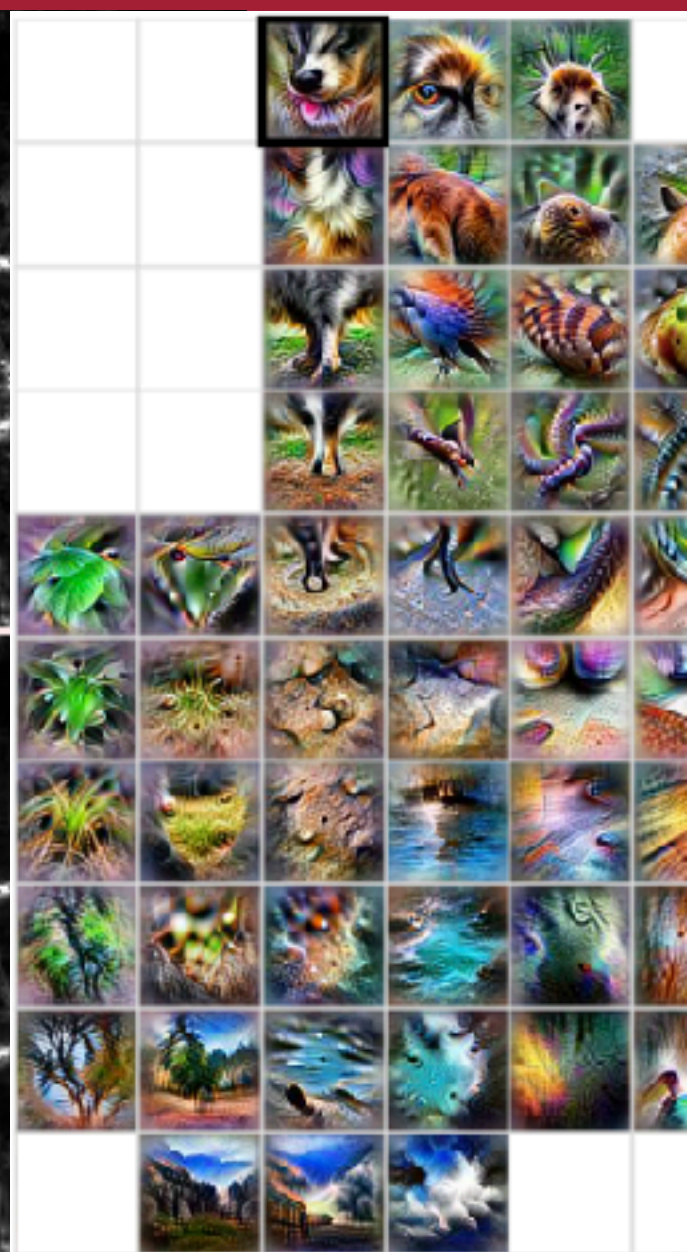
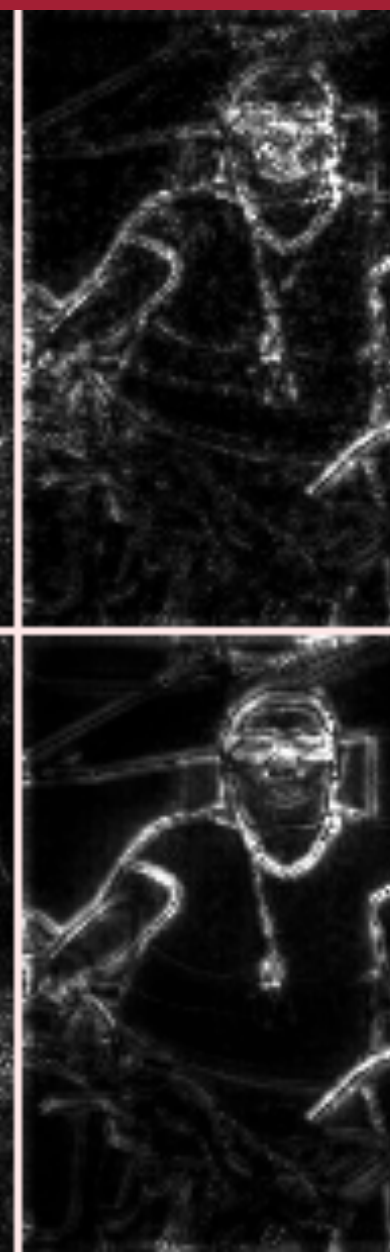
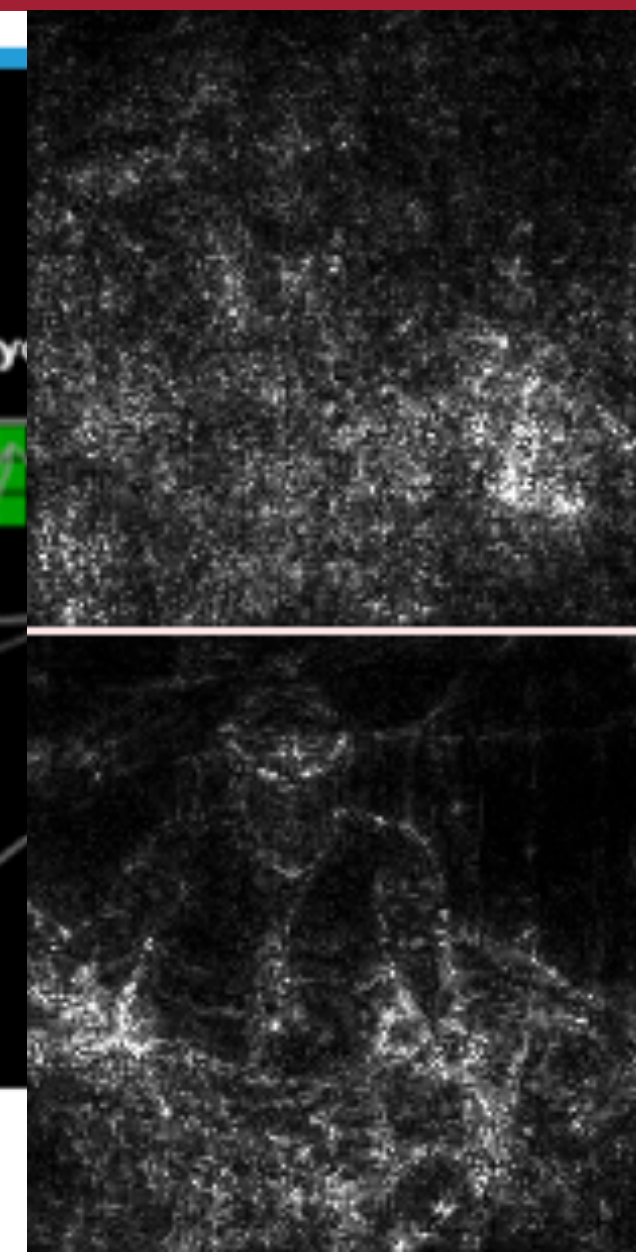
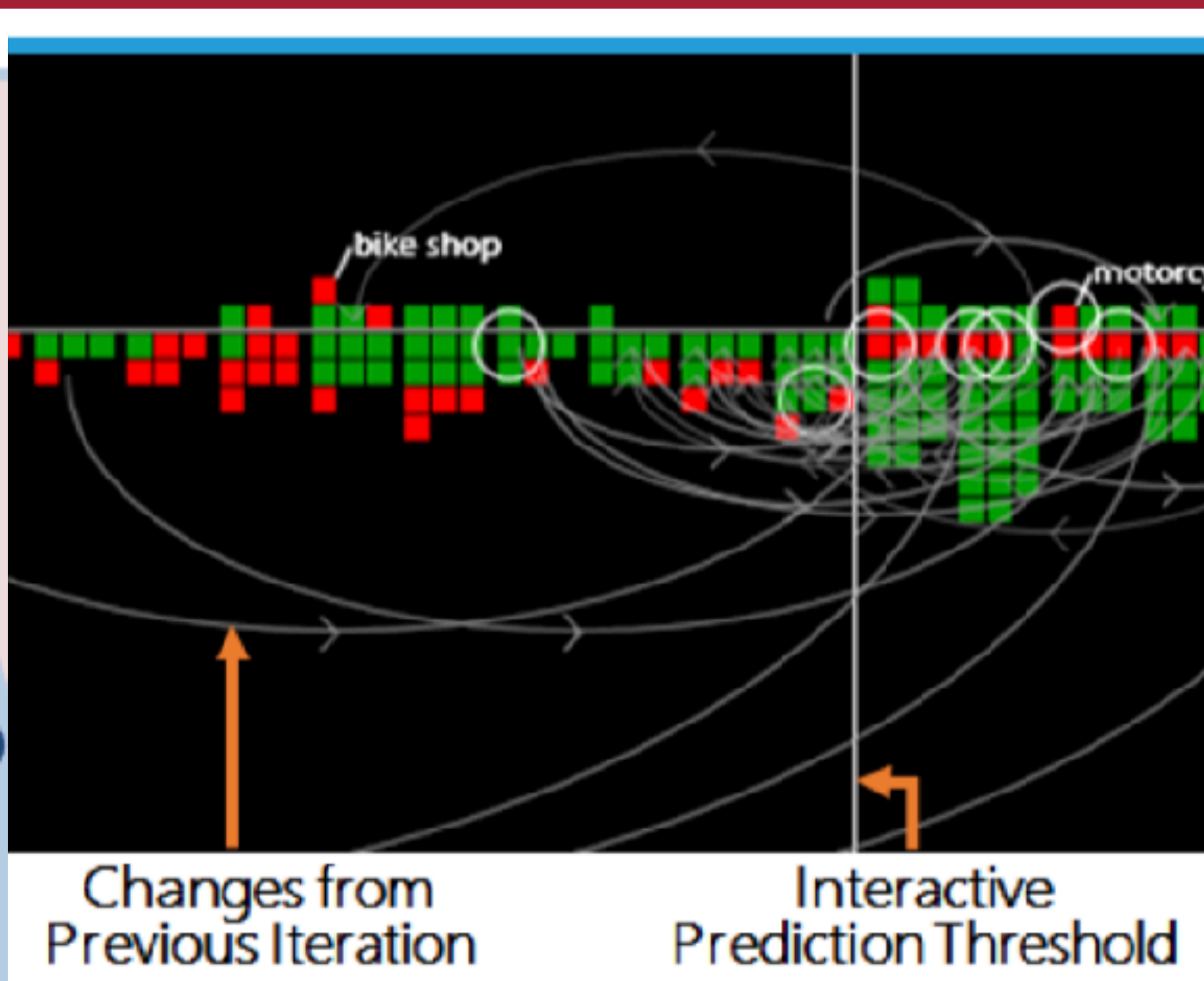
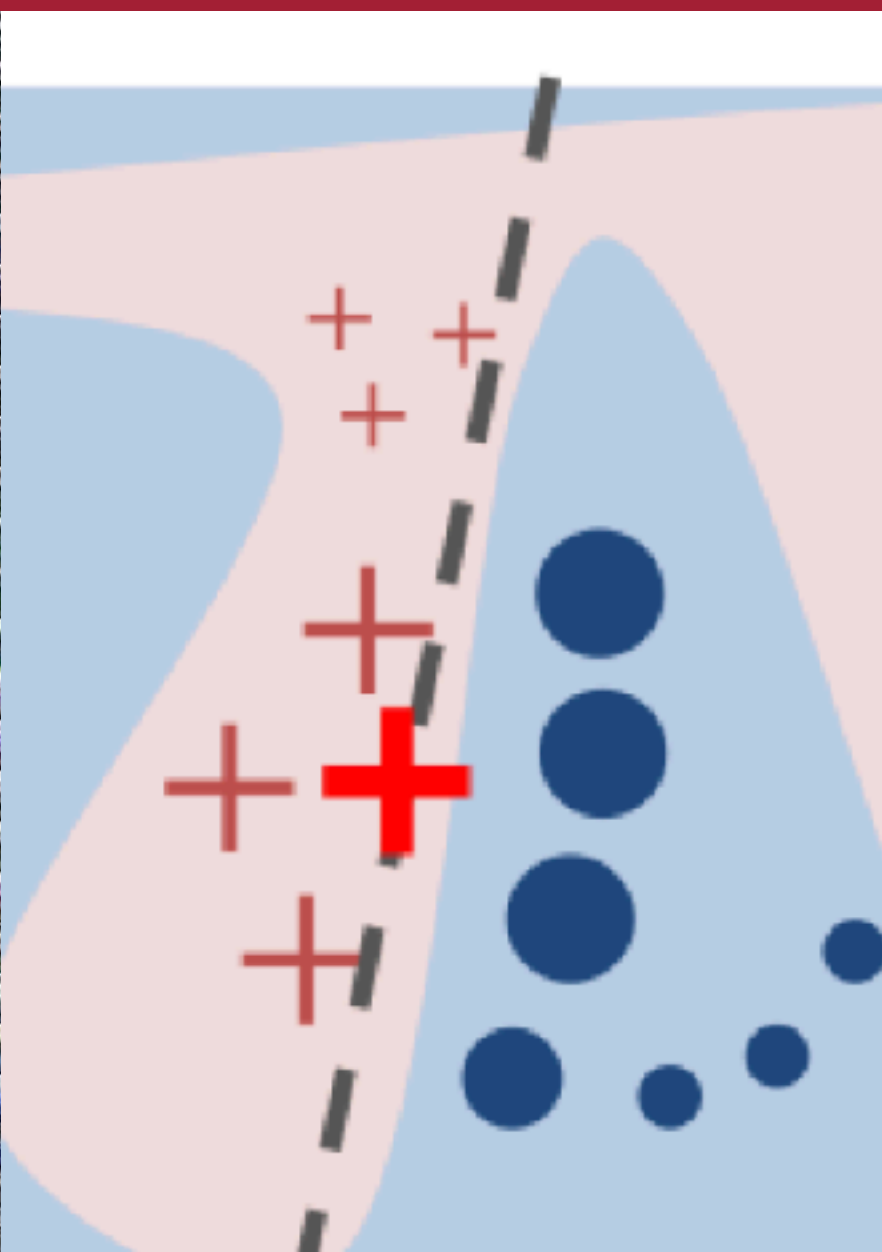


# 6.859: Interactive Data Visualization

## Research Frontiers: ML Interpretability

Arvind Satyanarayan



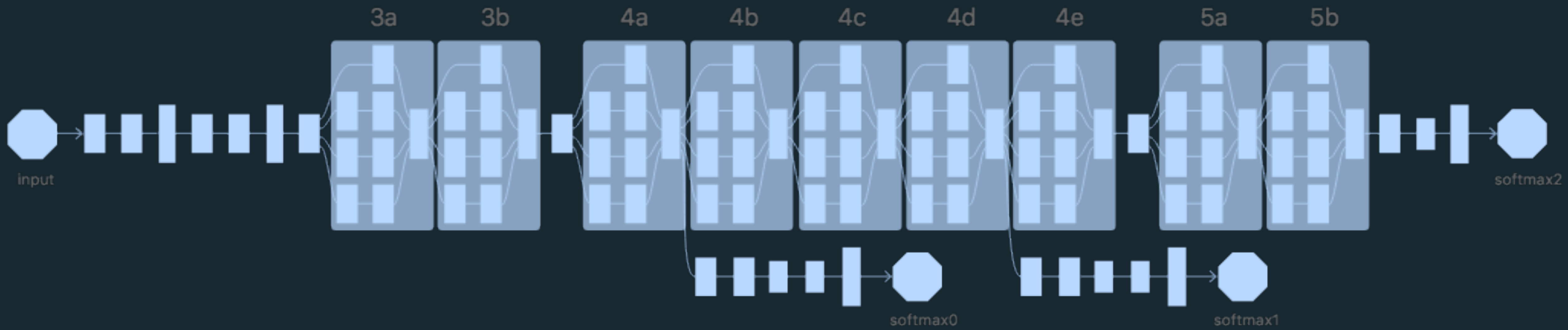




“Labrador Retriever”



“Golden Retriever”



“Labrador Retriever”



“Golden Retriever”

Label	Probability
Labrador retriever	69.2%
Golden retriever	11.6%
Tennis ball	2.6%

Source: ImageNet

# Why do we want interpretability?

**Fairness:** is the model biased? is it discriminating?

**Causality:** using models to infer properties about the natural world.

**Reliability:** how well does this model generalize, or transfer to a new domain?

**Trust:** how much confidence do I have in the model?

**Transparency:** how do I audit a model's decision-making?

- Browse...** **New...**
- Cycling Keywords**  
cycling, cycle, cyclist, cyclists, bic...
- Cycling Clubs and Organizations**  
cycling club, cycle club, club, cycl...
- Cycling Lessons**  
lessons, lesson, class, classes, tes...
- Bike Safety and Maintenance**  
safety, bike safety, maintenance, ...
- Cycling Tours and Events**  
cycling tour, biking tour, bike tou...
- motorcycle**  
motorcycles, honda, yamaha, red...

- Inactive Features
- Shopping and Rentals**  
shop, store, rental, rentals, produ...

Find Items

Sample

Found items around score of 0.75.

<http://www.tillamookrodeo.com>



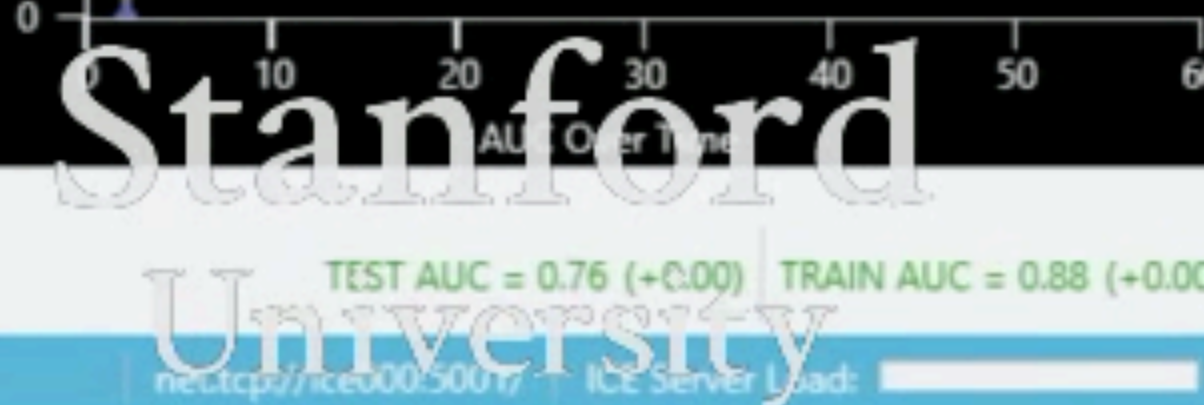
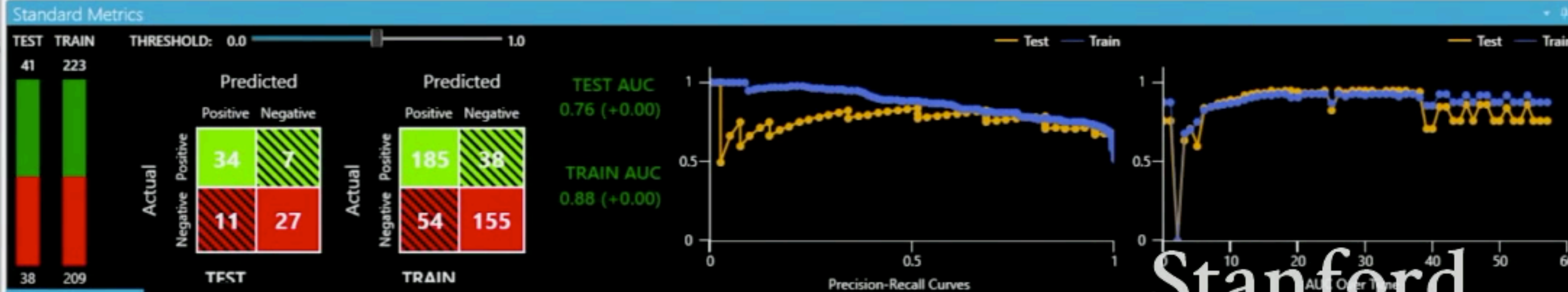
**0.25**

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



**0.75**

Find Items | Featuring | Feature Suggestions | Review Labels



FEATURES: 6 ITEMS: 511

Start

Session Features

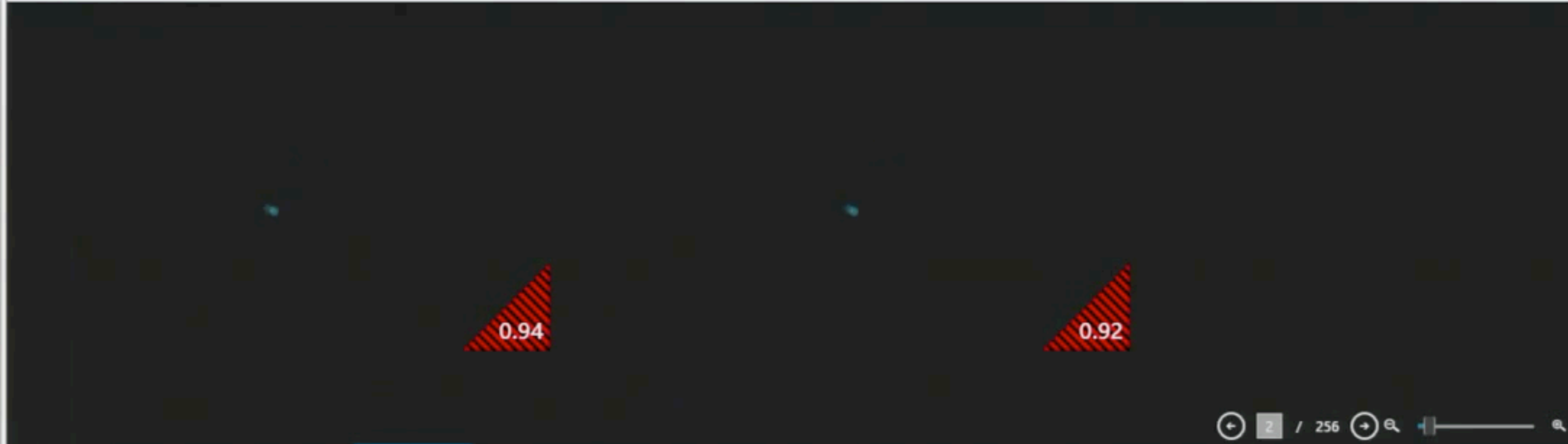
- Cycling Keywords**  
cycling, cycle, cyclist, cyclists, bic...
- Cycling Clubs and Organizations**  
cycling club, cycle club, club, cycl...
- Cycling Lessons**  
lessons, lesson, class, classes, tea...
- Bike Safety and Maintenance**  
safety, bike safety, maintenance...
- Cycling Tours and Events**  
cycling tour, biking tour, bike tou...
- motorcycle**  
motorcycles, honda, yamaha, rid...

Inactive Features

- Shopping and Rentals**  
shop, store, rental, rentals, produ...

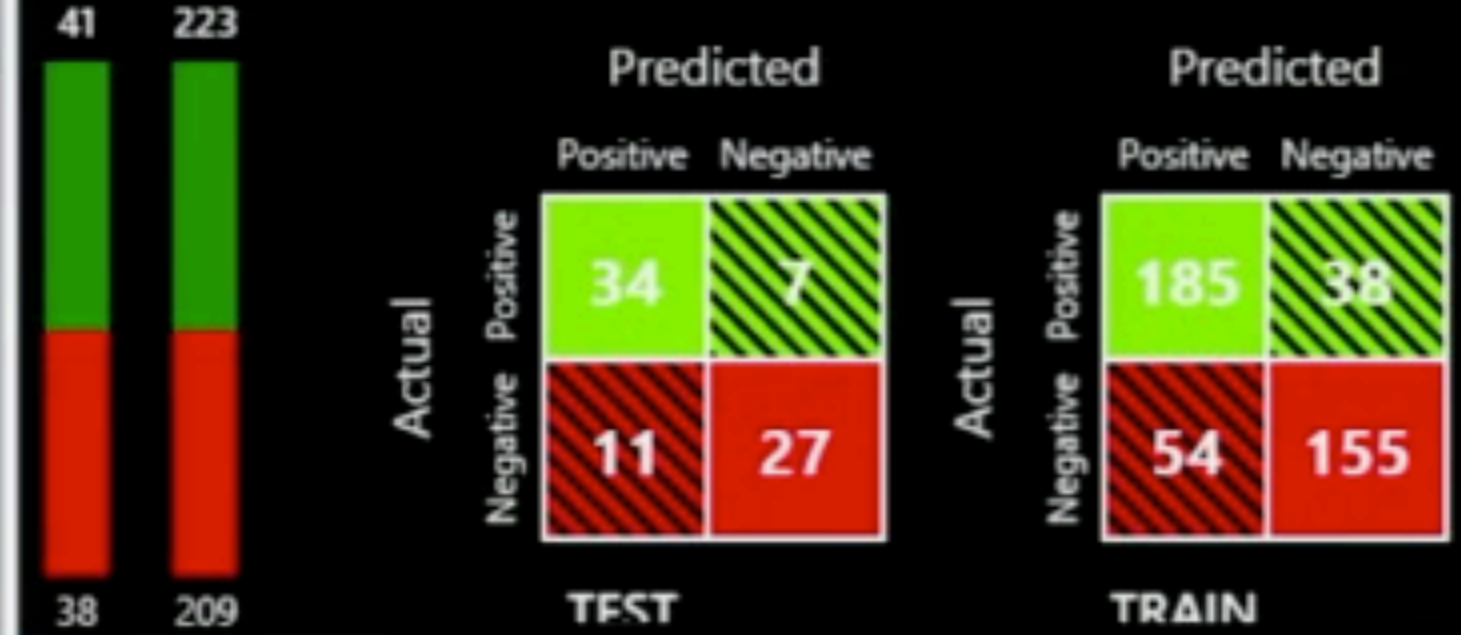
Review Labels

Review Labels toolbar: Refresh, Filter, Sort (Score ↓), Filter (All Data), Errors, Refresh View, Save Changes

