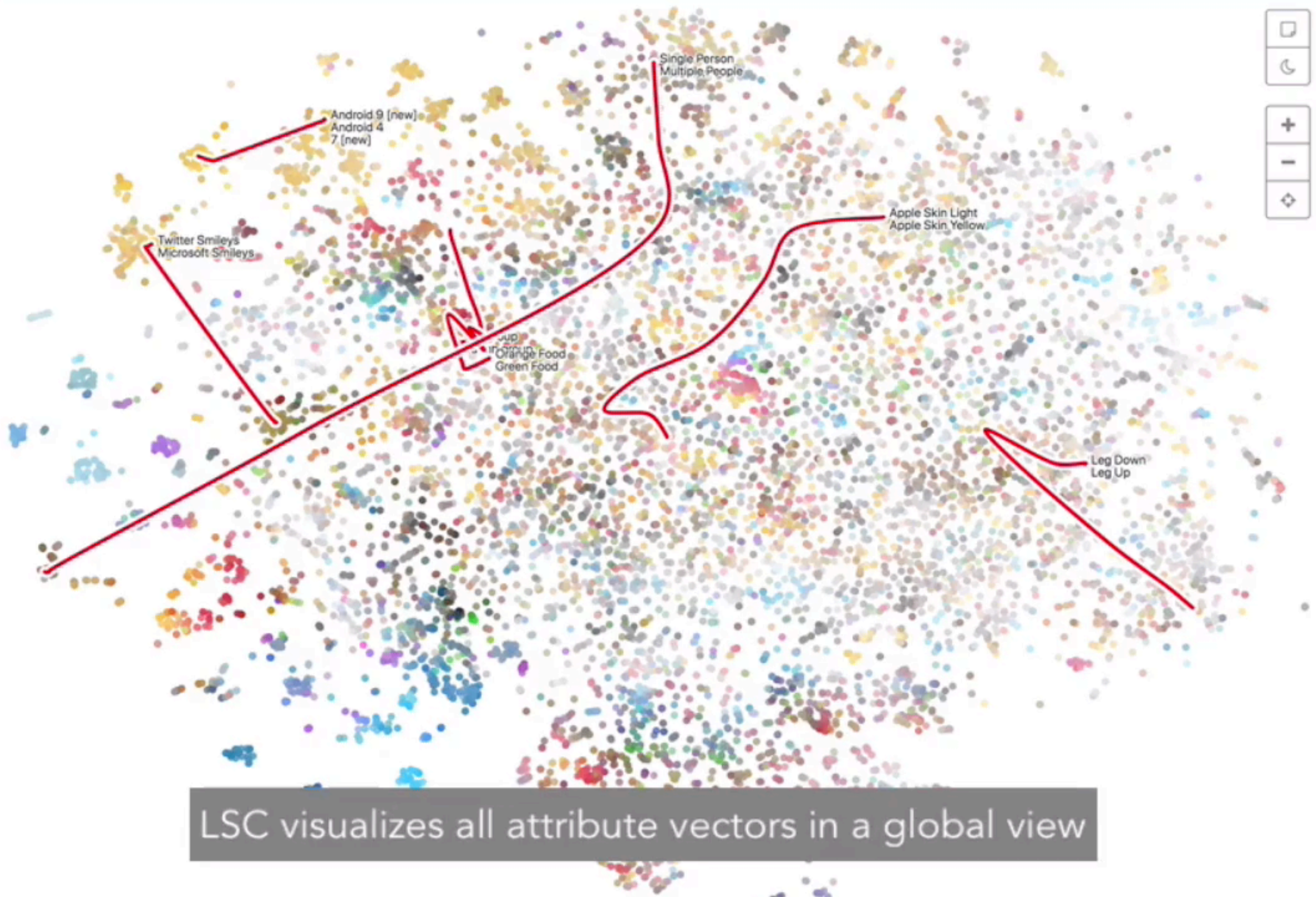
A randomized set of one million images is fed through the network, collecting one random spatial activation per image.

The activations are fed through UMAP to reduce them to two dimensions. They are then plotted, with similar activations placed near each other.

We then draw a grid and average the activations that fall within a cell and run feature inversion on the averaged activation. We also optionally size the grid cells according to the density of the number of activations that are averaged within.

 TRY IN A NOTEBOOK





Navigation controls: a square icon, a moon icon, a plus icon, a minus icon, and a diamond icon.

Groups Vectors

ADD VECTOR

○ — Add — ○

- VECTOR LIST
- 😄😄😄😄 Android 9 [new]
  - 😄😄😄😄 Android 4-7 [new]
  - 😭😭😭😭 Cry group
  - 😂😂😂😂 Laugh group
  - 🍊🍊🍊🍊 Orange Food
  - 🍏🍏🍏🍏 Green Food
  - 🦵🦵🦵🦵 Leg Down
  - 🦵🦵🦵🦵 Leg Up
  - 😄😄😄😄 Twitter Smileys
  - 😄😄😄😄 Microsoft Smileys
  - 🍏🍏🍏🍏 Apple Skin Light
  - 🍏🍏🍏🍏 Apple Skin Yellow
  - 🧑🧑🧑🧑 Man
  - 🧑🧑🧑🧑 Woman
  - 🧑🧑🧑🧑 Single Person
  - 🧑🧑🧑🧑 Multiple People

Search: 🔍 Latent Dimensions: 32 ▾ Projection: t-SNE ▾ Perplexity: 30 ▾ Category: All ▾ Platform: All ▾