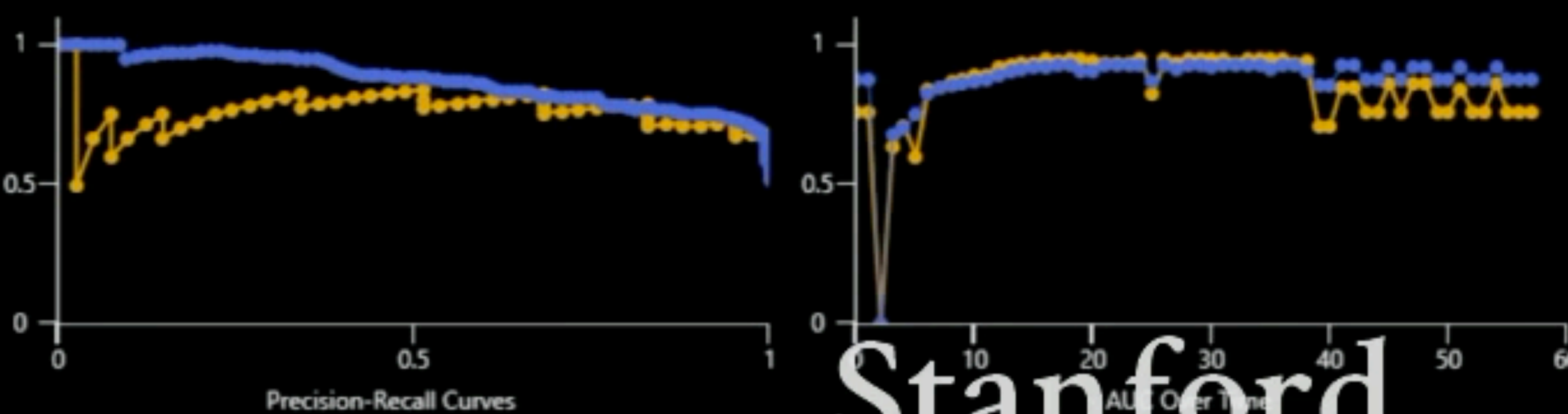


Standard Metrics

Standard Metrics: TEST 41, TRAIN 223, THRESHOLD: 0.0 (slider), 1.0



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



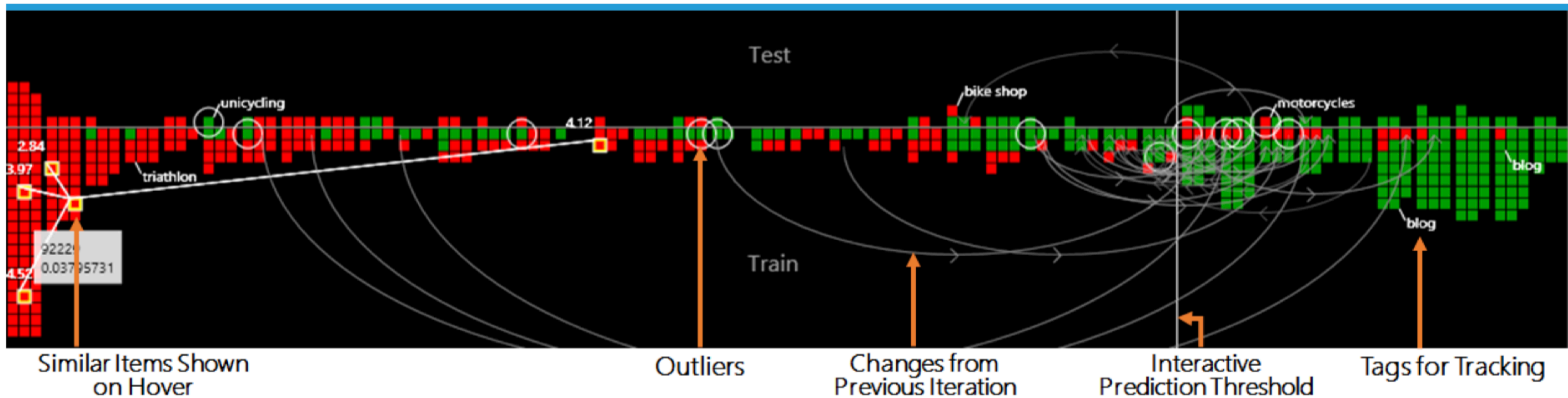
FEATURES: 6 ITEMS: 511



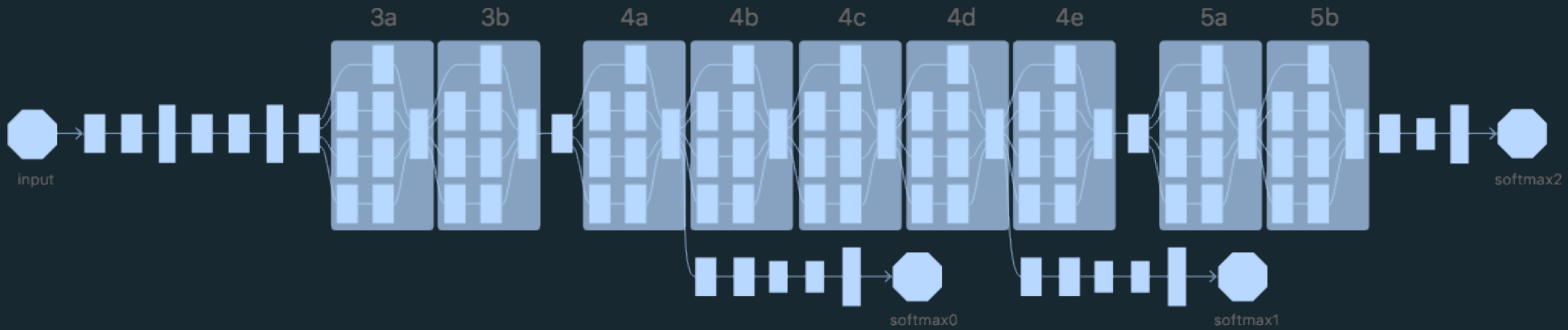
Stanford  
University



# ModelTracker



[Amershi et al., CHI 2015]



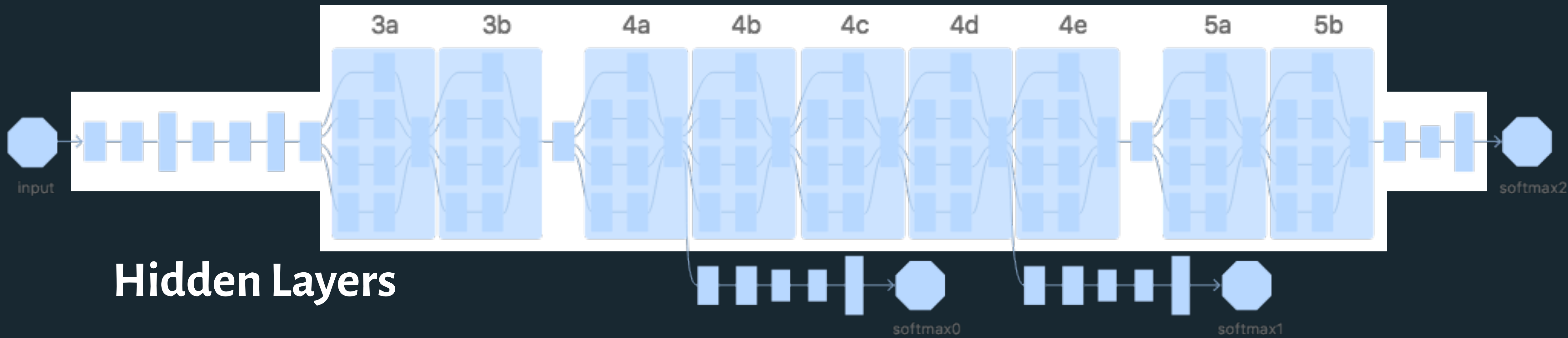
“Labrador Retriever”



“Golden Retriever”

Label	Probability
Labrador retriever	69.2%
Golden retriever	11.6%
Tennis ball	2.6%

Source: ImageNet



# Hidden Layers



“Labrador Retriever”

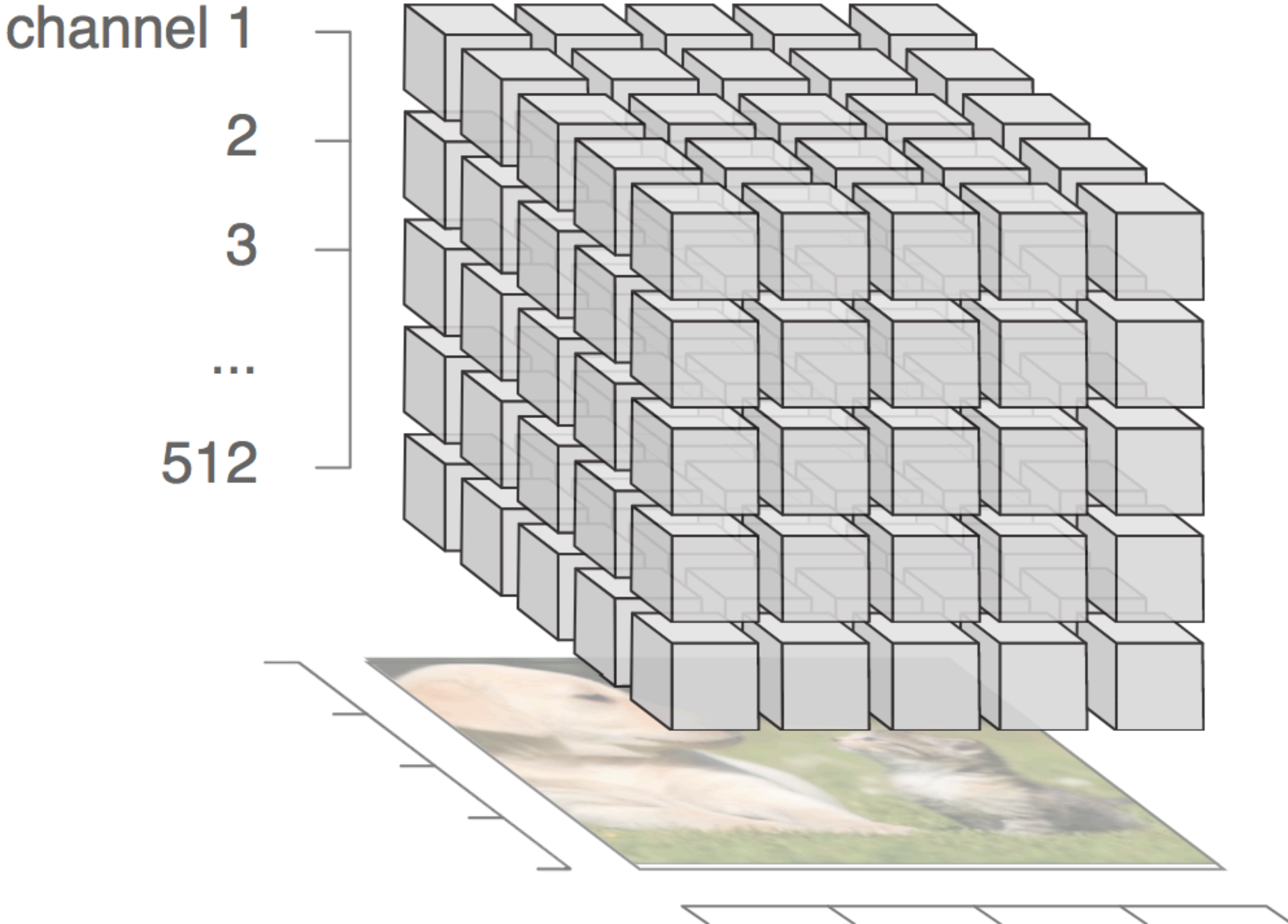


“Golden Retriever”

Label	Probability
Labrador retriever	69.2%
Golden retriever	11.6%
Tennis ball	2.6%

Source: ImageNet

# Individual Neurons



# Dimensionality Reduction

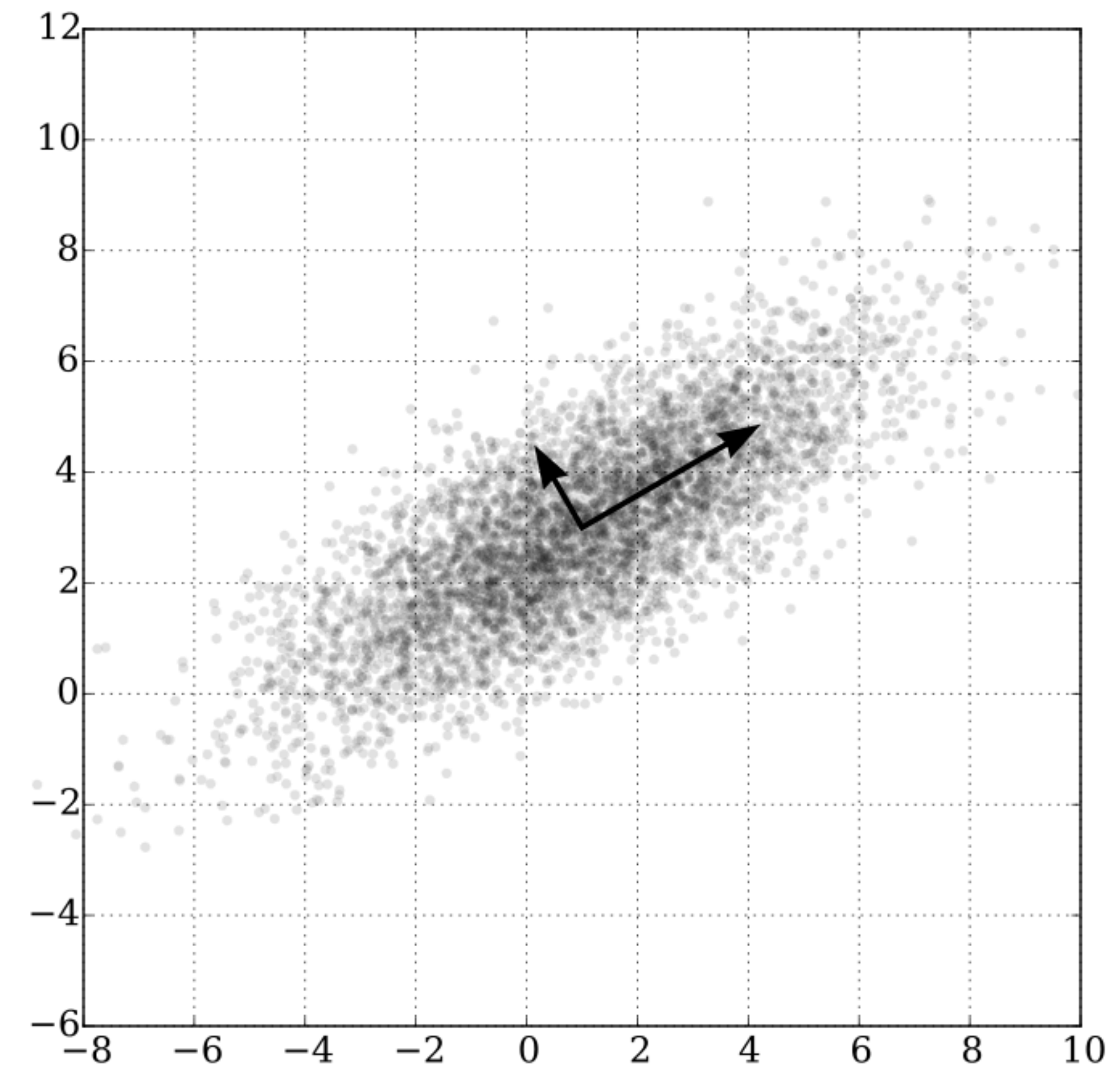
Project nD data to 2D or 3D. Interpret/sanity check learned representations.

But **subject to their own interpretation issues**.

Different techniques make different trade-offs:

**PCA (Principal Component Analysis):**

roughly fit a p-dimensional ellipsoid to the data, order axes by amount of data variance they explain. **Preserves global structure.**



# Dimensionality Reduction

Project nD data to 2D or 3D. Interpret/sanity check learned representations.

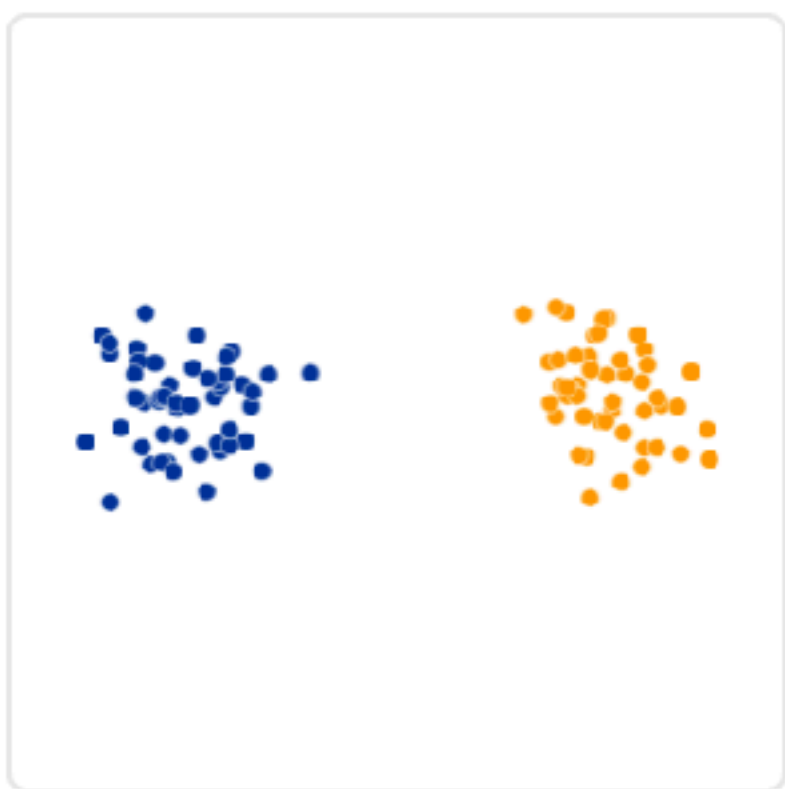
But **subject to their own interpretation issues**.

Different techniques make different trade-offs:

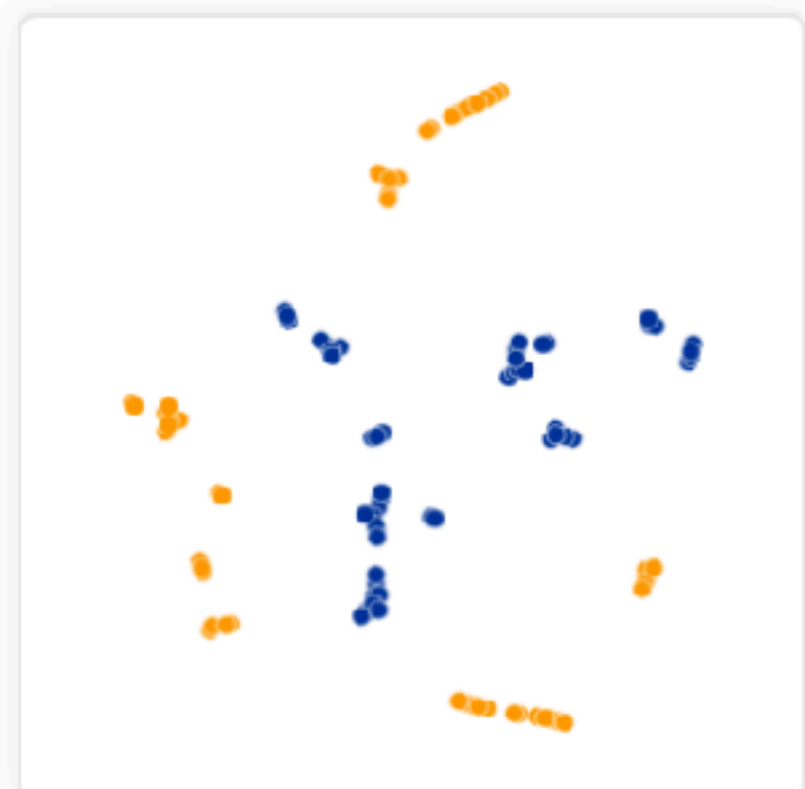
**PCA (Principal Component Analysis)**: roughly fit a p-dimensional ellipsoid to the data, order axes by amount of data variance they explain.

**Preserves global structure.**

**t-SNE (t-Dist. Stochastic Neighbor Embedding)**: probabilistic distribution that adapts and performs different transformations on different regions.



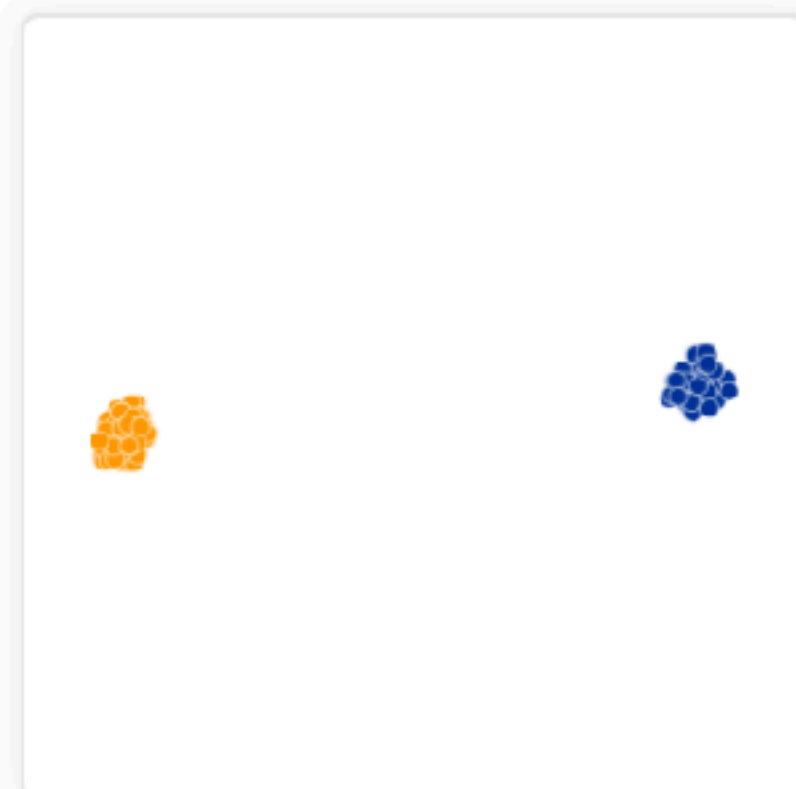
*Original*



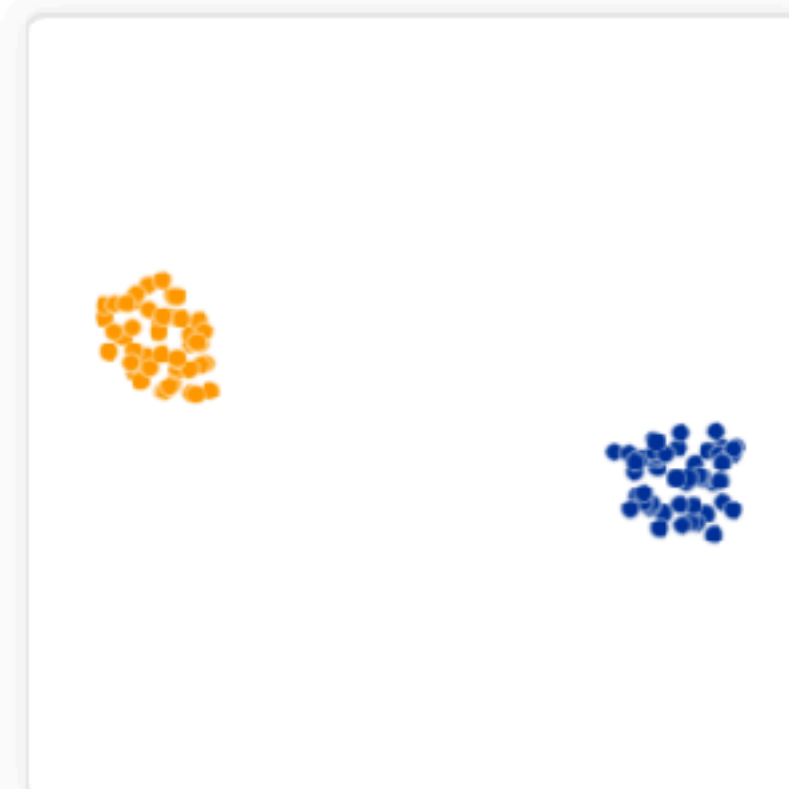
Perplexity: 2  
Step: 5,000



Perplexity: 5  
Step: 5,000



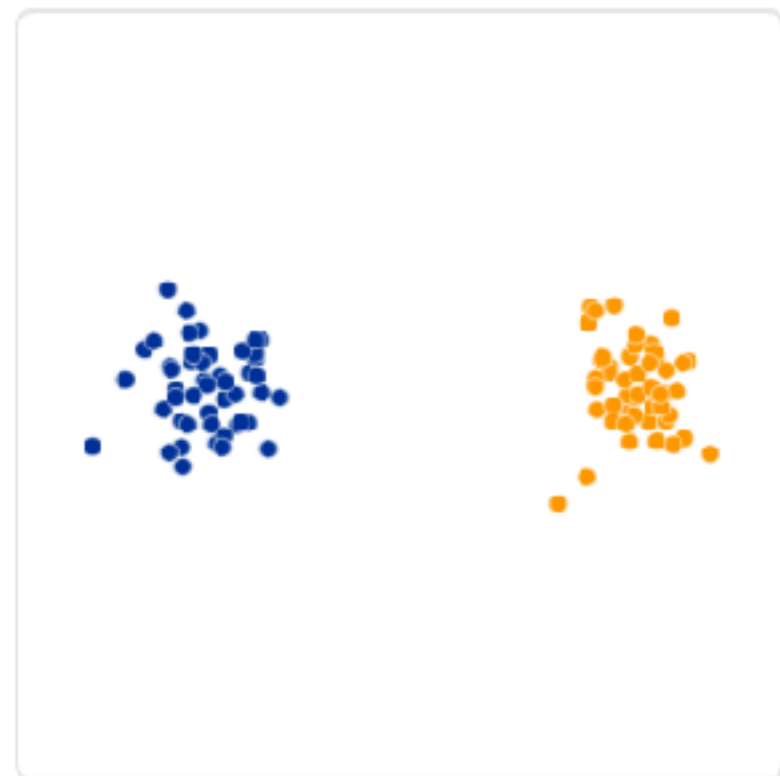
Perplexity: 30  
Step: 5,000



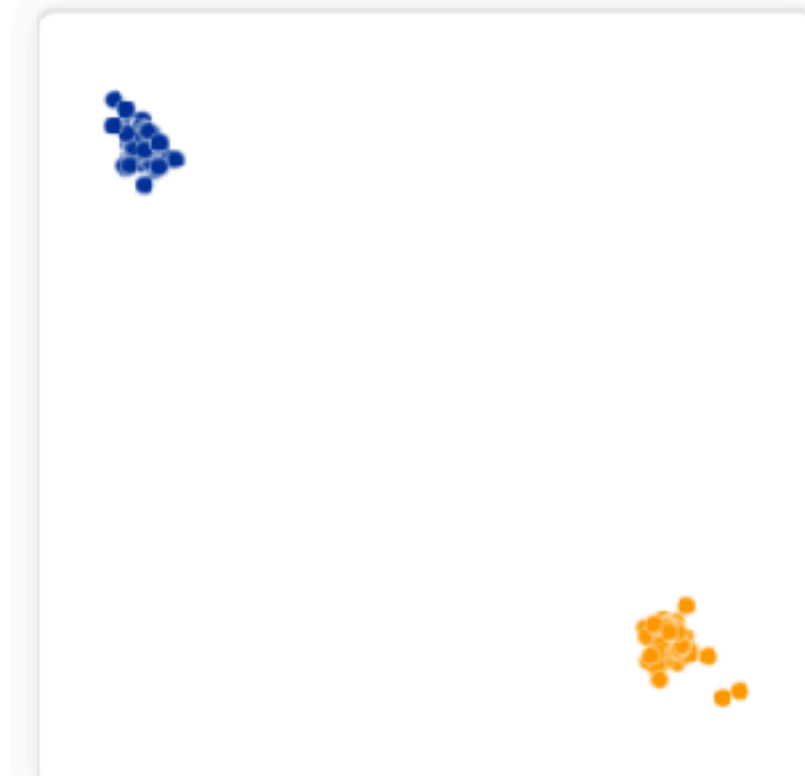
Perplexity: 50  
Step: 5,000



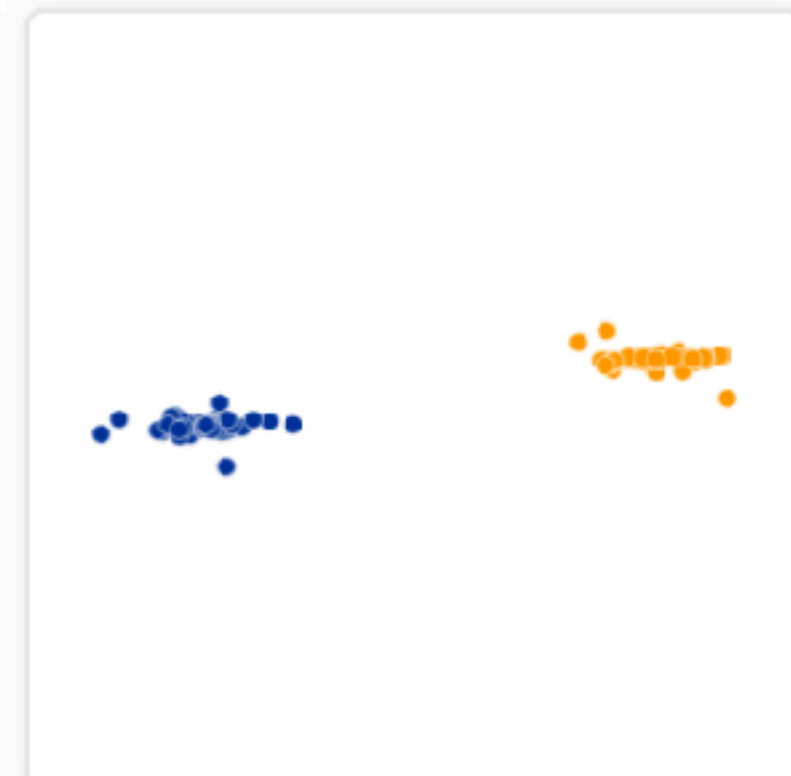
Perplexity: 100  
Step: 5,000



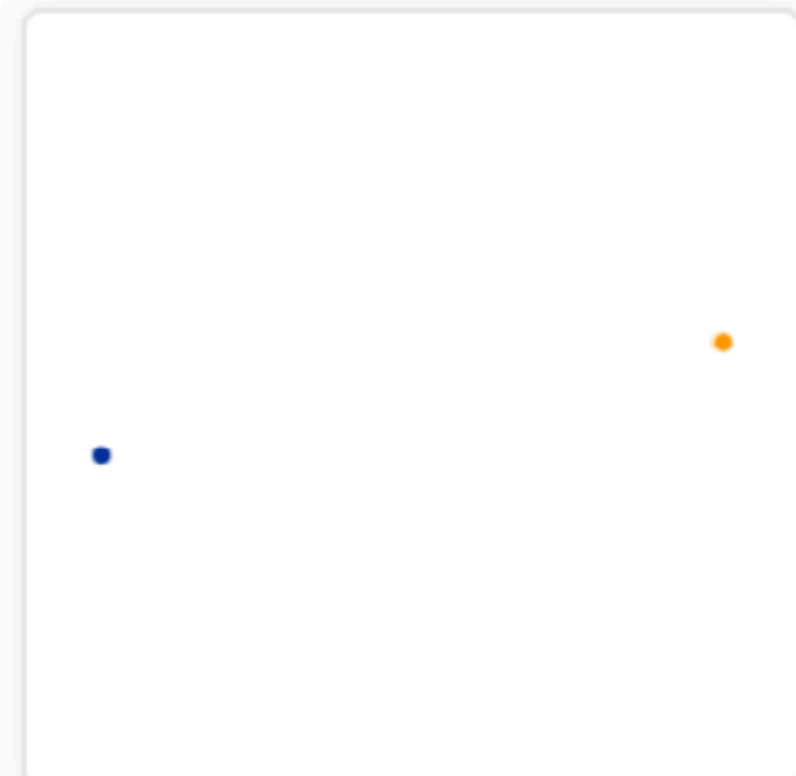
*Original*



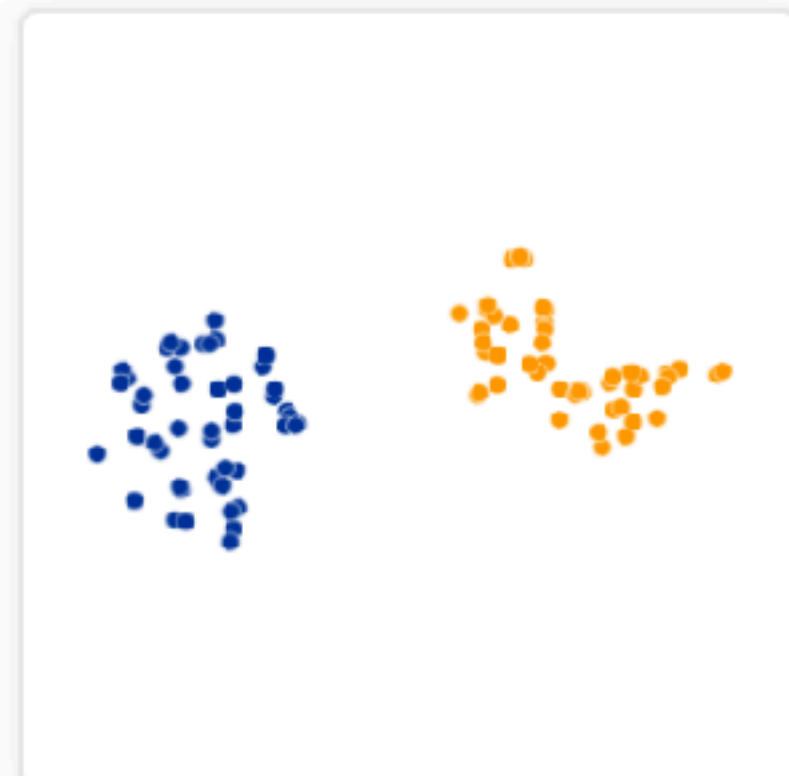
Perplexity: 30  
Step: 10



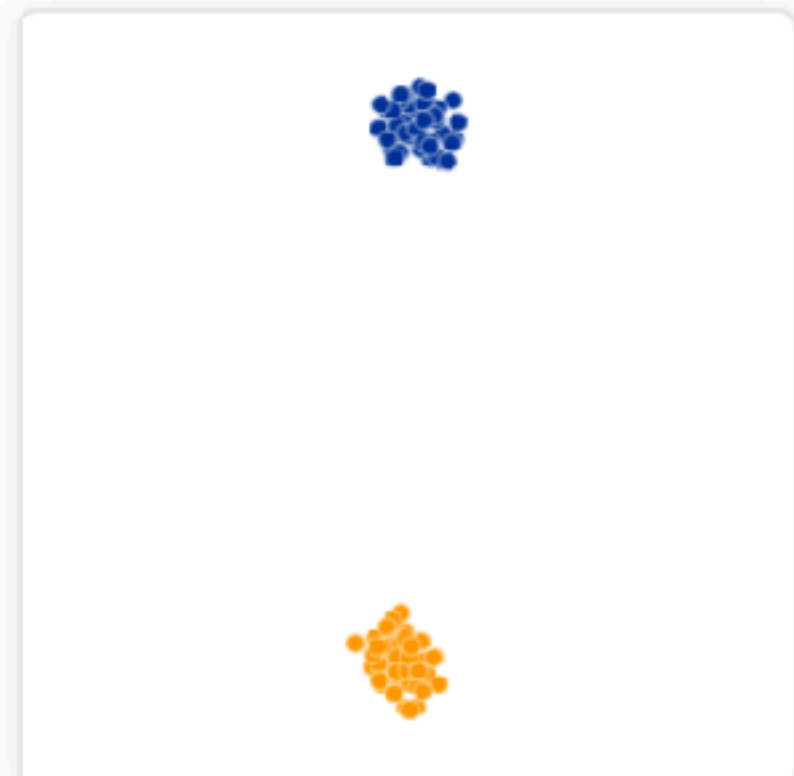
Perplexity: 30  
Step: 20



Perplexity: 30  
Step: 60

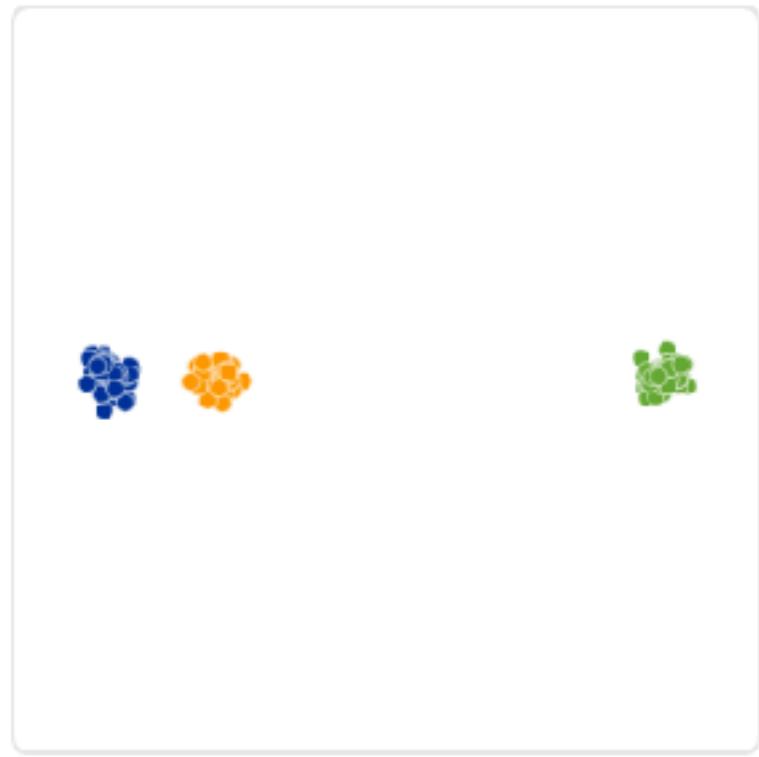


Perplexity: 30  
Step: 120

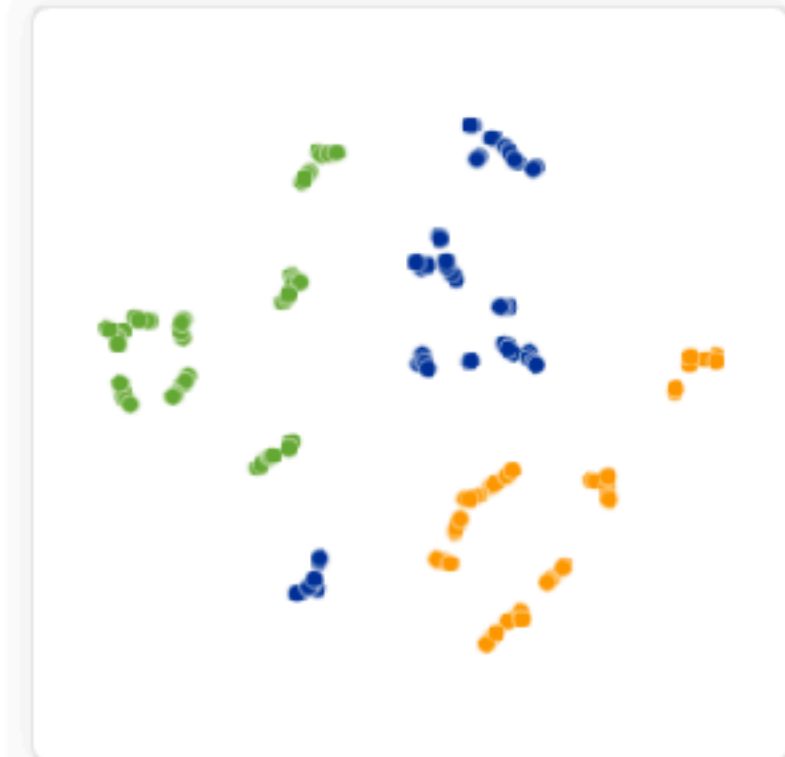


Perplexity: 30  
Step: 1,000

[Wattenberg et al., Distill 2016]



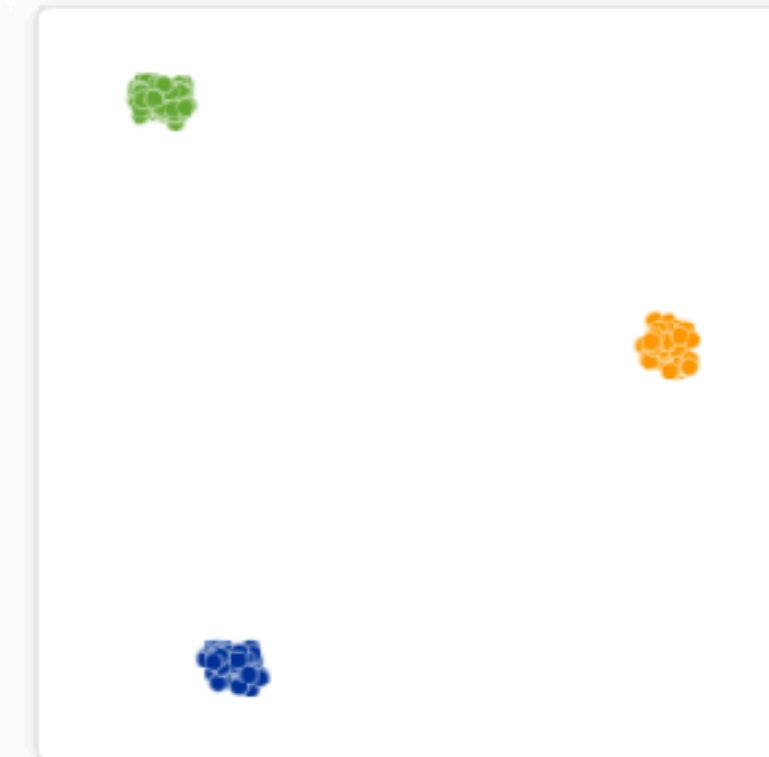
*Original*



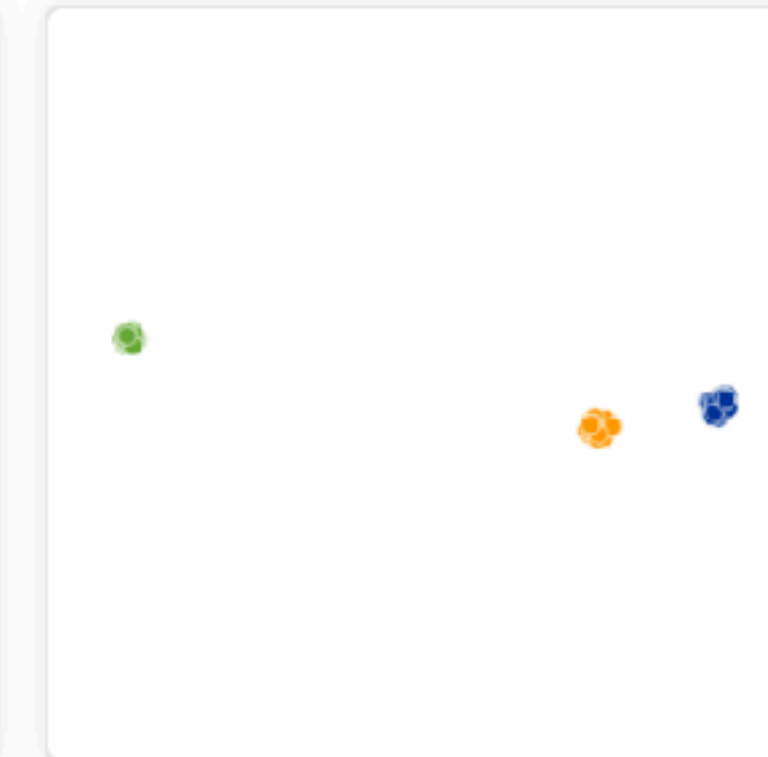
Perplexity: 2  
Step: 5,000



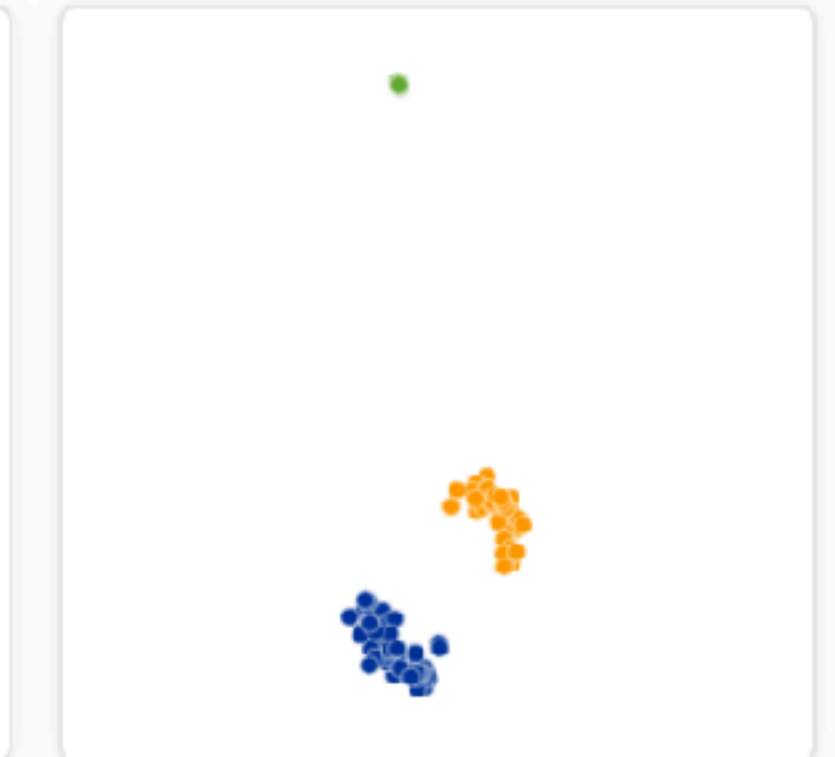
Perplexity: 5  
Step: 5,000



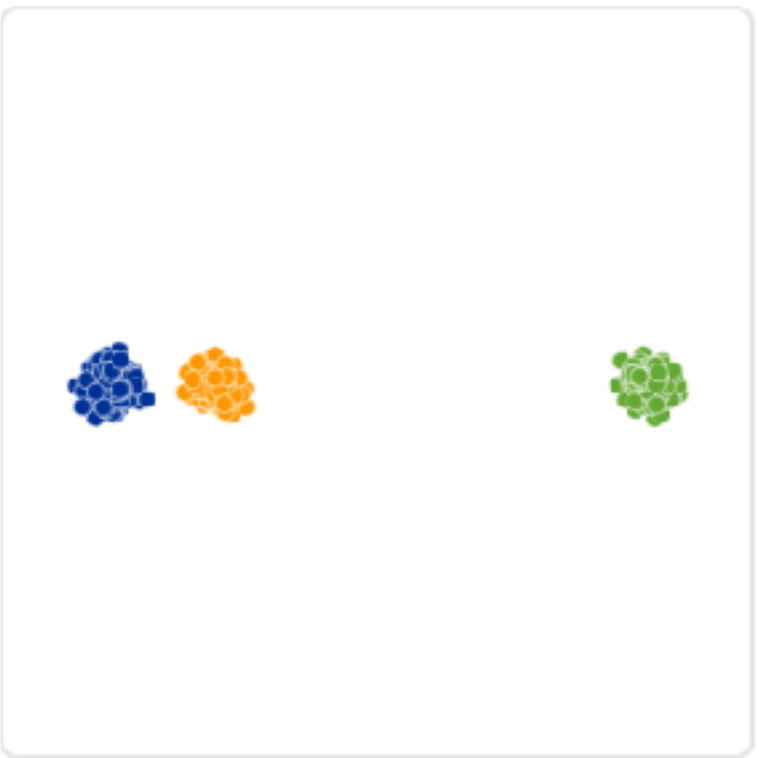
Perplexity: 30  
Step: 5,000



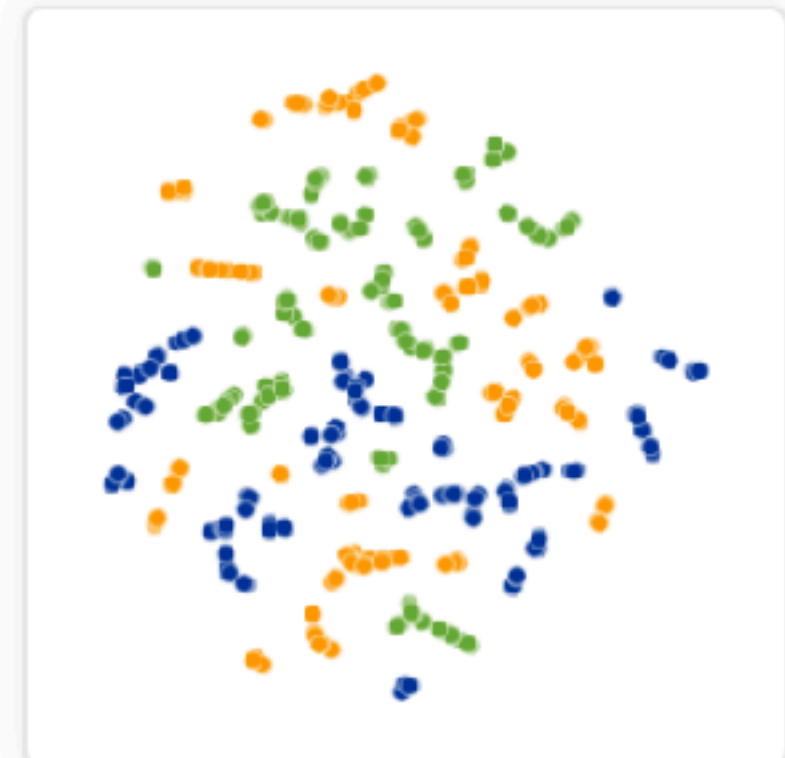
Perplexity: 50  
Step: 5,000



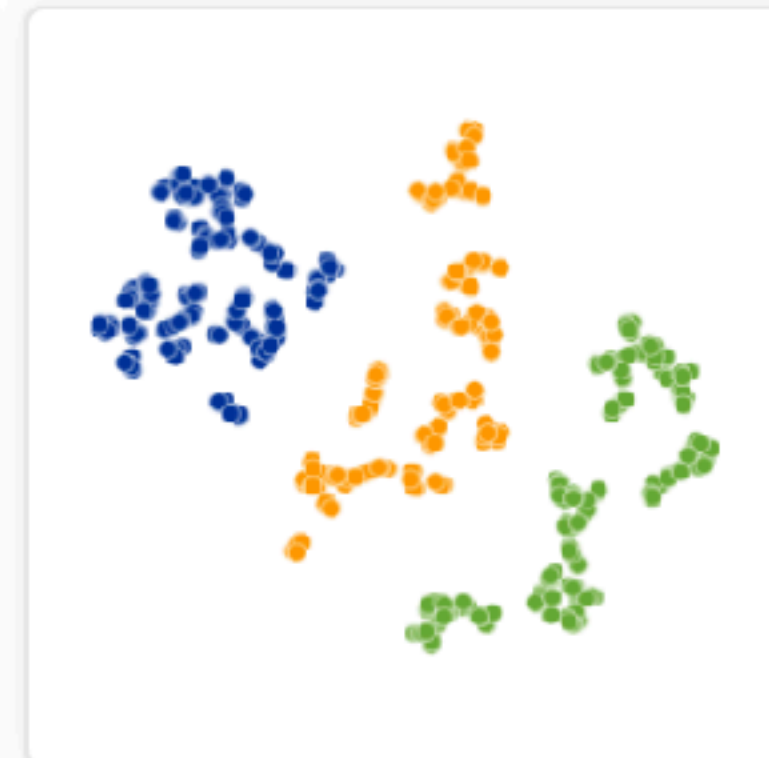
Perplexity: 100  
Step: 5,000



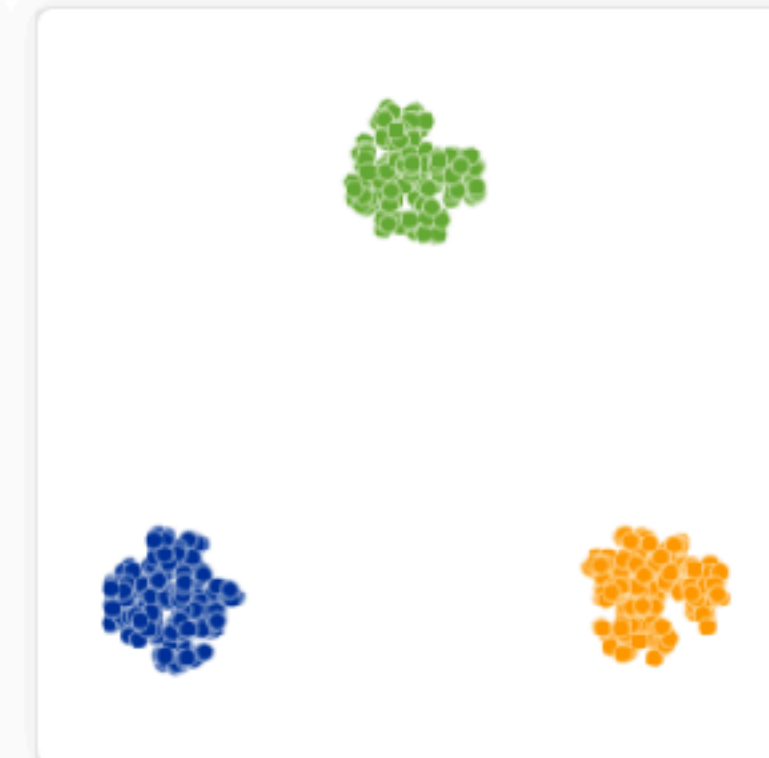
*Original*



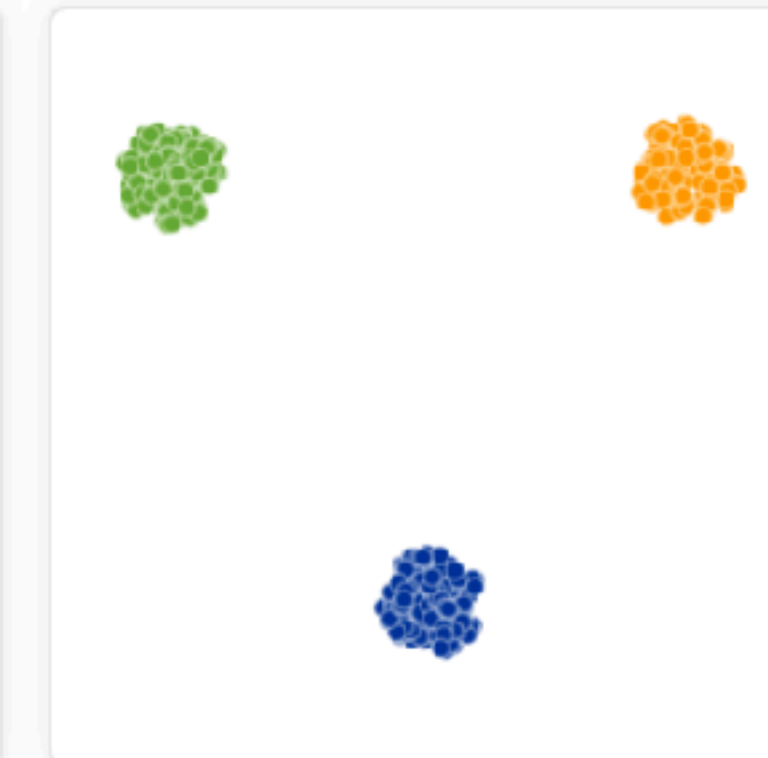
Perplexity: 2  
Step: 5,000



Perplexity: 5  
Step: 5,000



Perplexity: 30  
Step: 5,000



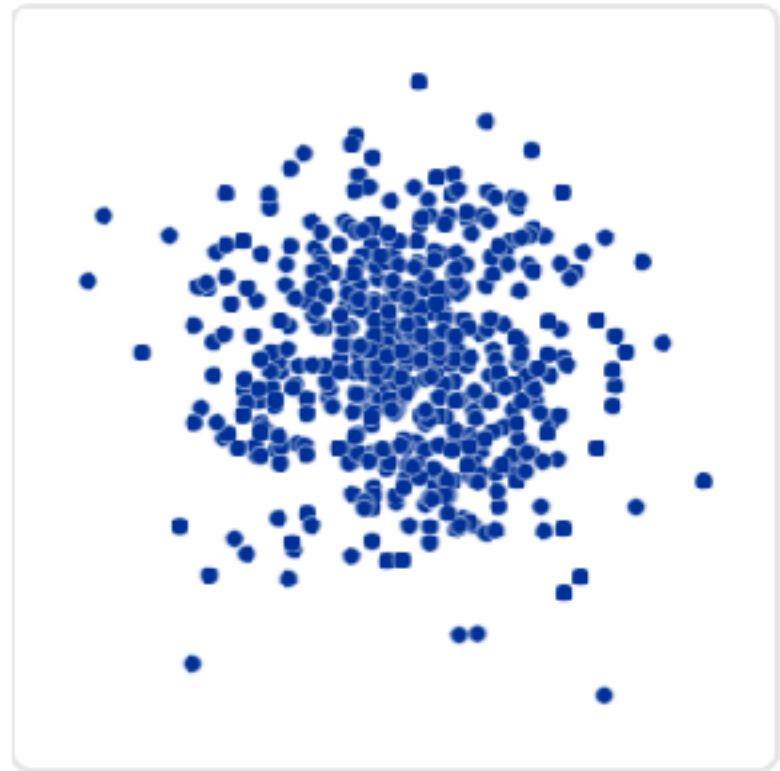
Perplexity: 50  
Step: 5,000



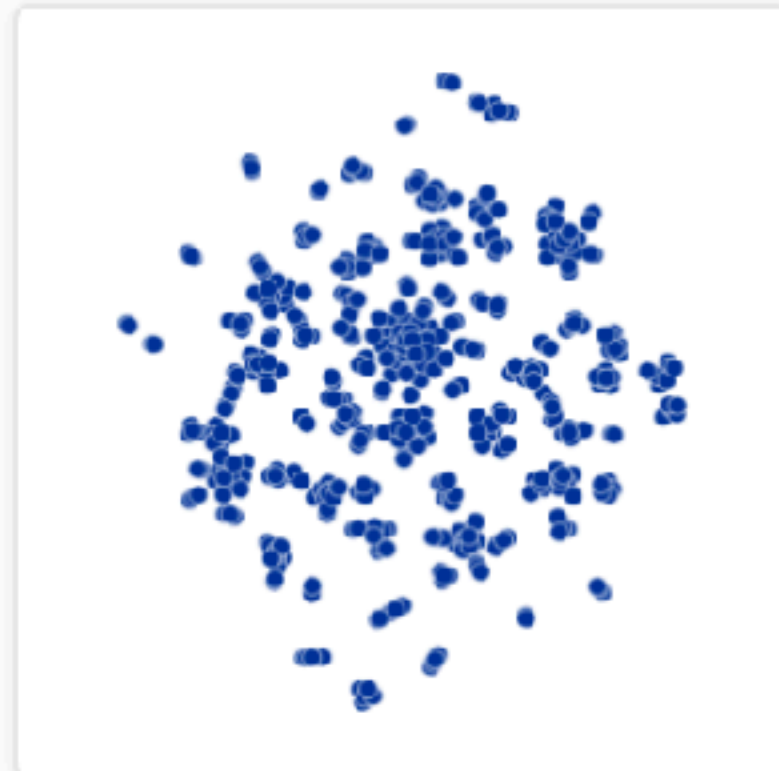
Perplexity: 100  
Step: 5,000

[Wattenberg et al., Distill 2016]





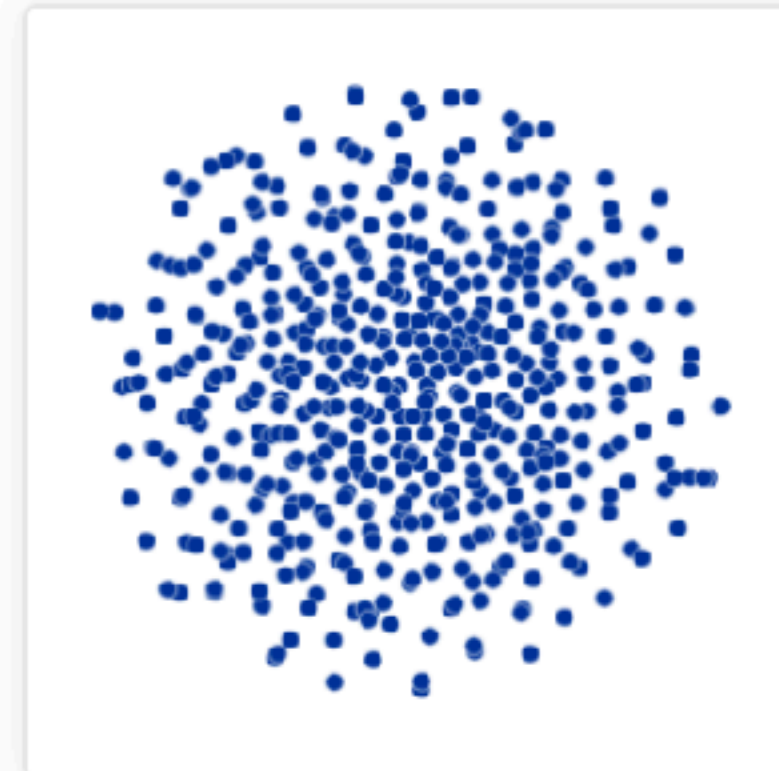
*Original*



Perplexity: 2  
Step: 5,000



Perplexity: 5  
Step: 5,000



Perplexity: 30  
Step: 5,000



Perplexity: 50  
Step: 5,000



Perplexity: 100  
Step: 5,000

[Wattenberg et al., Distill 2016]

# Dimensionality Reduction

Project nD data to 2D or 3D. Interpret/sanity check learned representations.

But **subject to their own interpretation issues**.

Different techniques make different trade-offs:

**PCA (Principal Component Analysis)**: roughly fit a p-dimensional ellipsoid to the data, order axes by amount of data variance they explain.

**Preserves global structure.**

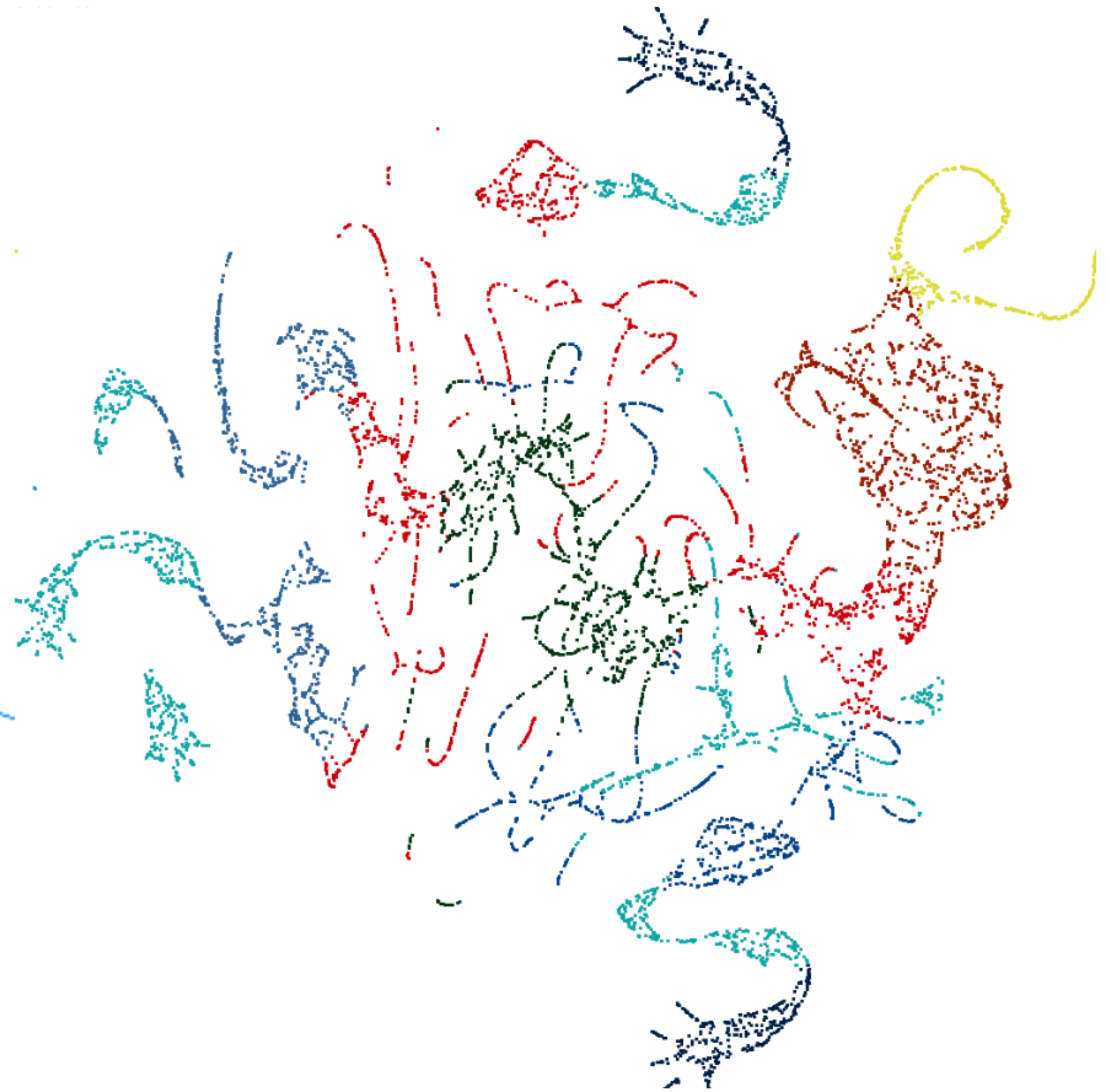
**t-SNE (t-Dist. Stochastic Neighbor Embedding)**: probabilistic distribution that adapts and performs different transformations on different regions.

**UMAP (Uniform Manifold Approx. & Projection)**: Identify local regions, stitch them together. **Tries to balance local/global trade-off.**

Original 3D Data

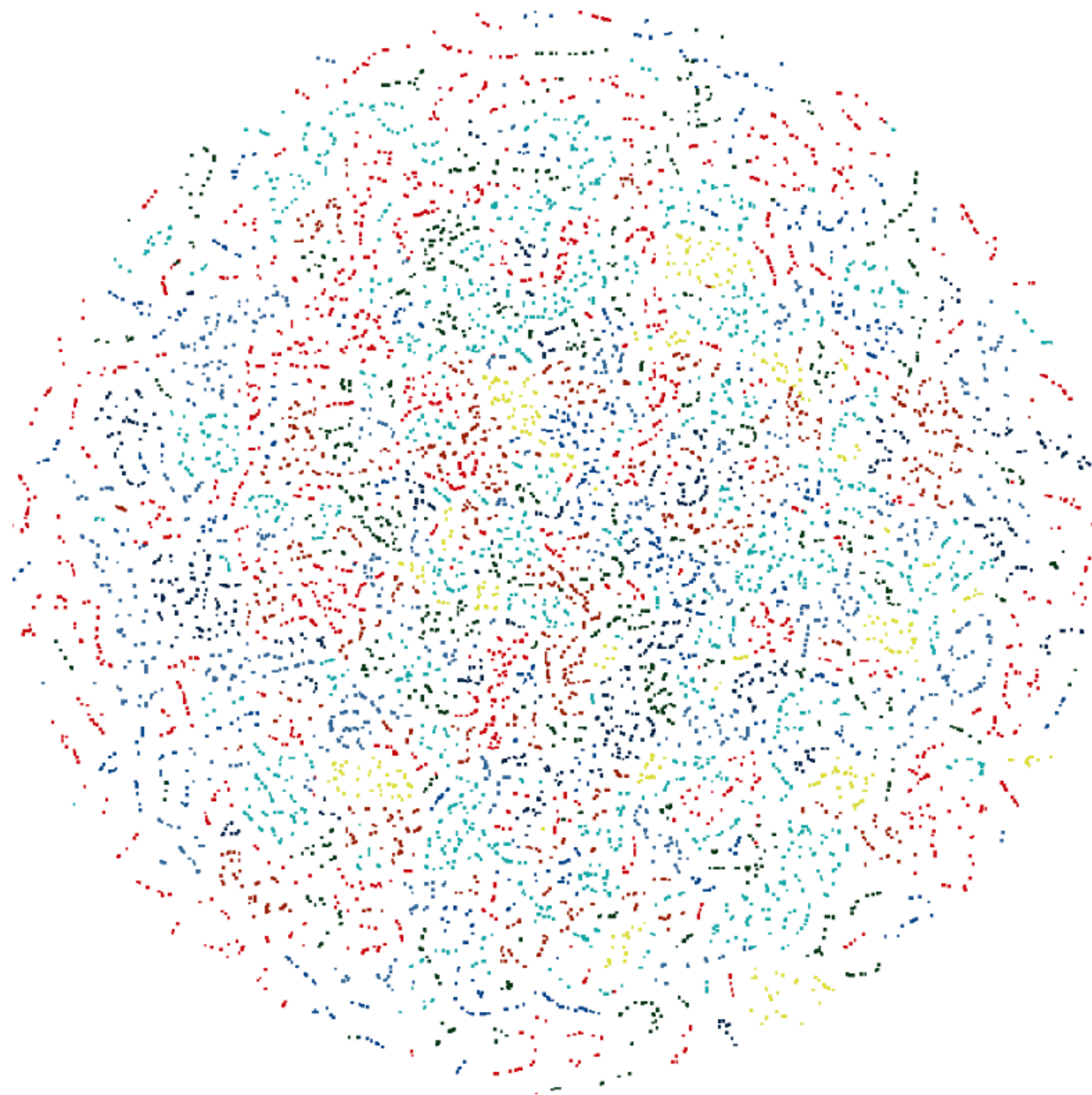


2D UMAP Projection



[Coenen & Pearce, 2019]

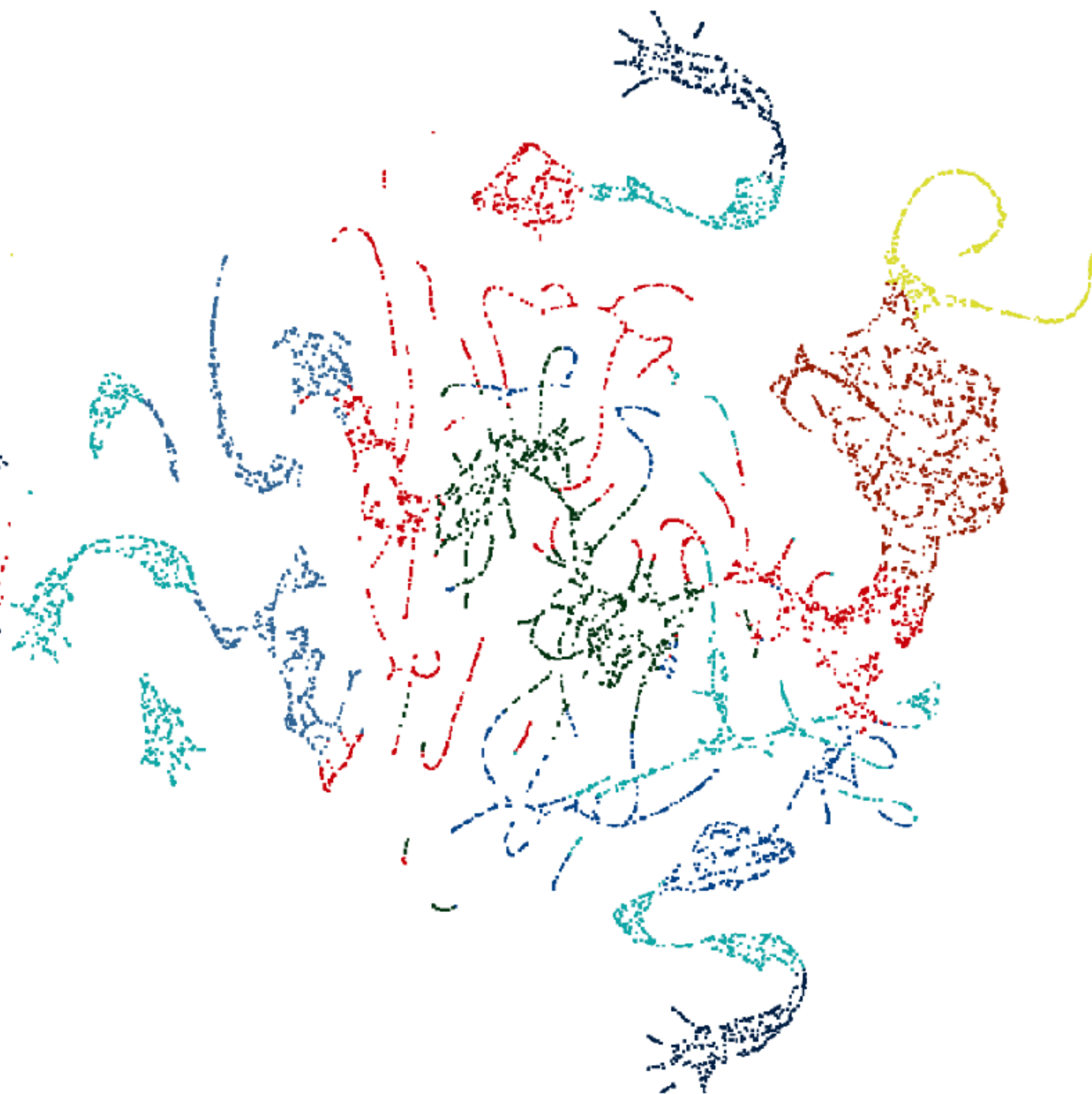
2D t-SNE projection



perplexity: 5  
time: 9m 18s



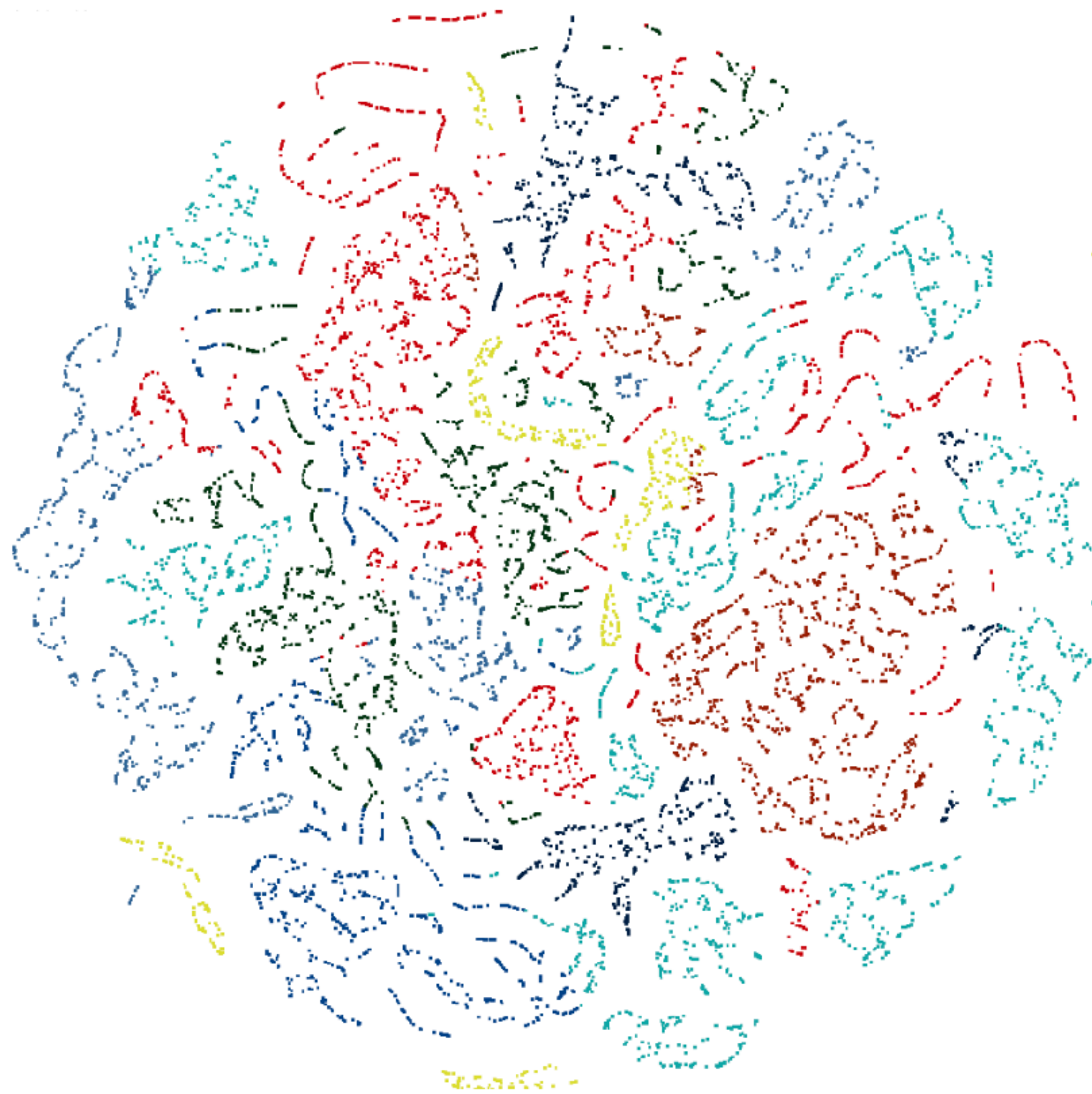
2D UMAP projection



n\_neighbors: 15  
min\_dist: 0.1  
time: 1m 2s



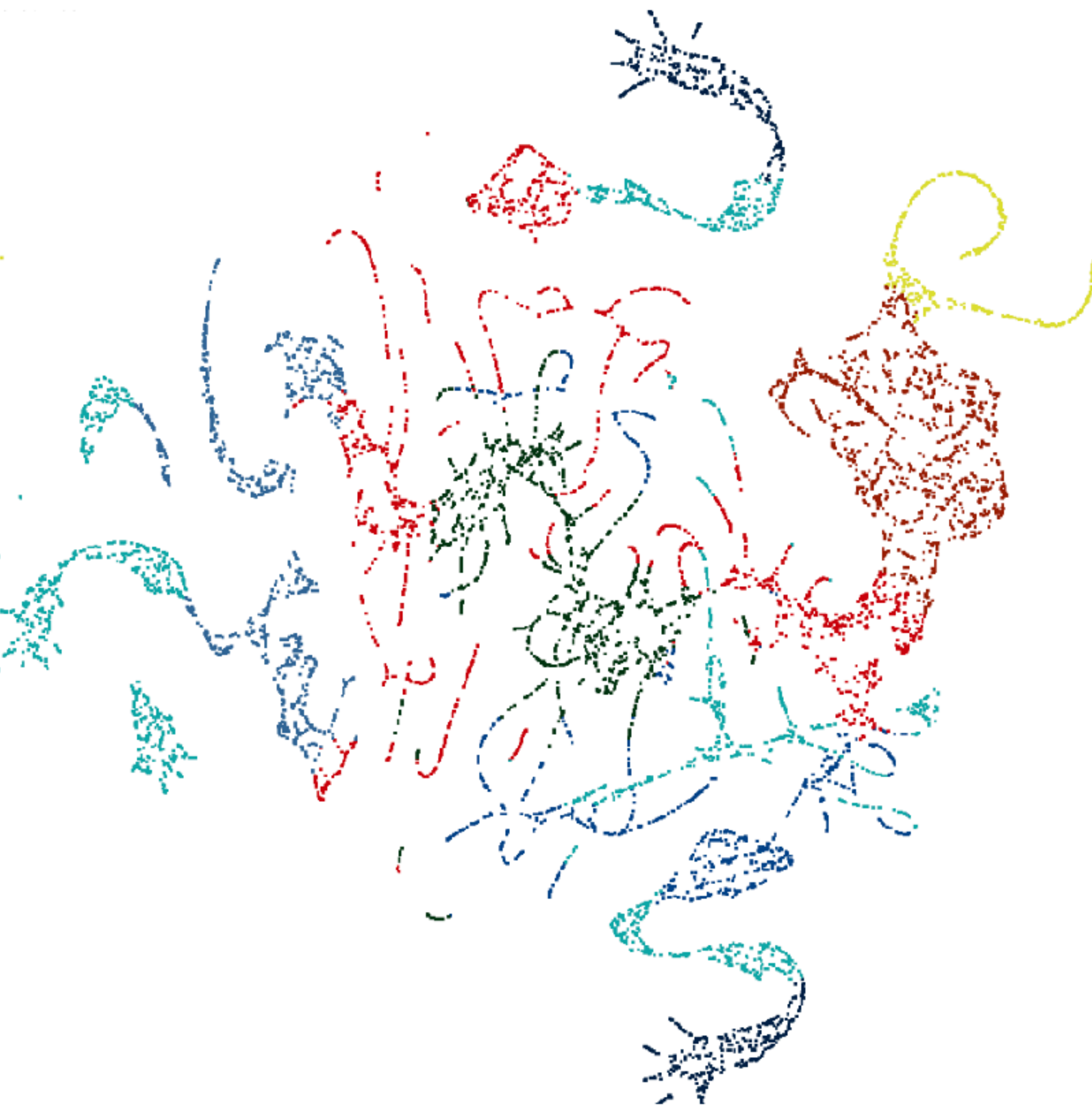
2D t-SNE projection



perplexity: 30  
time: 11m 5s



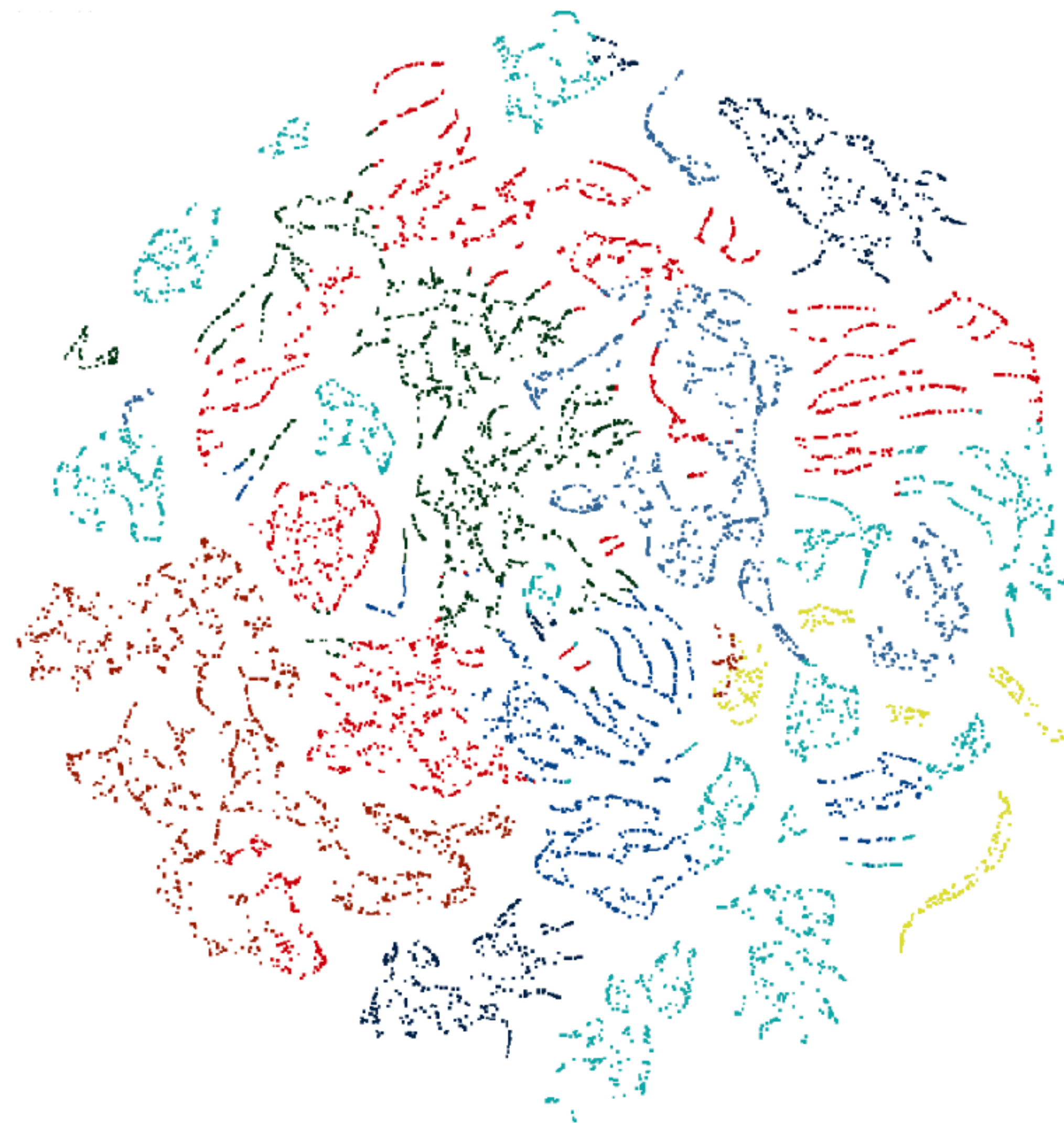
2D UMAP projection



n\_neighbors: 15  
min\_dist: 0.1  
time: 1m 2s



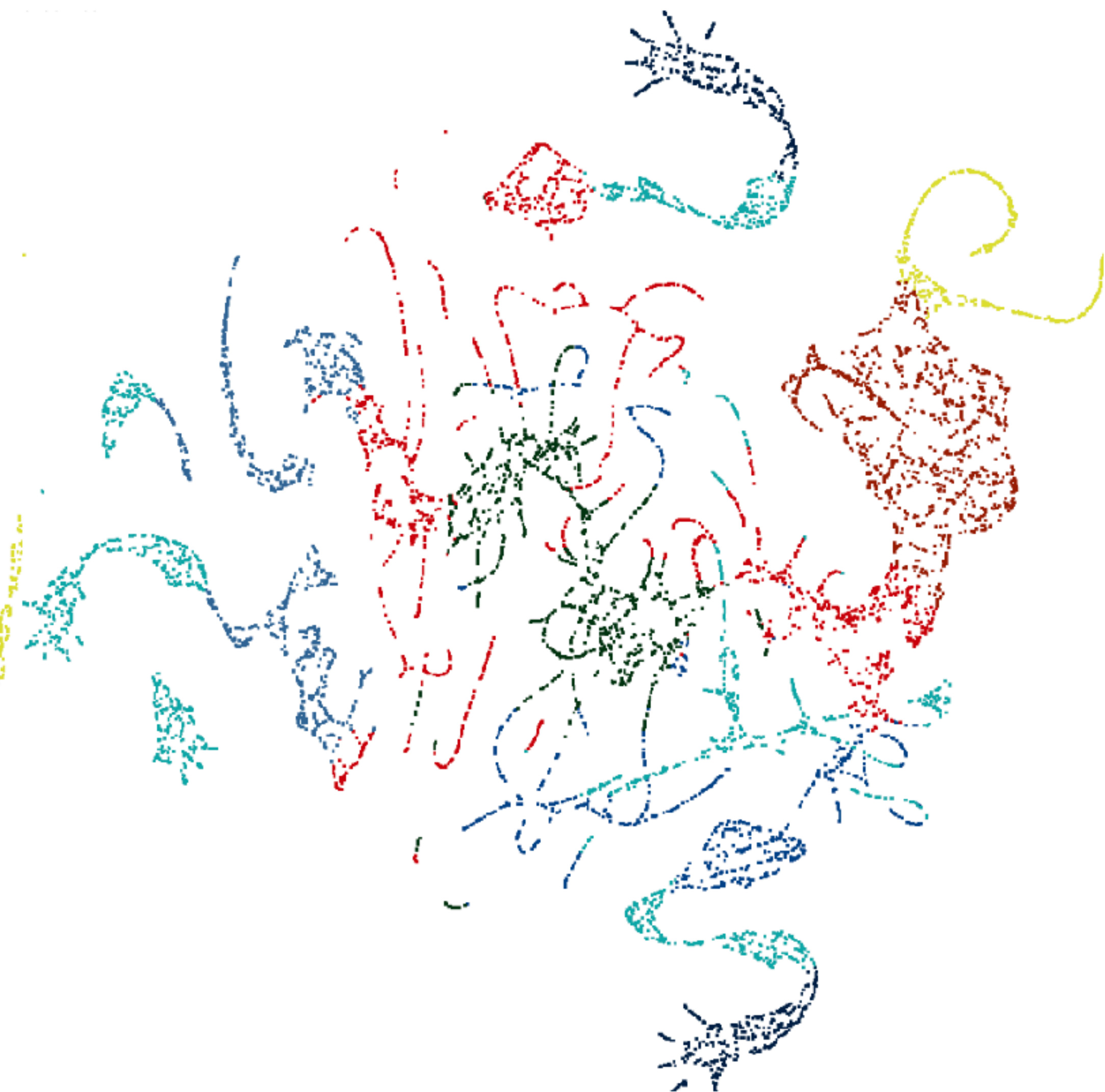
2D t-SNE projection



perplexity: 50  
time: 12m 42s



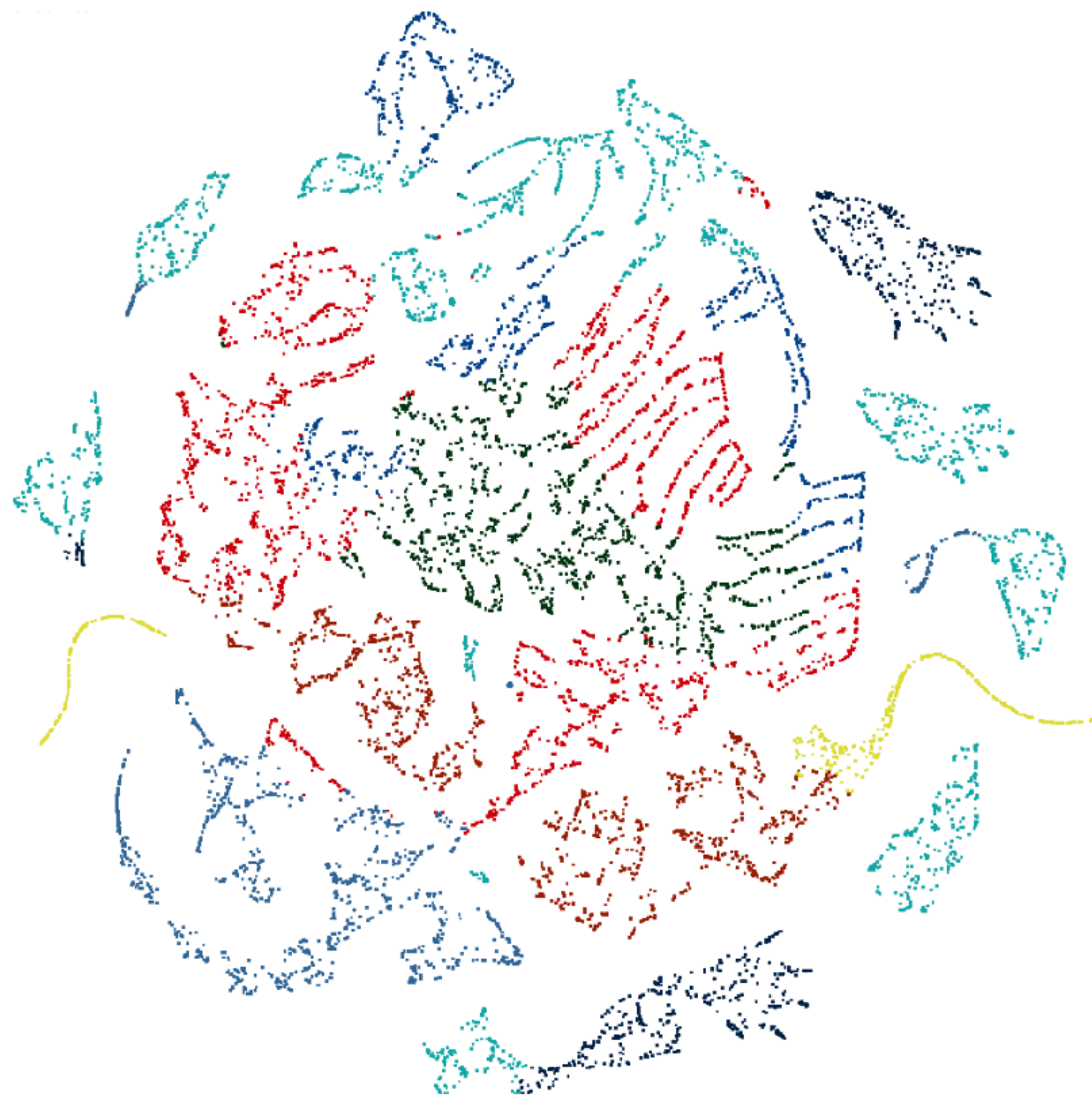
2D UMAP projection



n\_neighbors: 15  
min\_dist: 0.1  
time: 1m 2s



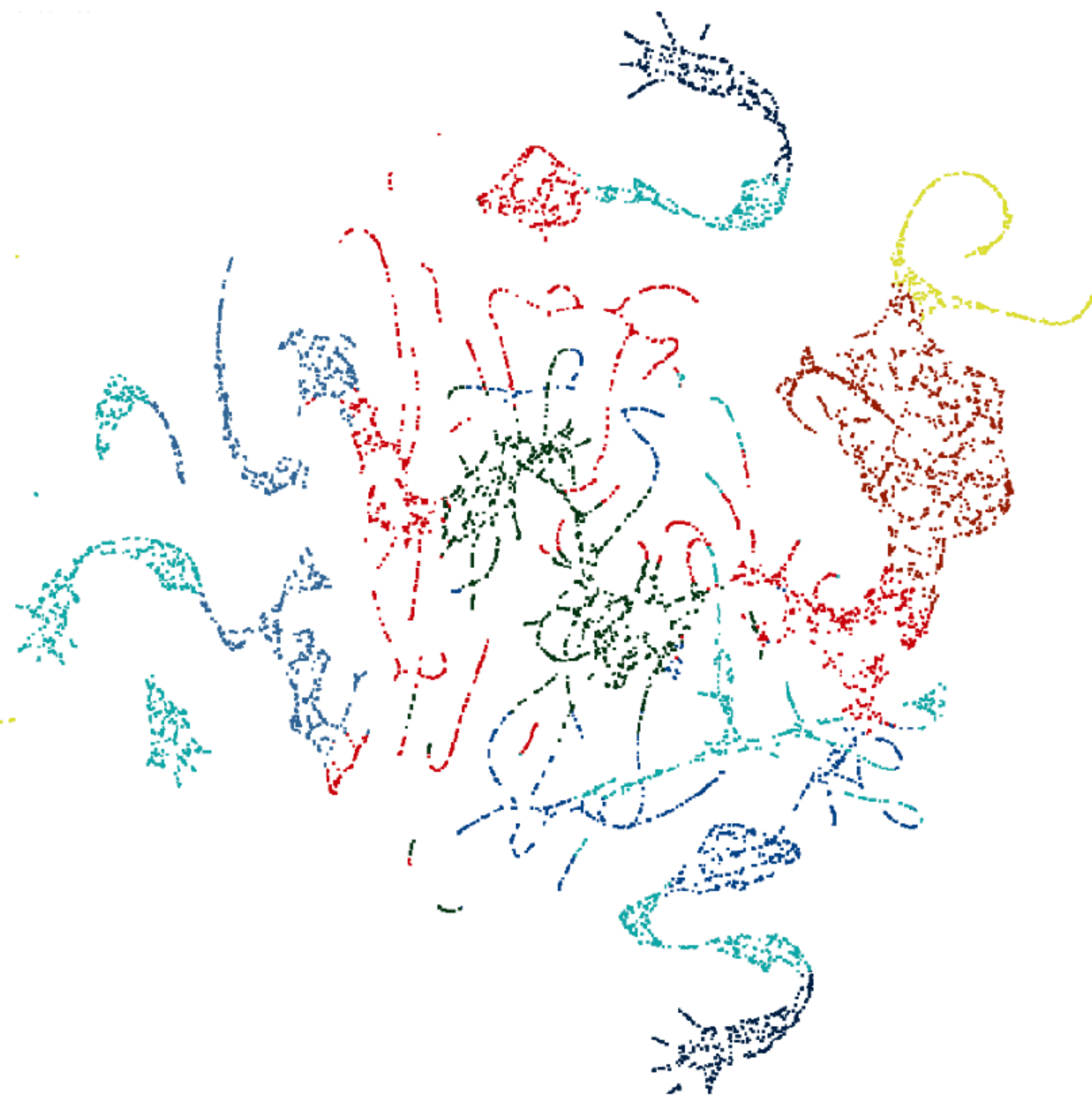
2D t-SNE projection



perplexity: 100  
time: 16m 1s



2D UMAP projection



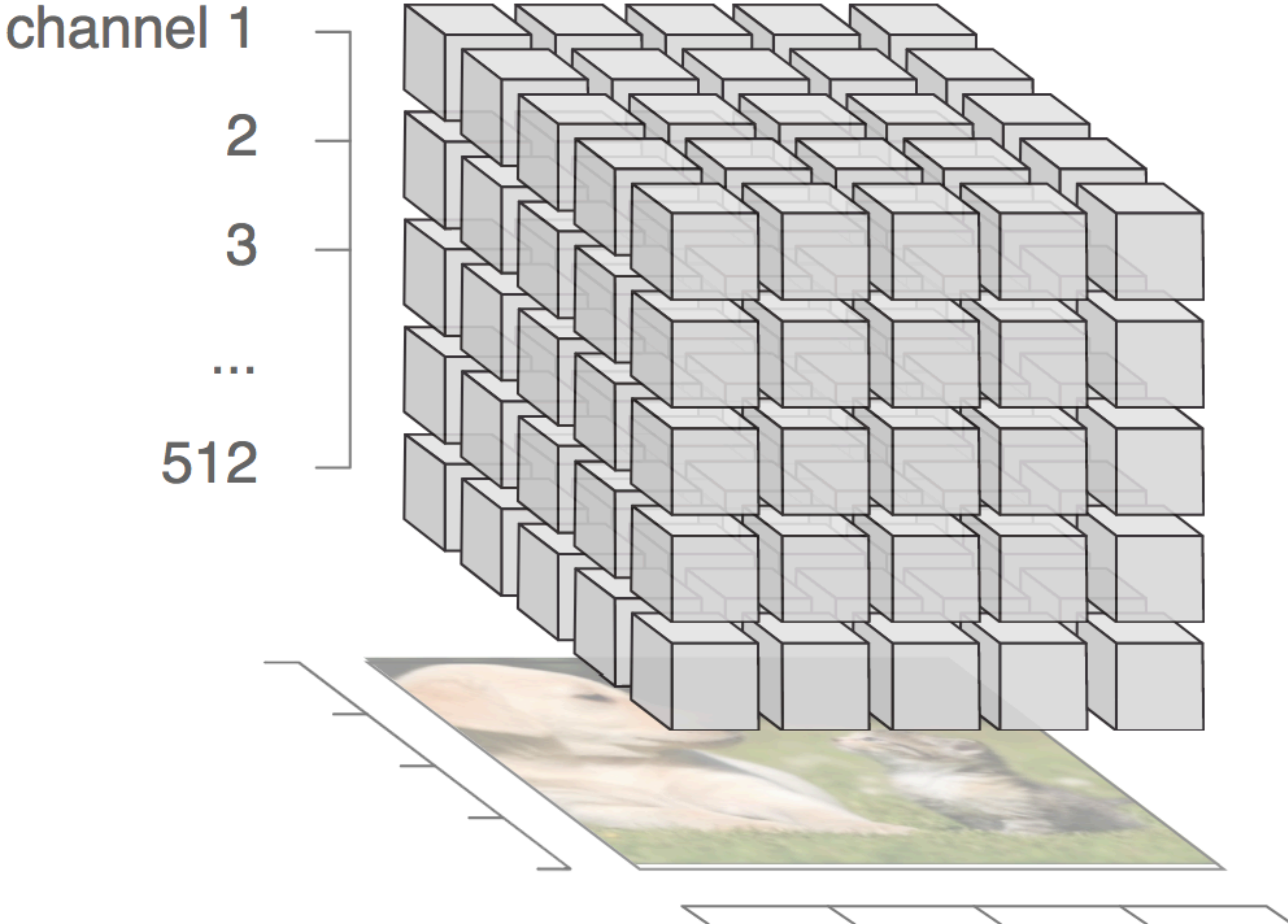
n\_neighbors: 15  
min\_dist: 0.1  
time: 1m 2s



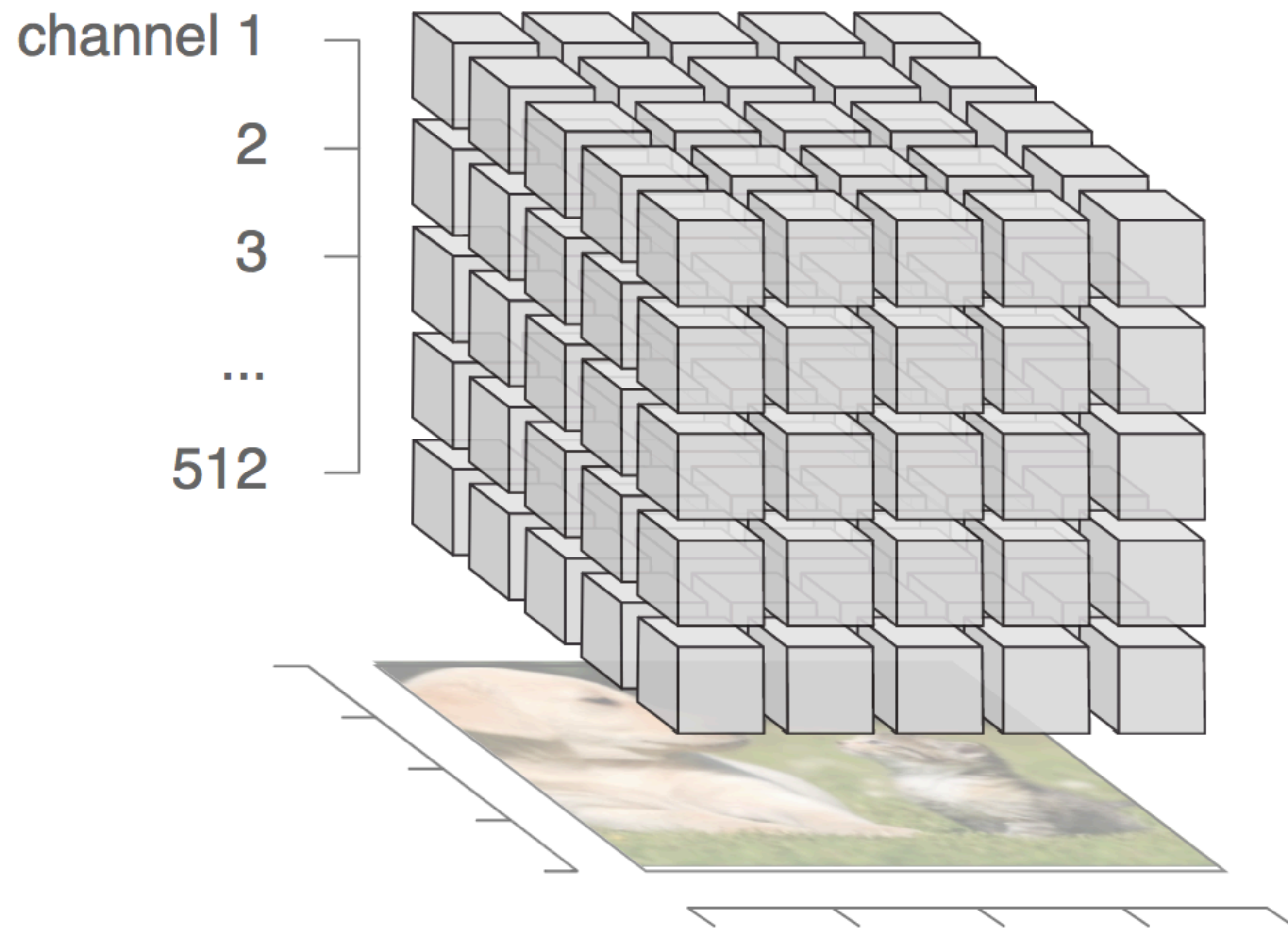




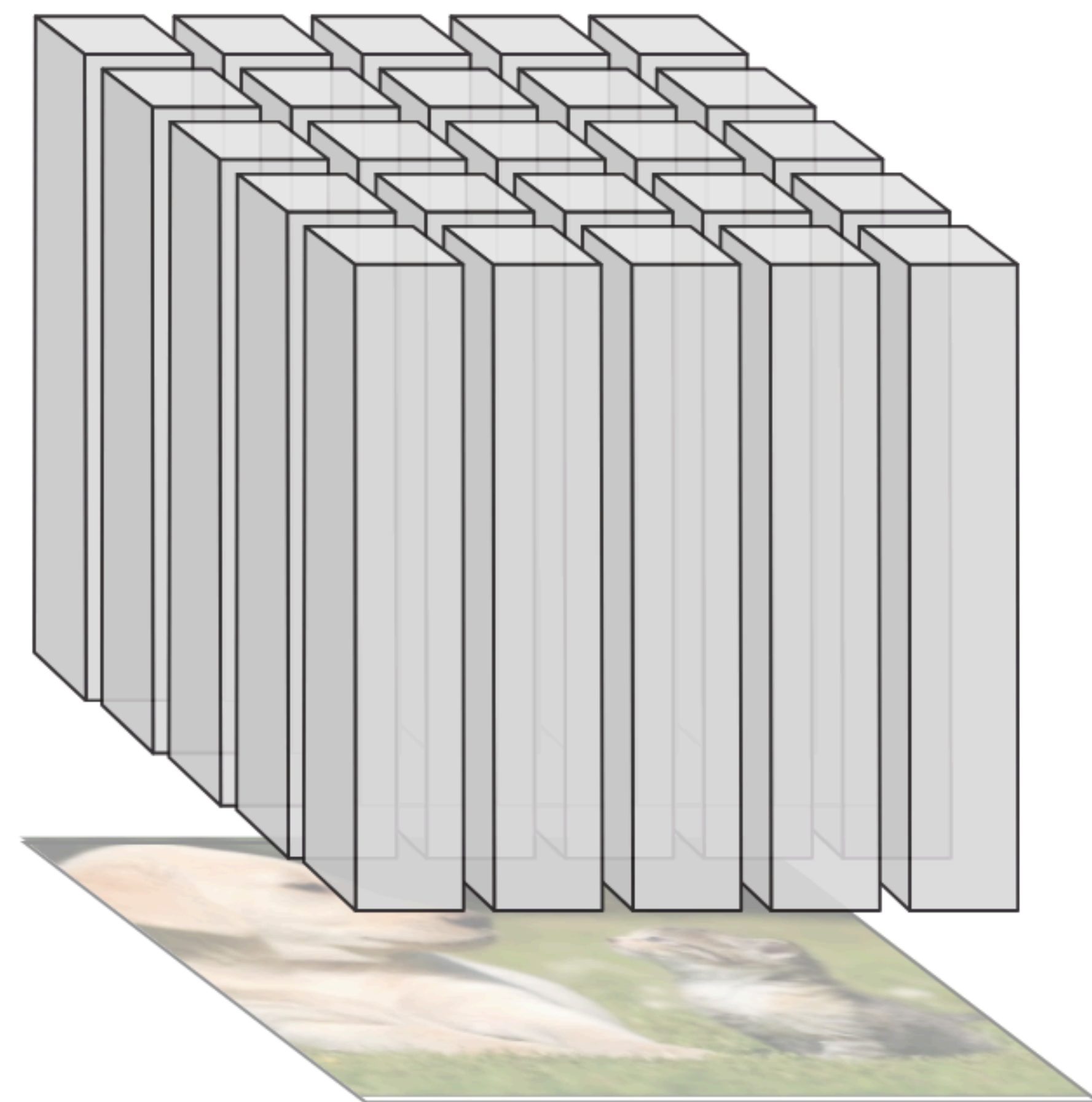
# Individual Neurons



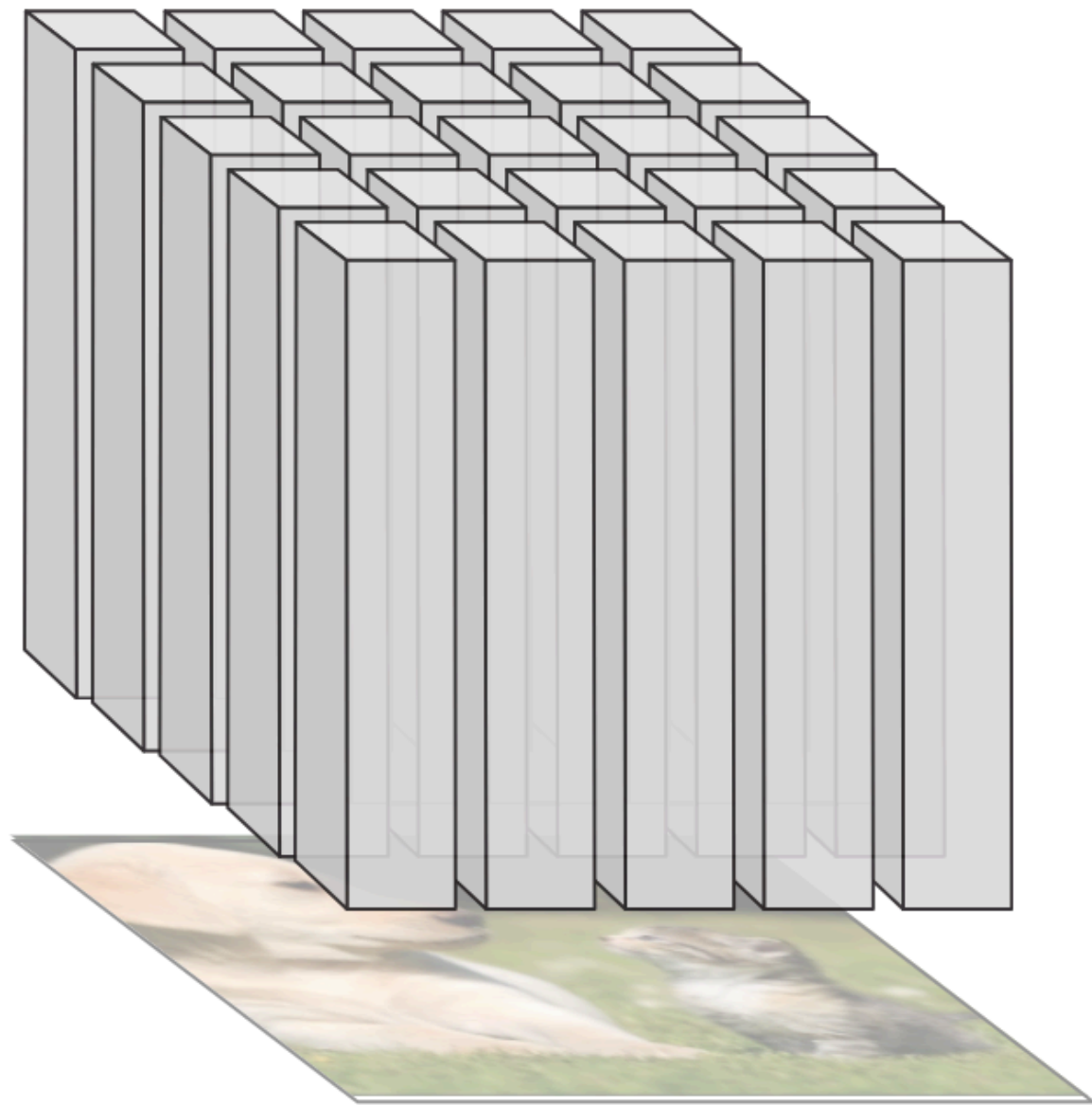
## Individual Neurons



## Spatial Activations

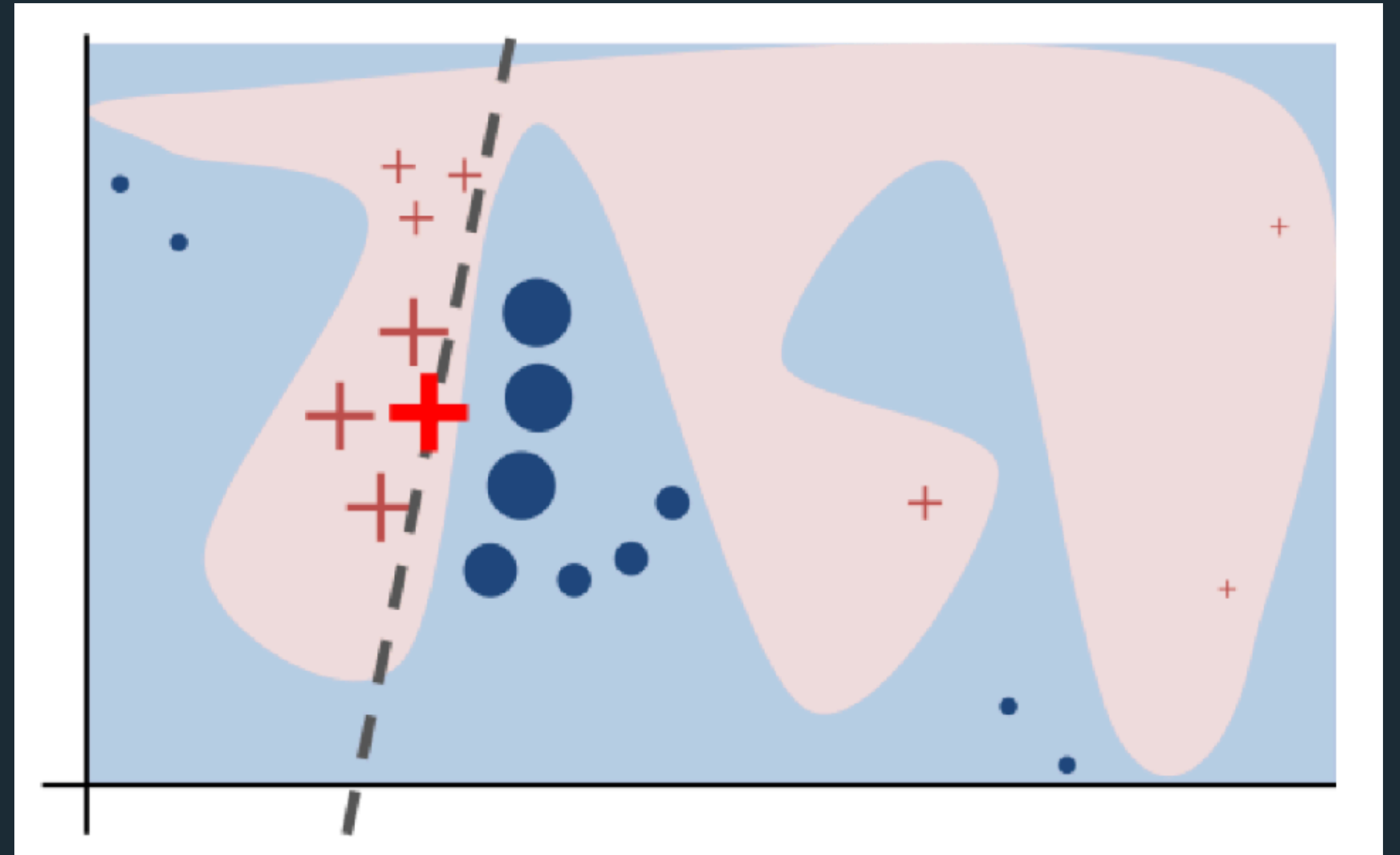


## Spatial Activations



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

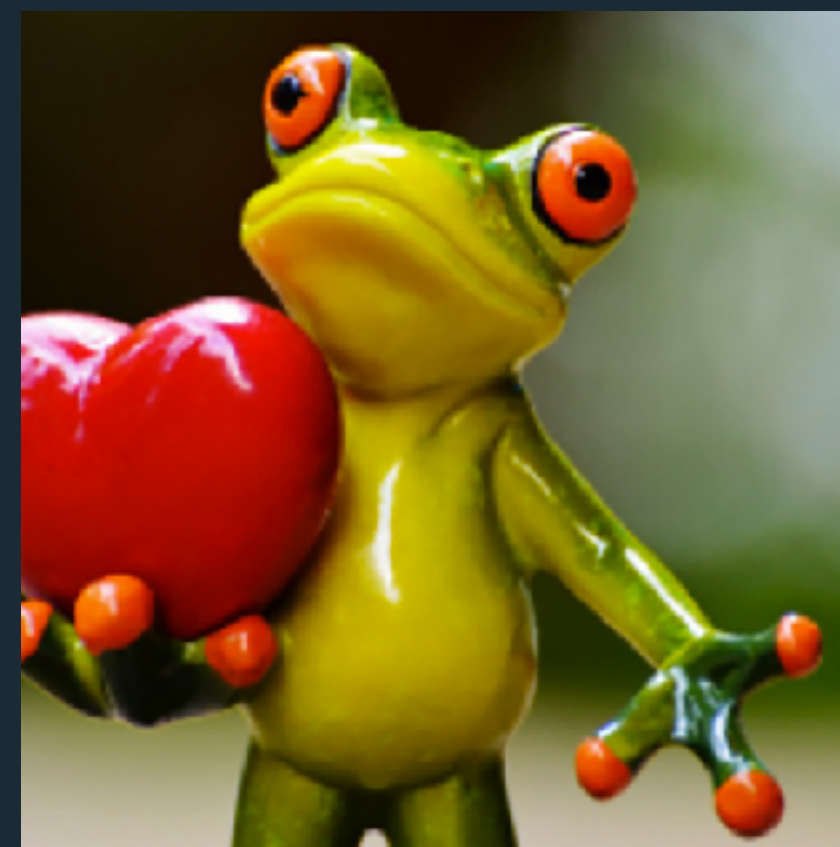
**L**ocal  
**I**nterpretable  
**M**odel-Agnostic  
**E**xplanations



[Ribeiro et al., KDD 2016]

# Local Interpretable Model-Agnostic Explanations

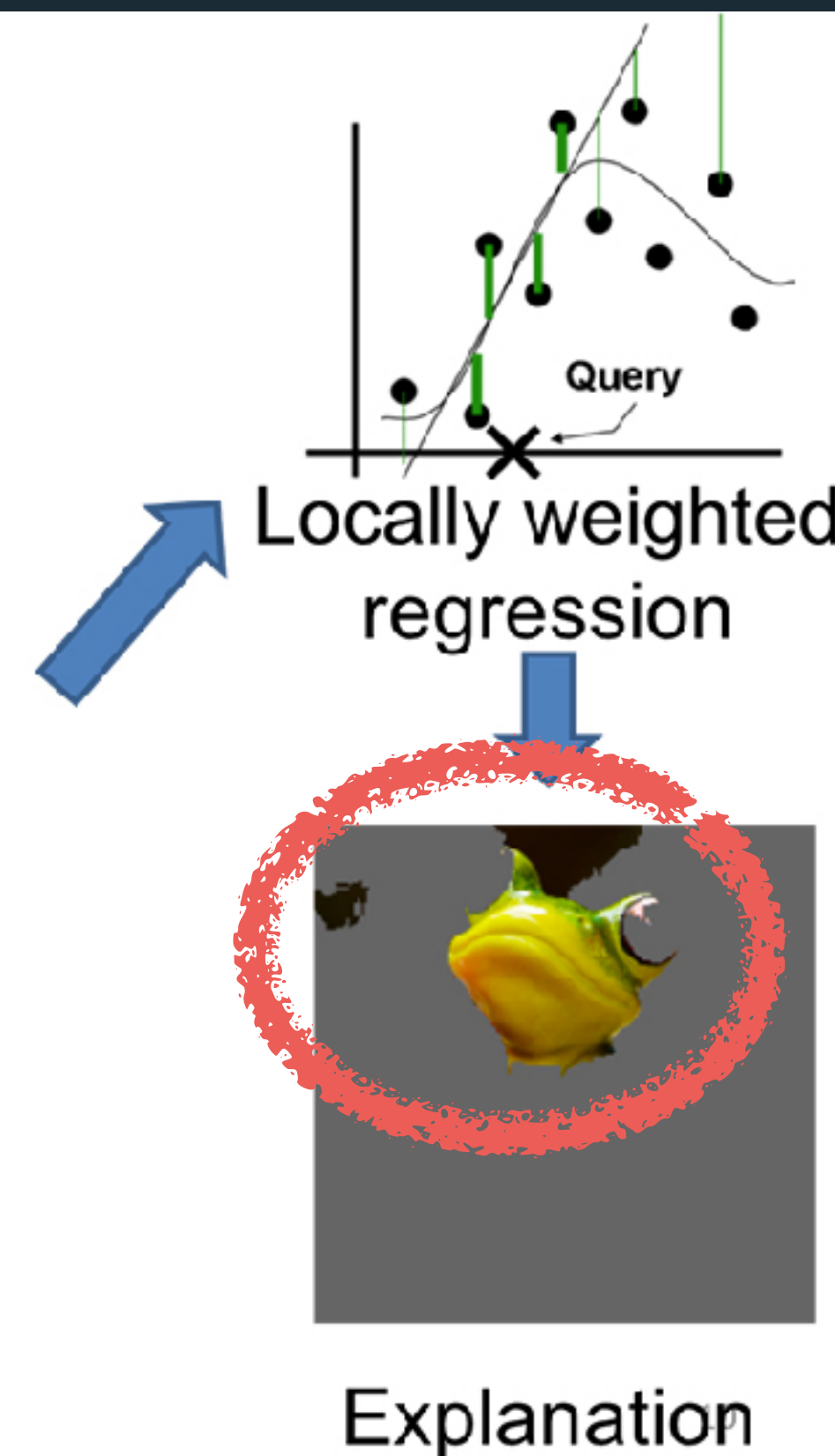
[Ribeiro et al., KDD 2016]



Identify subcomponents



Perturbed Instances	P(tree frog)
A fragmented version of the tree frog image where parts are scattered.	 0.85
A version of the tree frog image where only individual pixels are scattered, not connected.	 0.00001
The original, intact tree frog image.	 0.52



Regions sufficient for "tree frog" classification.

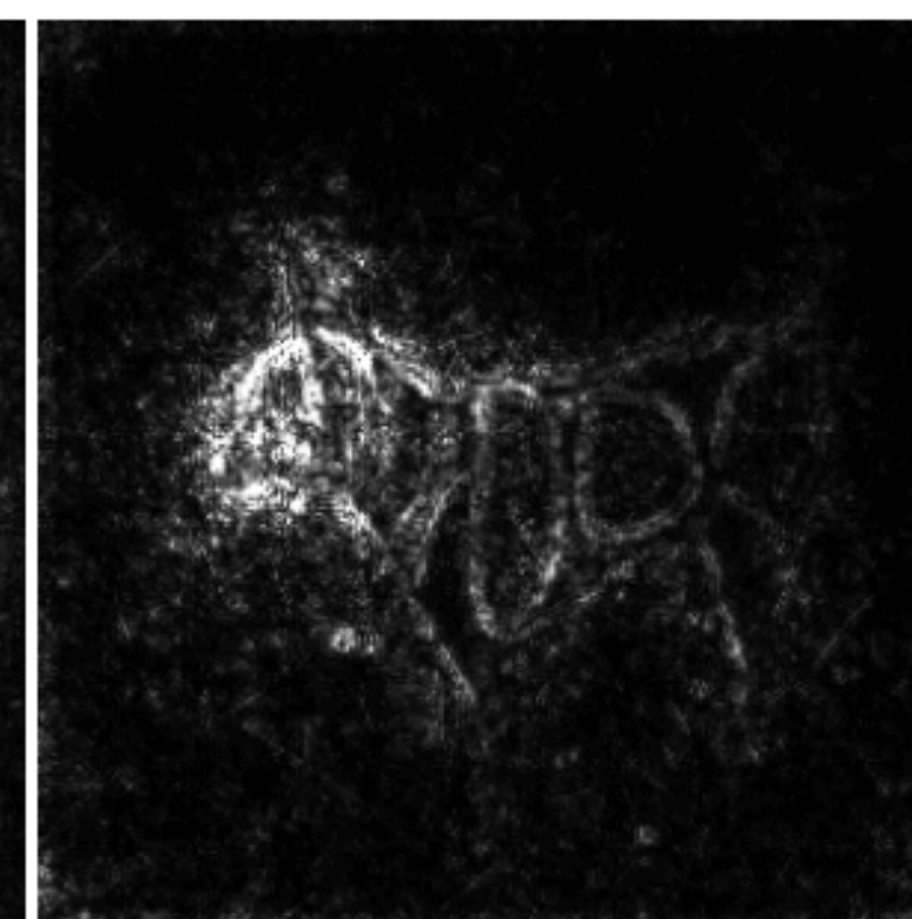
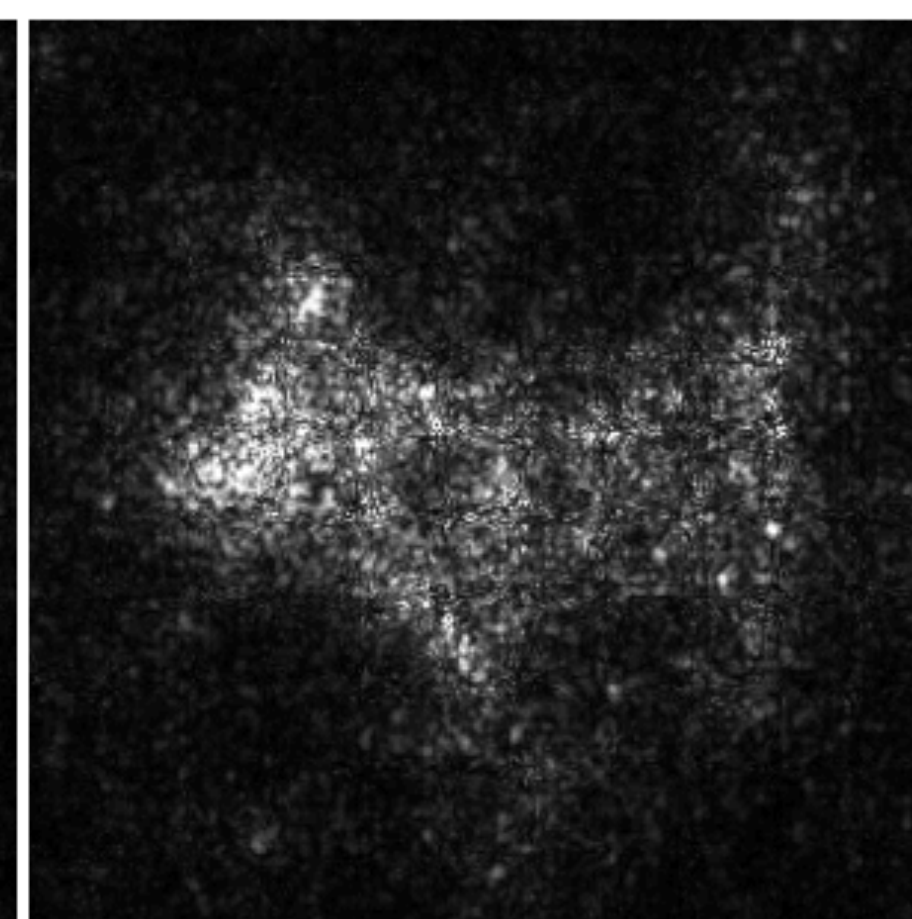
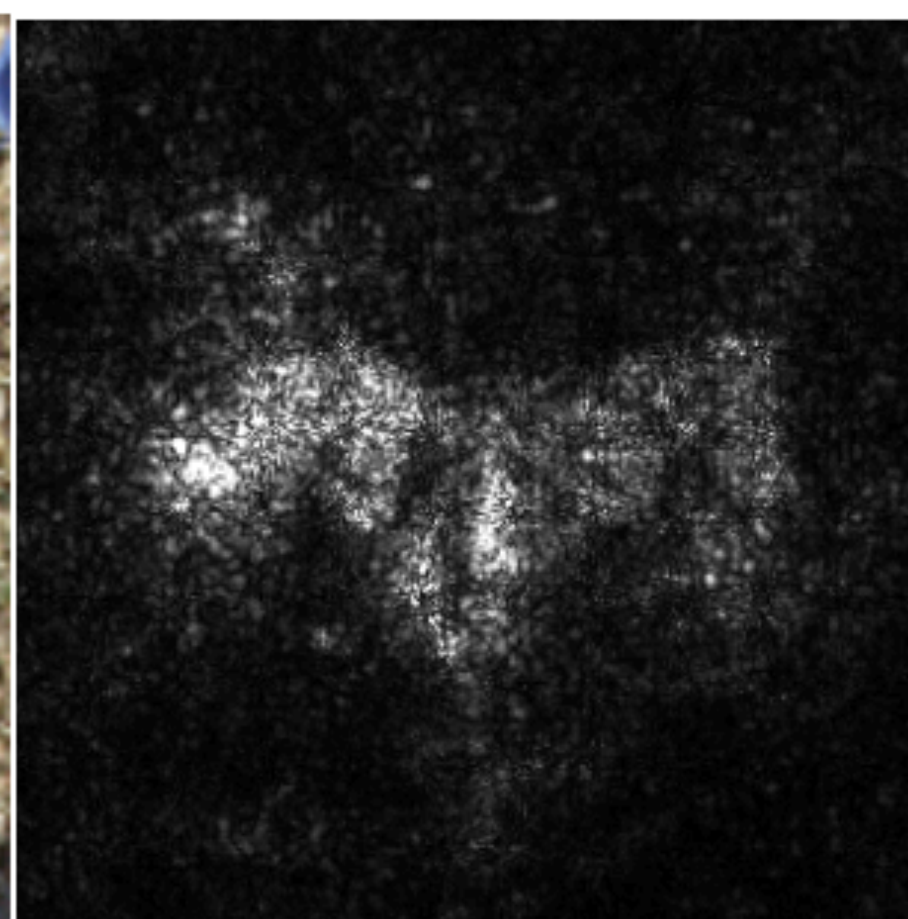
Image

Gradient

Integrated

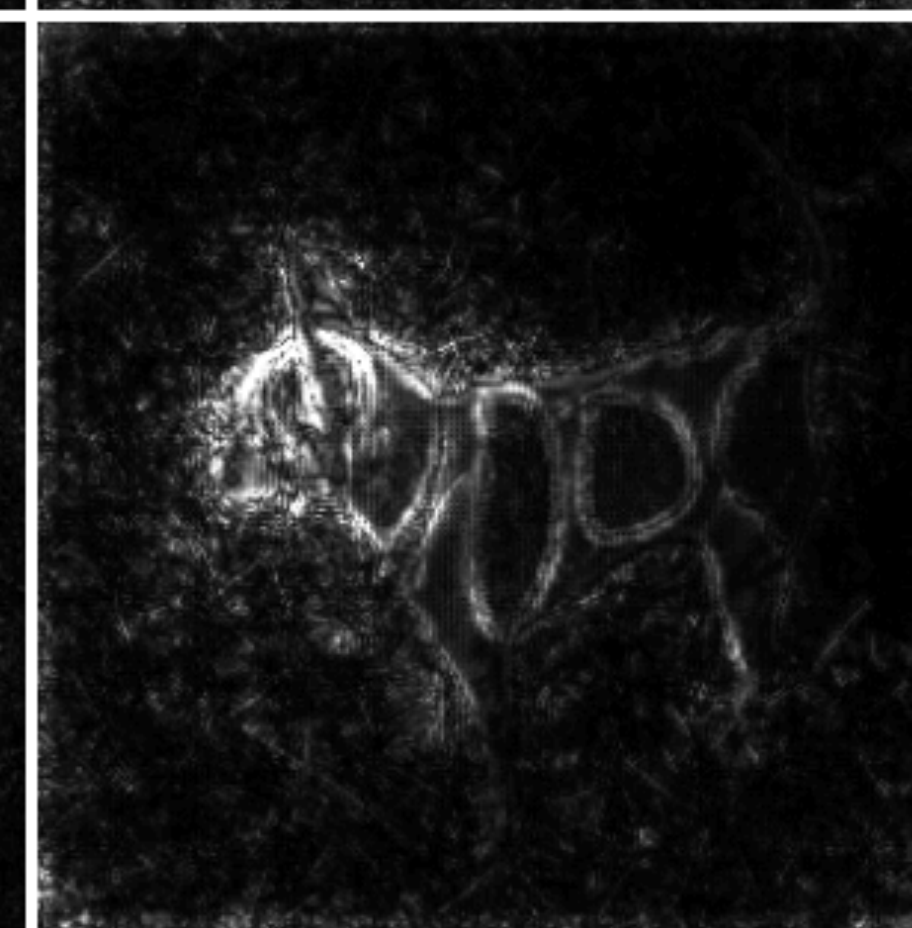
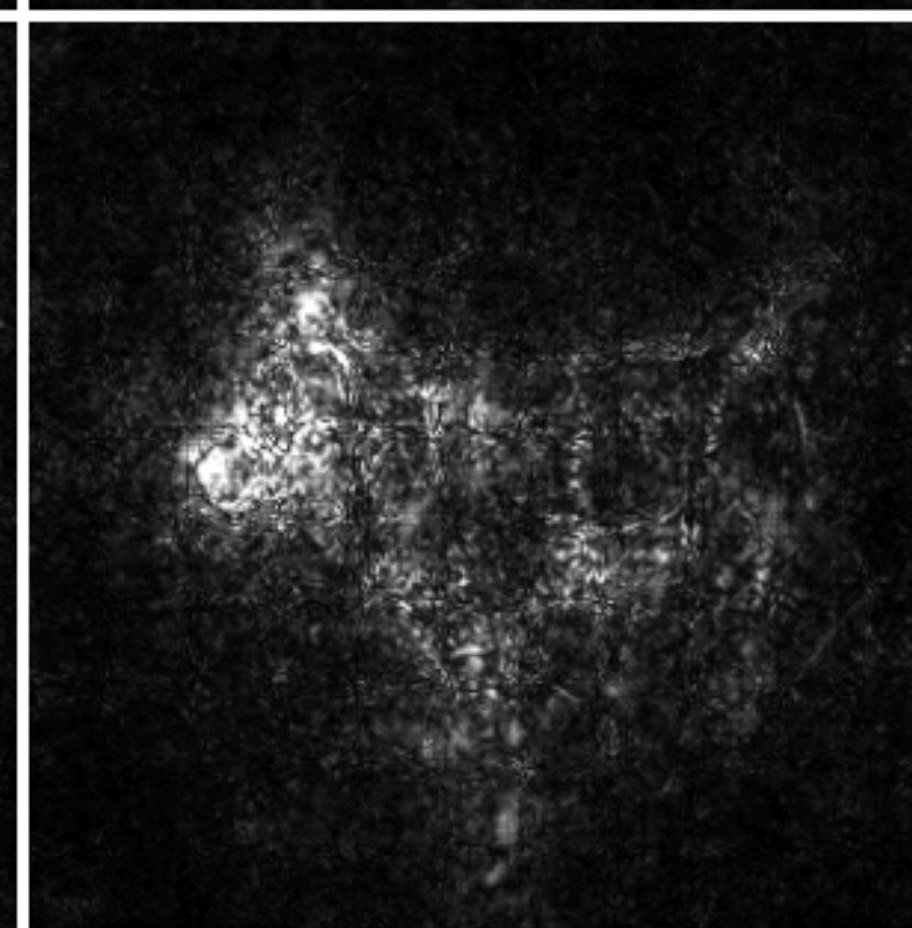
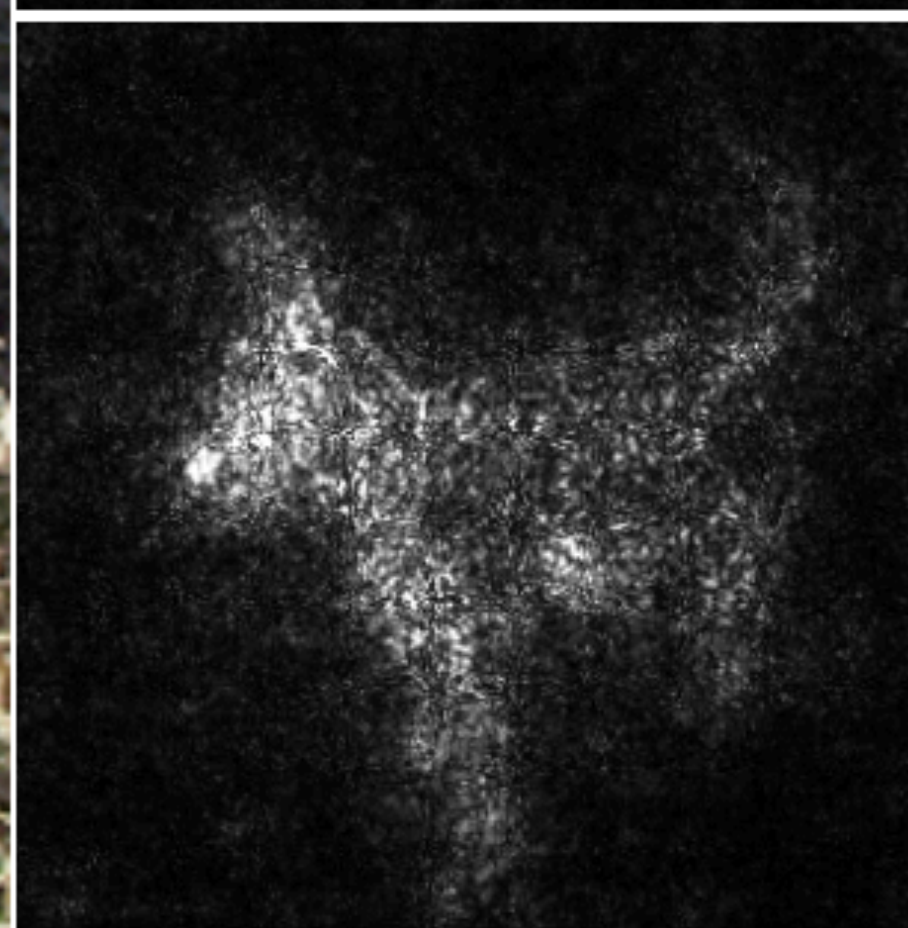
Guided Backprop

*Label: toy terrier*



Plain

SmoothGrad

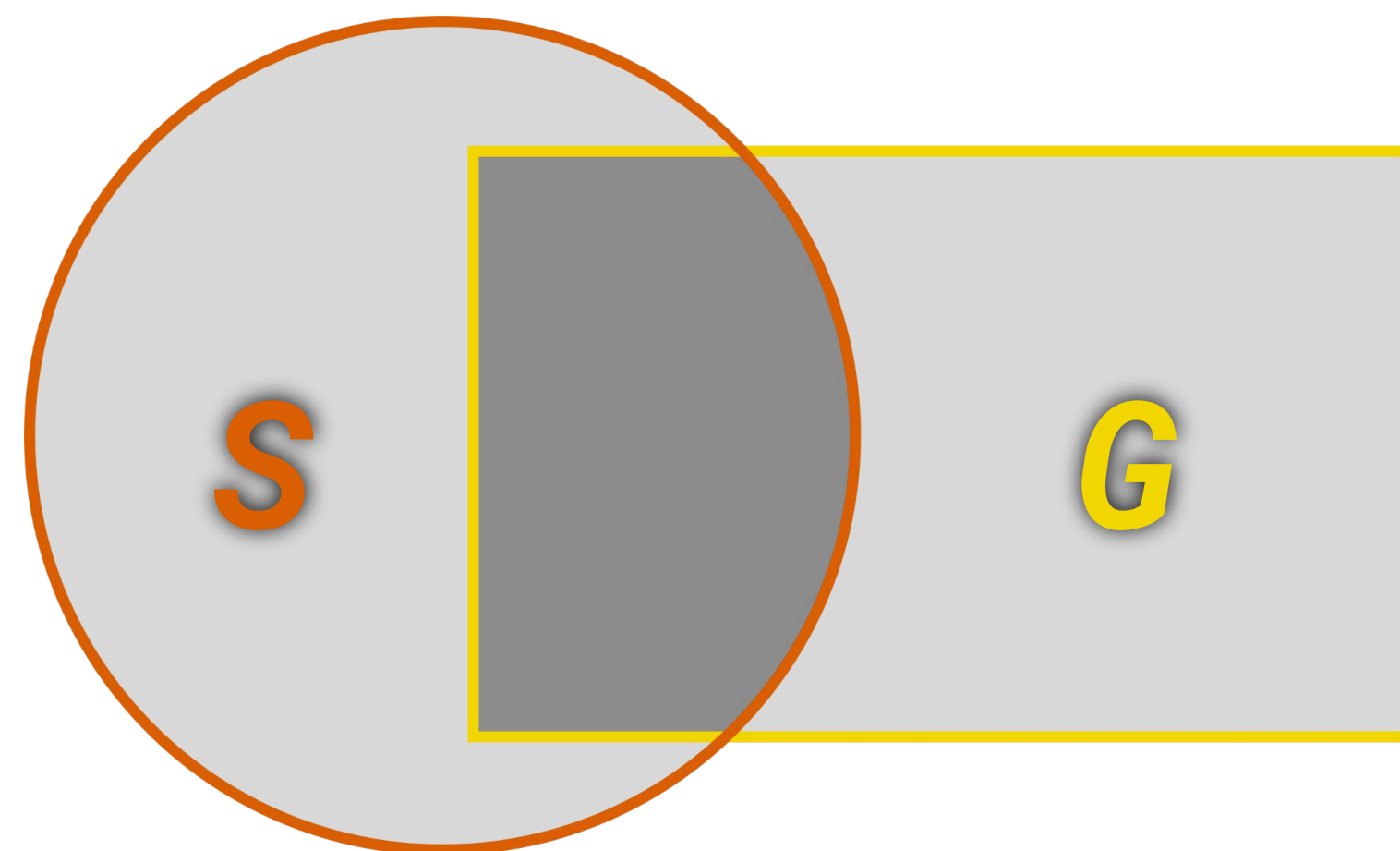


42

[SmoothGrad. Smilkov et al., 2017]

# Shared Interest: Measuring Human-AI Alignment

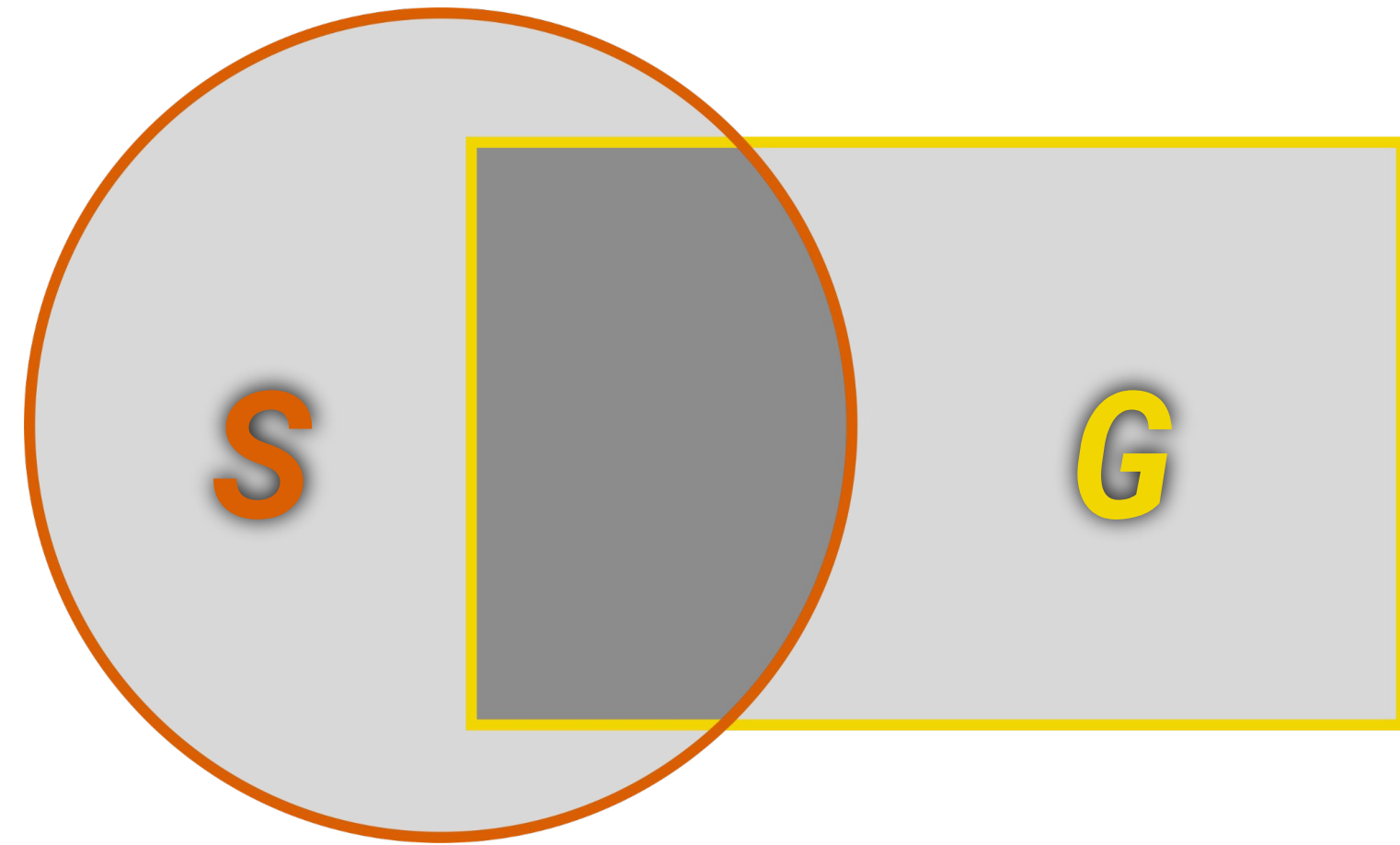
with



**Saliency  
Features**

**Ground Truth  
Features**  
(Human Annotation)

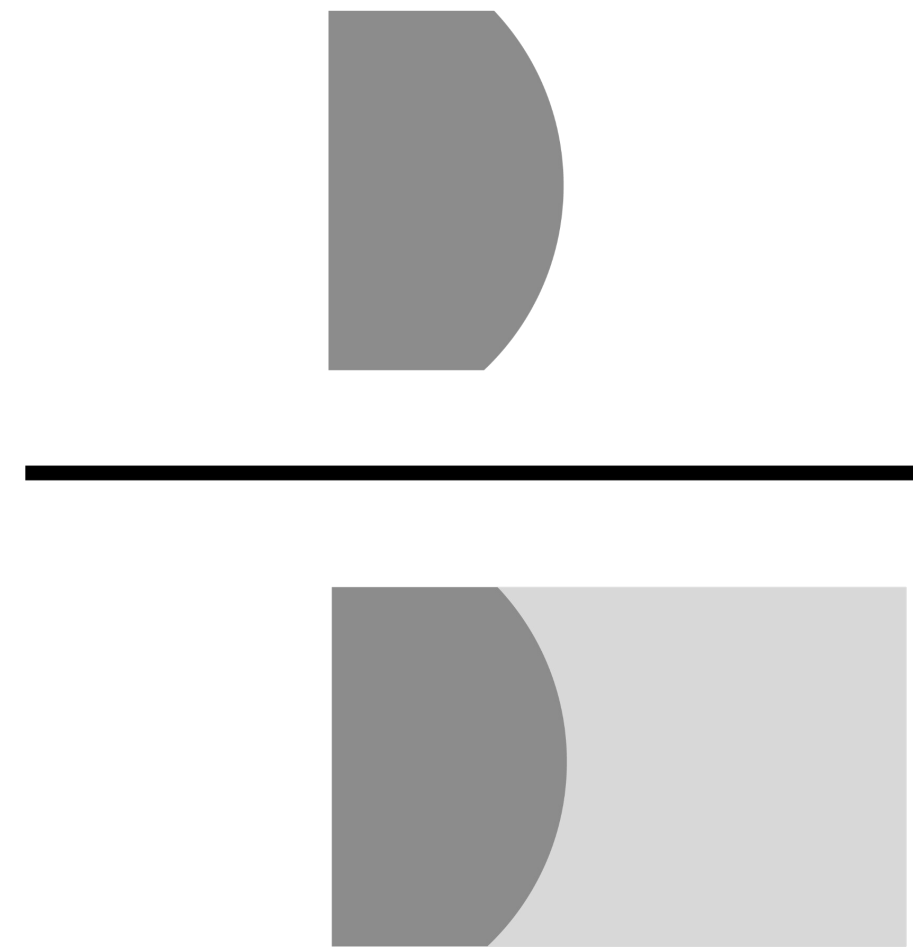
# Shared Interest: Measuring Human-AI Alignment



**Saliency Features**

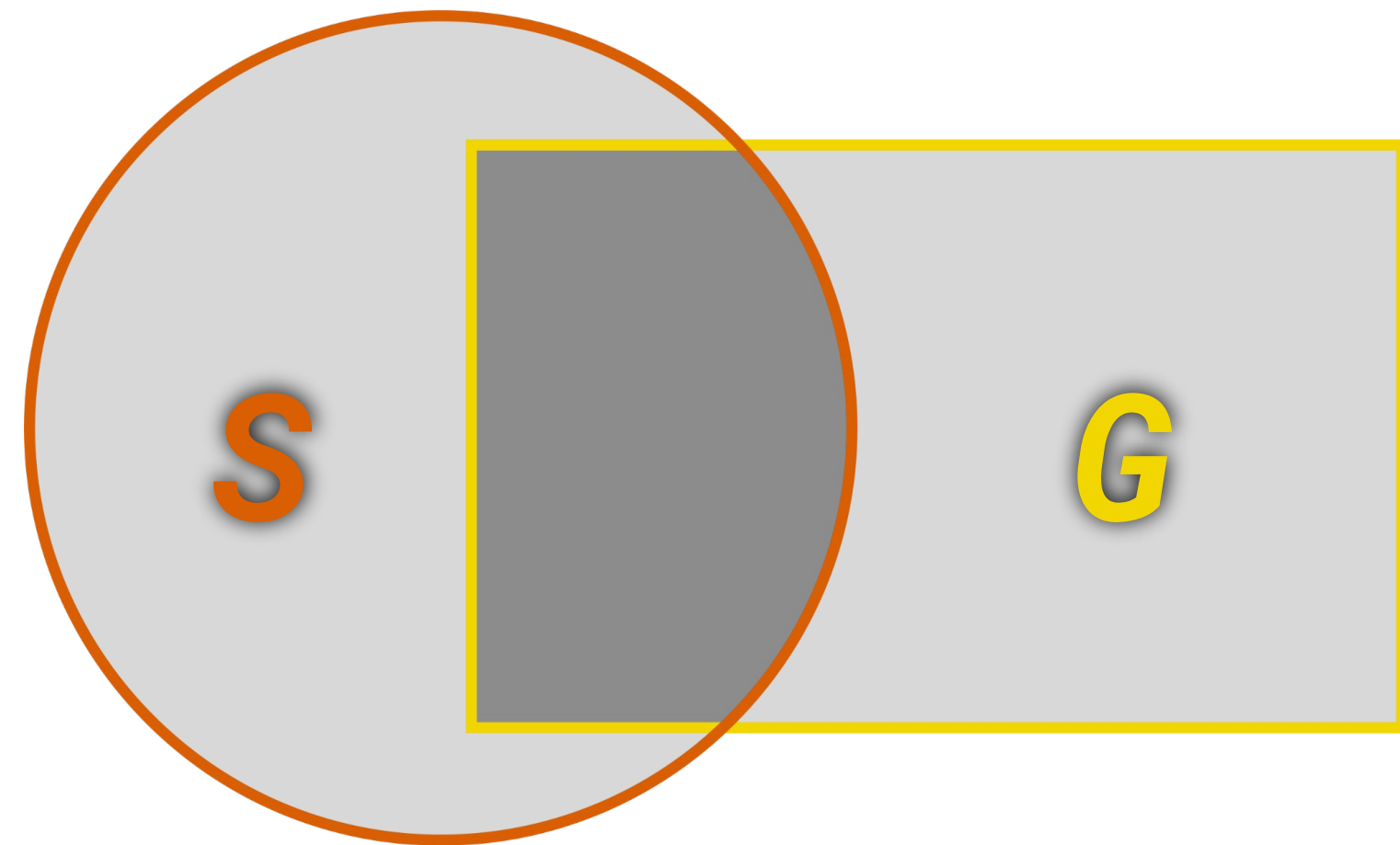
**Ground Truth Features**  
(Human Annotation)

**Ground Truth Coverage**





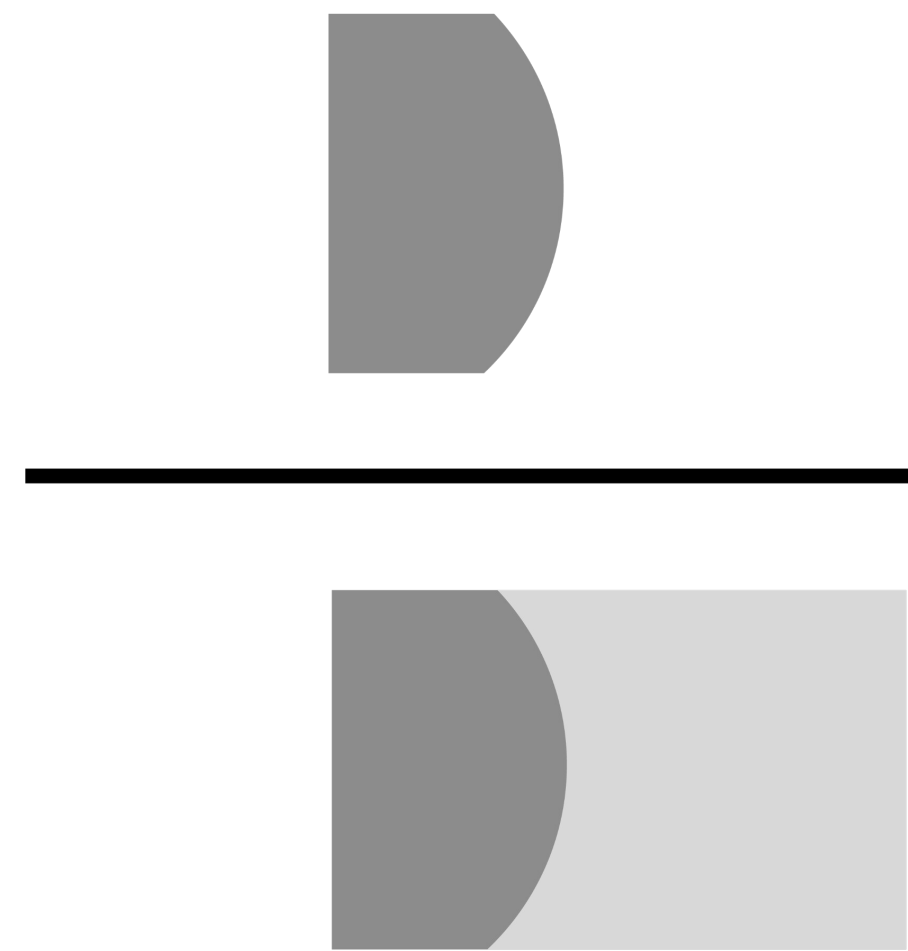
# Shared Interest: Measuring Human-AI Alignment



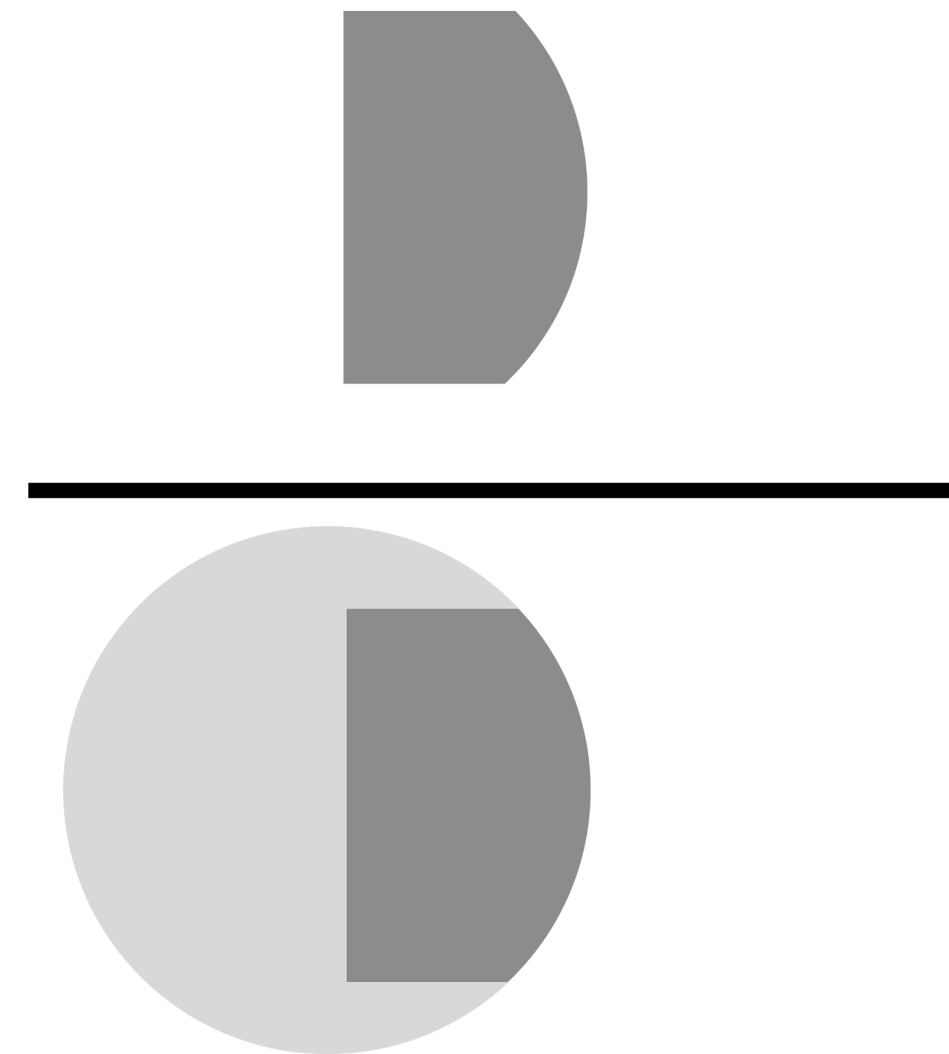
**Saliency Features**

**Ground Truth Features**  
(Human Annotation)

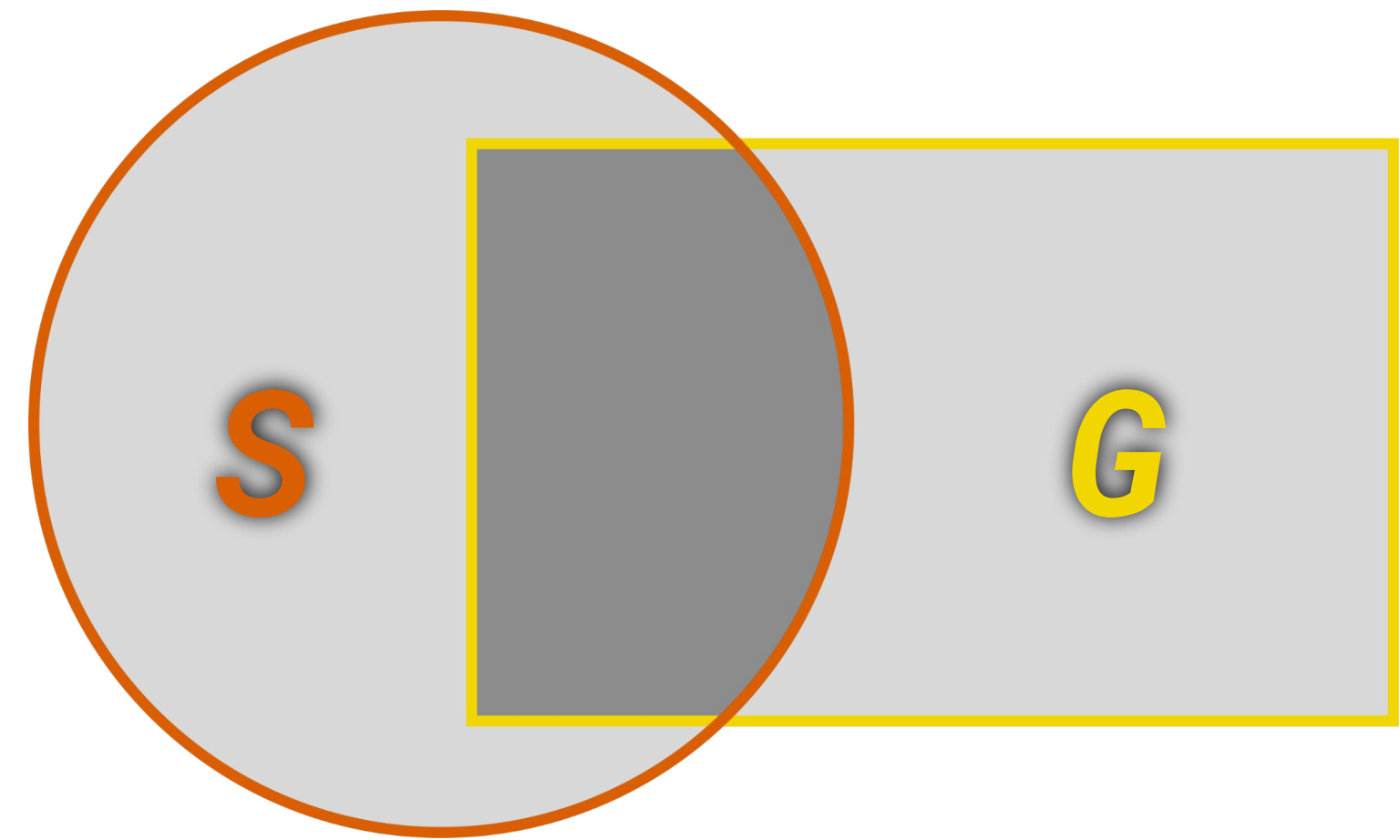
**Ground Truth Coverage**



**Saliency Coverage**



# Shared Interest: Measuring Human-AI Alignment



**Saliency Features**

**Ground Truth Features**  
(Human Annotation)

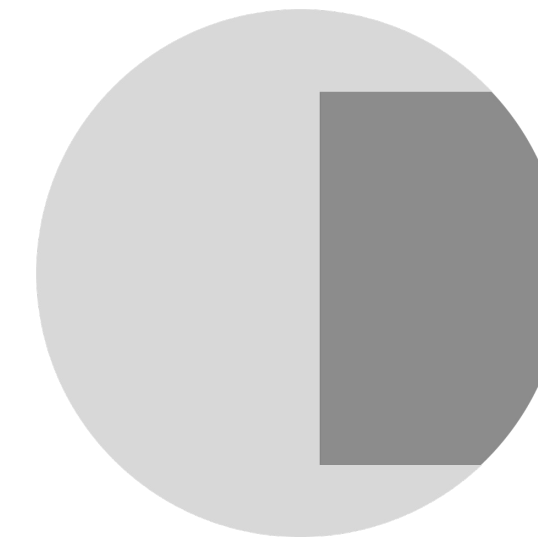
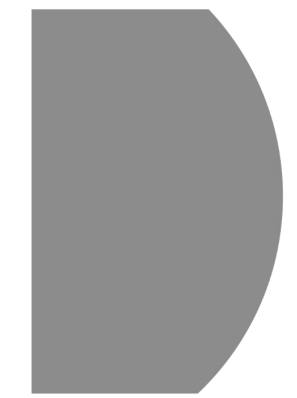
**Ground Truth Coverage**



**Saliency Coverage**



**IoU Coverage**



# Shared Interest: Measuring Human-AI Alignment



## LOW SHARED INTEREST SCORE

Incorrect

Correct

## HIGH SHARED INTEREST SCORE

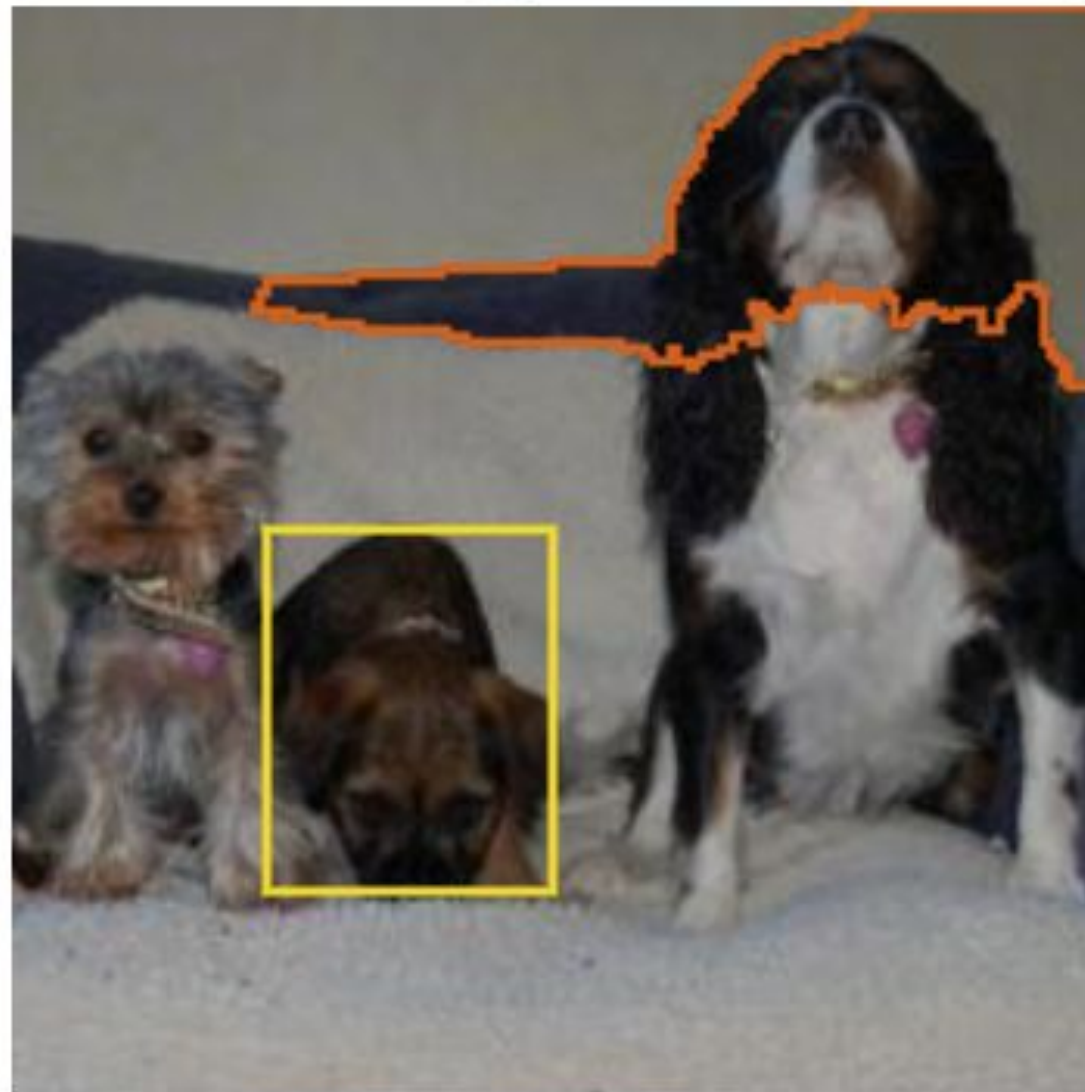
Incorrect

Correct

GROUND TRUTH COVERAGE

IoU: 0.00 GTC: 0.00 SC: 0.00

border terrier blenheim spaniel



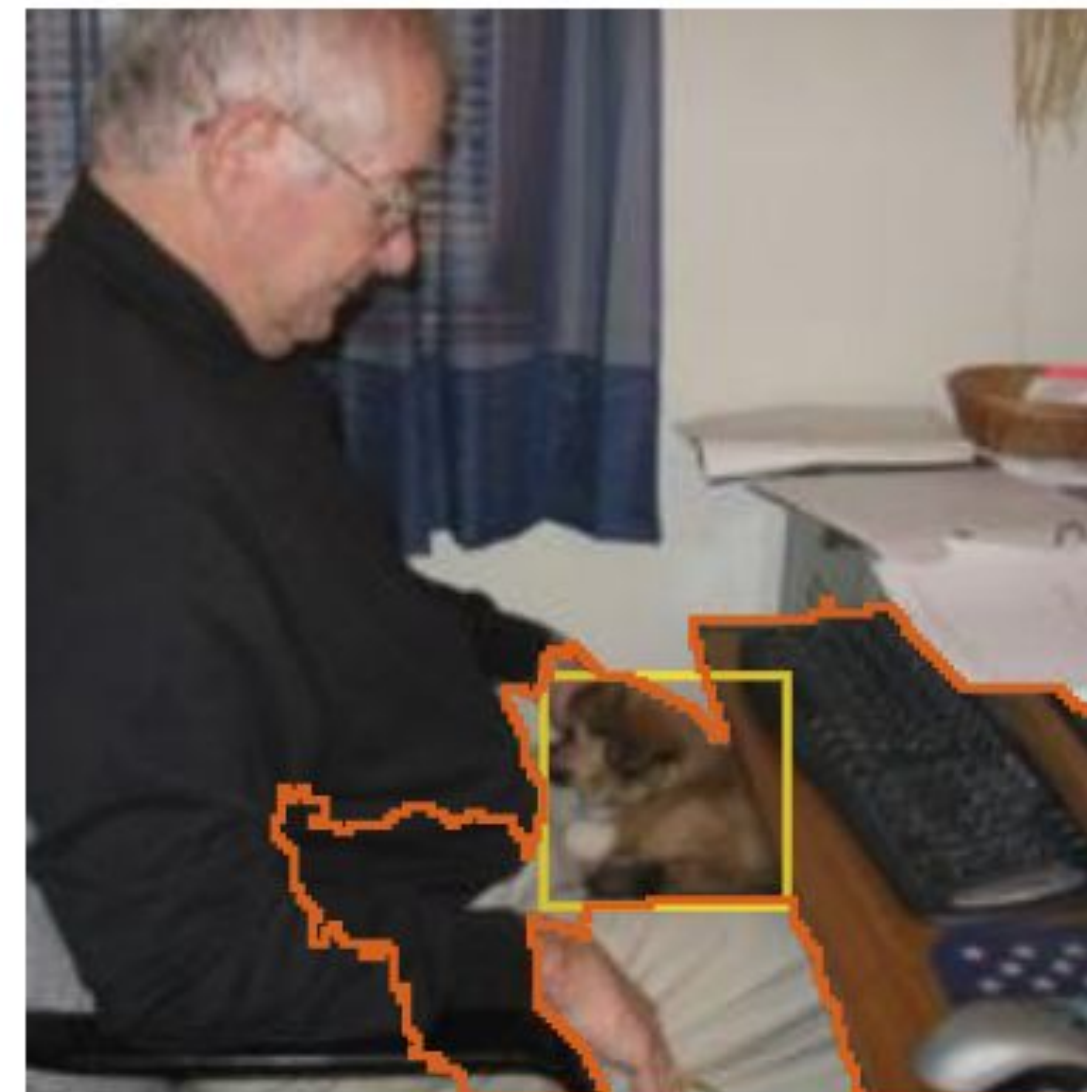
IoU: 0.01 GTC: 0.02 SC: 0.04

mountain bike mountain bike



IoU: 0.21 GTC: 0.95 SC: 0.21

chihuahua laptop

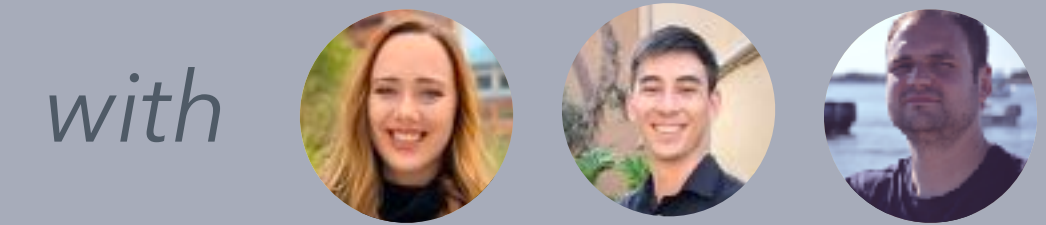


IoU: 0.38 GTC: 1.00 SC: 0.38

cab cab



# Shared Interest: Measuring Human-AI Alignment



## LOW SHARED INTEREST SCORE

Incorrect

Correct

## HIGH SHARED INTEREST SCORE

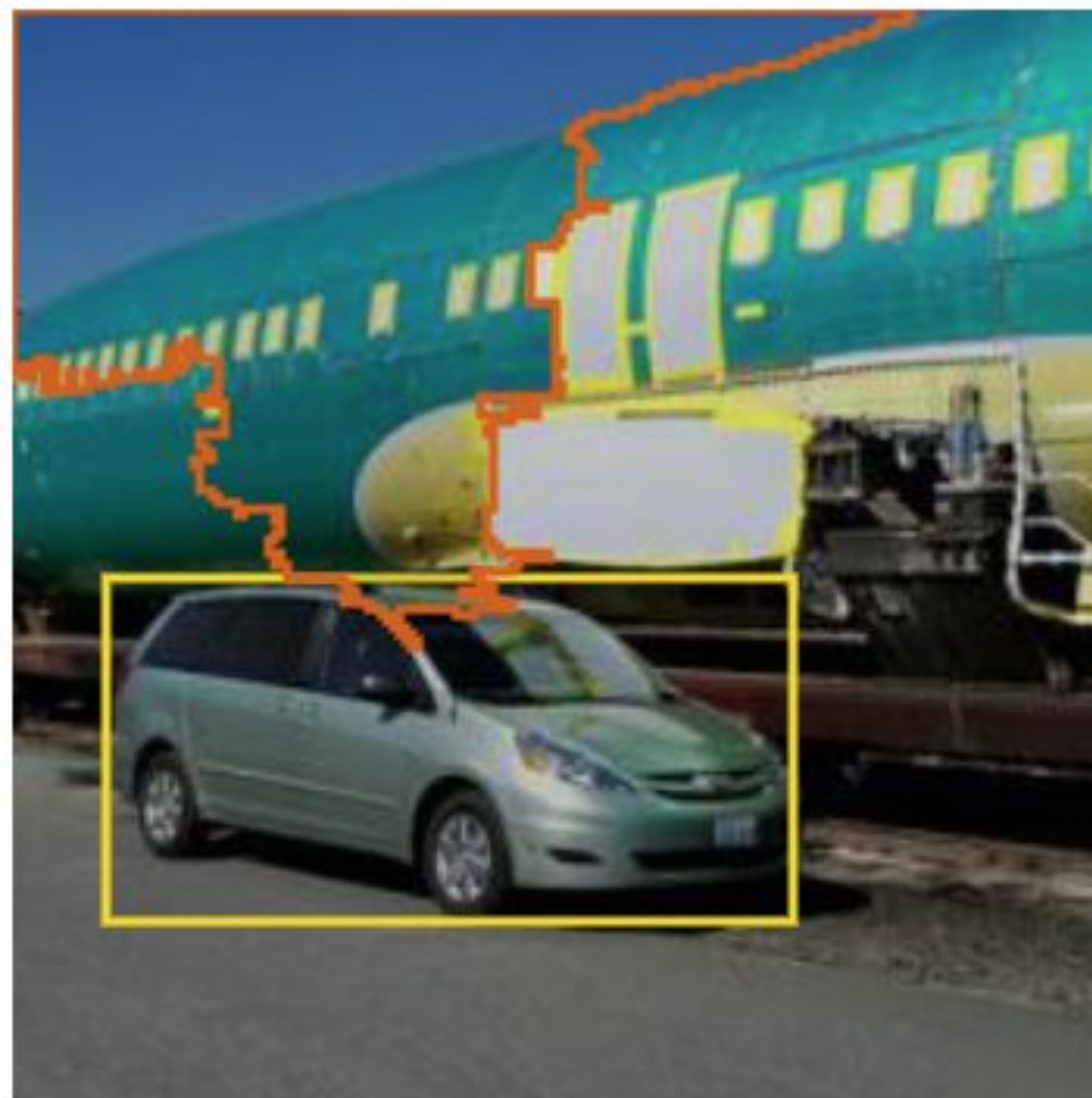
Incorrect

Correct

SALIENCY COVERAGE

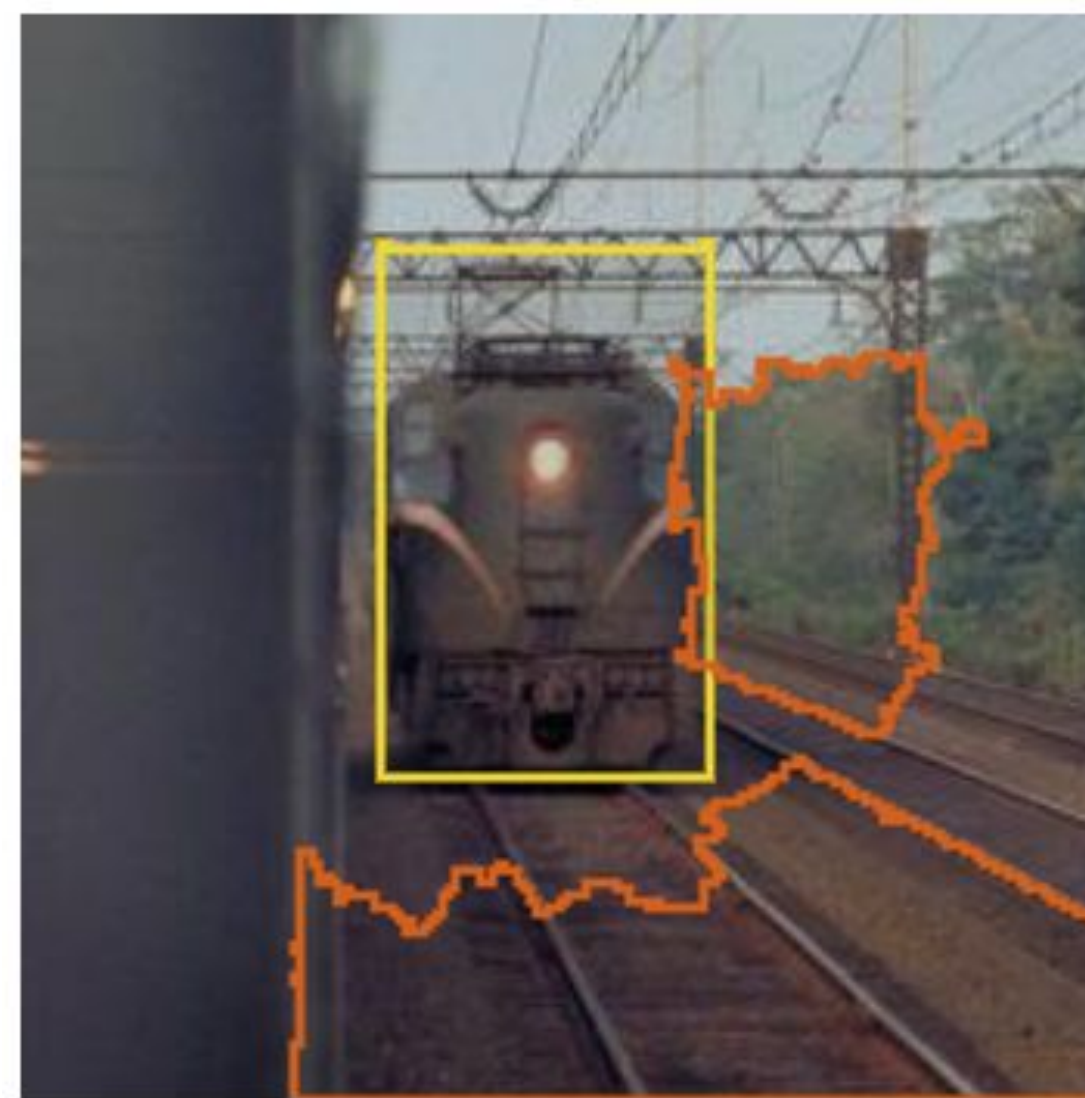
IoU: 0.01 GTC: 0.02 SC: 0.02

minivan airliner



IoU: 0.02 GTC: 0.04 SC: 0.02

electric loco... electric loco...



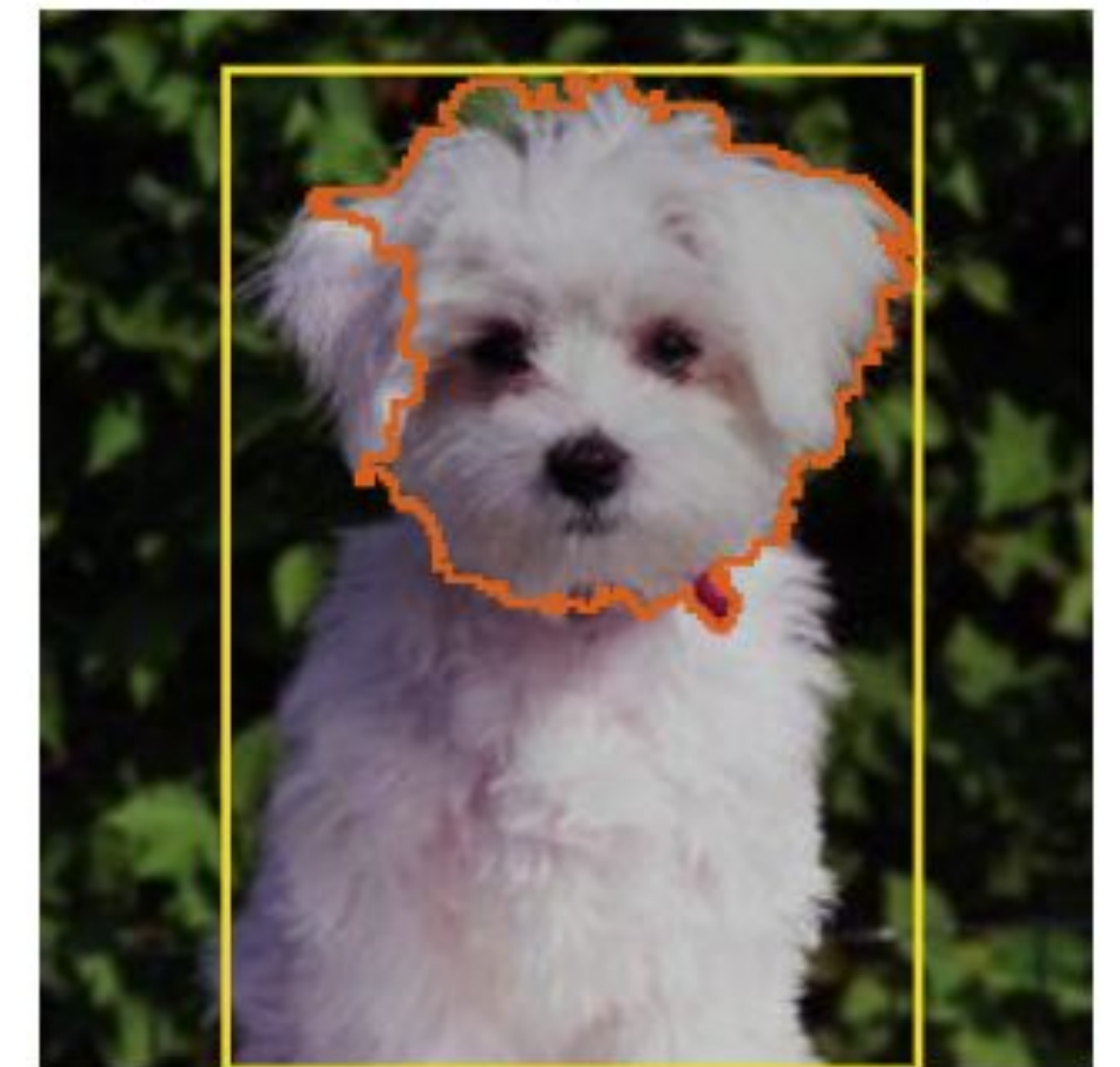
IoU: 0.14 GTC: 0.14 SC: 1.00

english setter dalmatian



IoU: 0.31 GTC: 0.31 SC: 1.00

maltese dog maltese dog



# Shared Interest: Measuring Human-AI Alignment



## LOW SHARED INTEREST SCORE

Incorrect

Correct

## HIGH SHARED INTEREST SCORE

Incorrect

Correct

IOU COVERAGE

IoU: 0.00 GTC: 0.00 SC: 0.00

samoyed arctic fox



IoU: 0.04 GTC: 0.05 SC: 0.17

horse cart horse cart



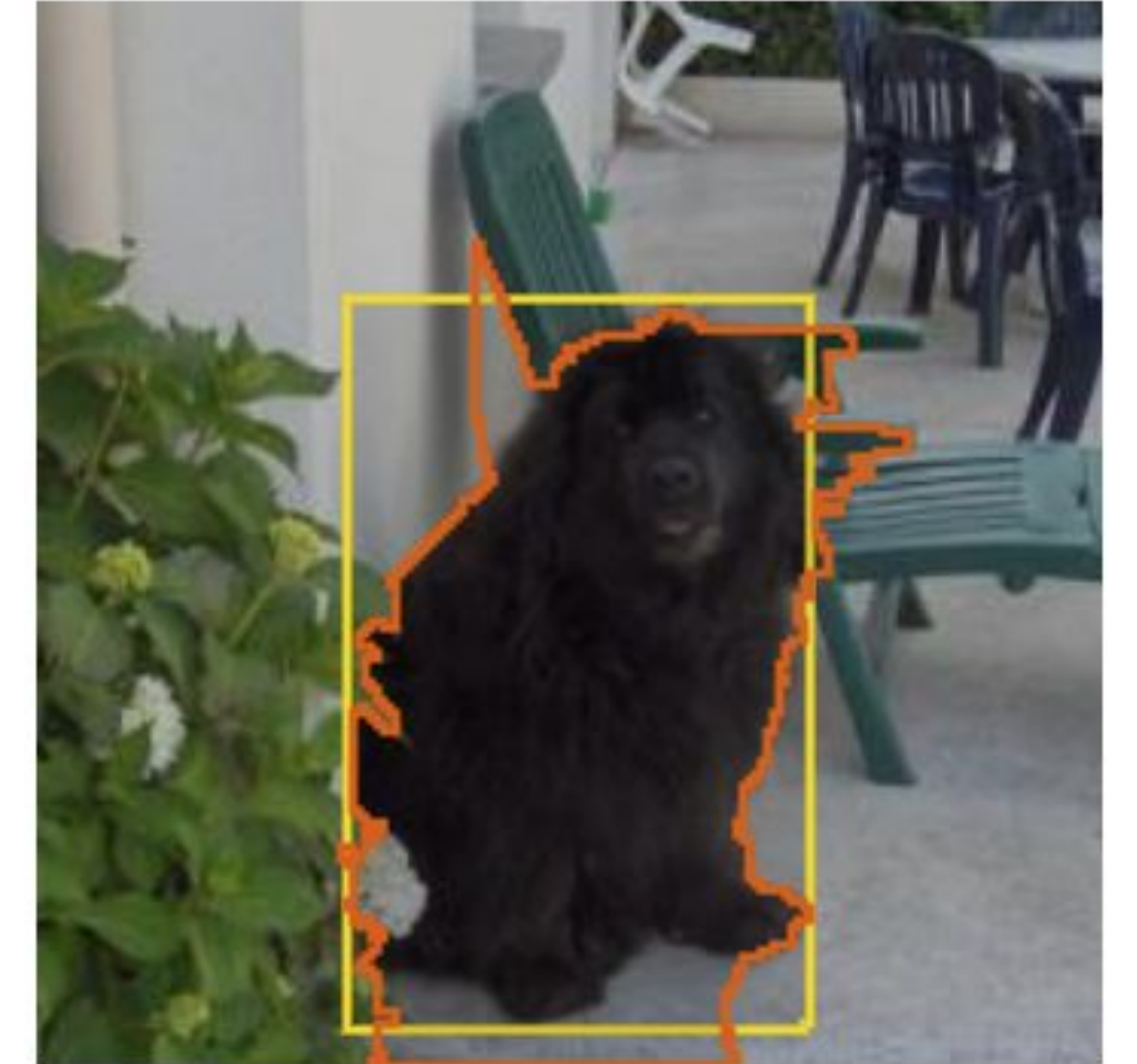
IoU: 0.70 GTC: 0.77 SC: 0.88

snowplow pickup



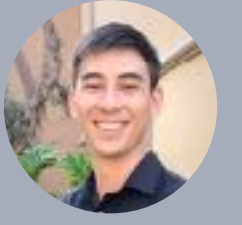
IoU: 0.73 GTC: 0.78 SC: 0.92

newfoundland newfoundland



# Shared Interest: Measuring Human-AI Alignment

with



Select a Score IoU

Select an Image



Reset

Submit

Paint Brush Size



## Gradient and Signal Methods

## Attribution Methods

Input

Gradients

Signal Methods

Gradient x Input

IG Reference Points

DTD Reference Points

[0 1]

[-1 0]

GB

PN

GI

Black

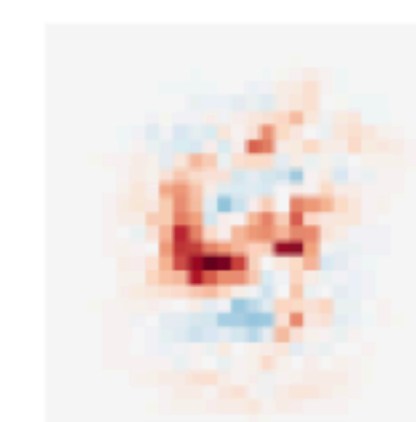
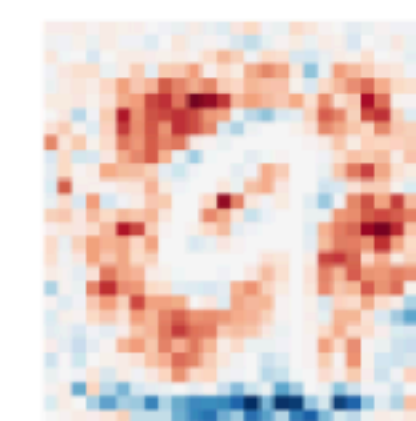
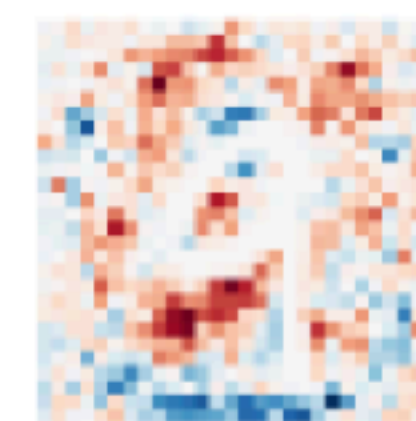
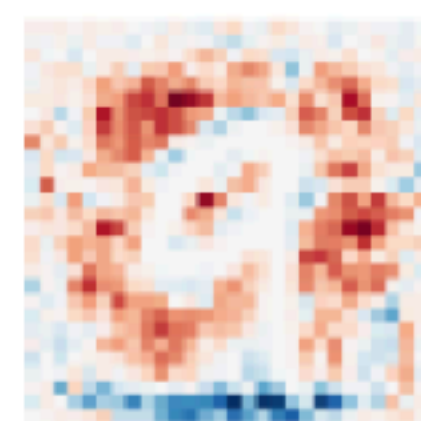
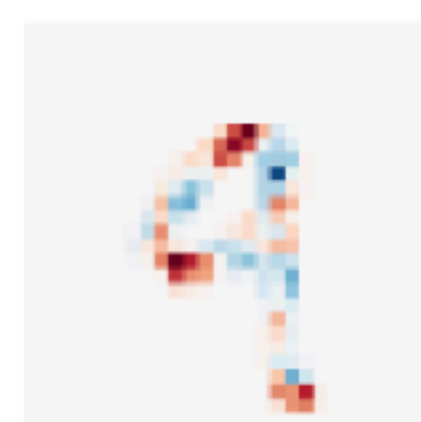
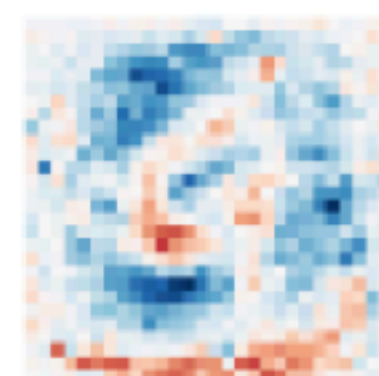
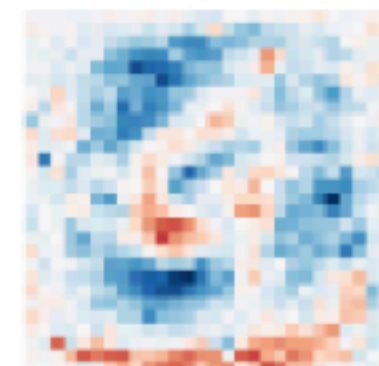
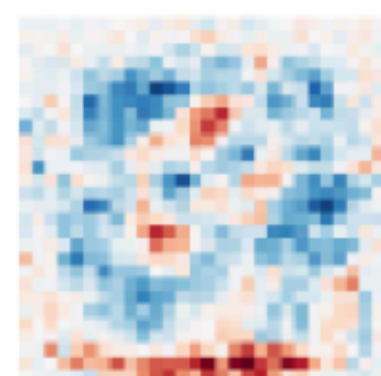
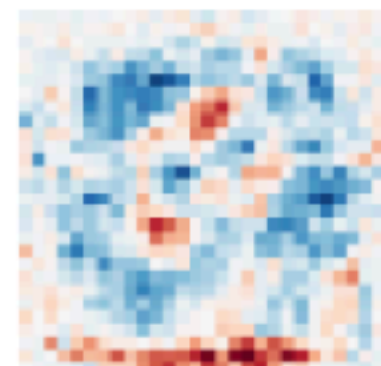
Zero

LRP

PA

Network 1

Network 2

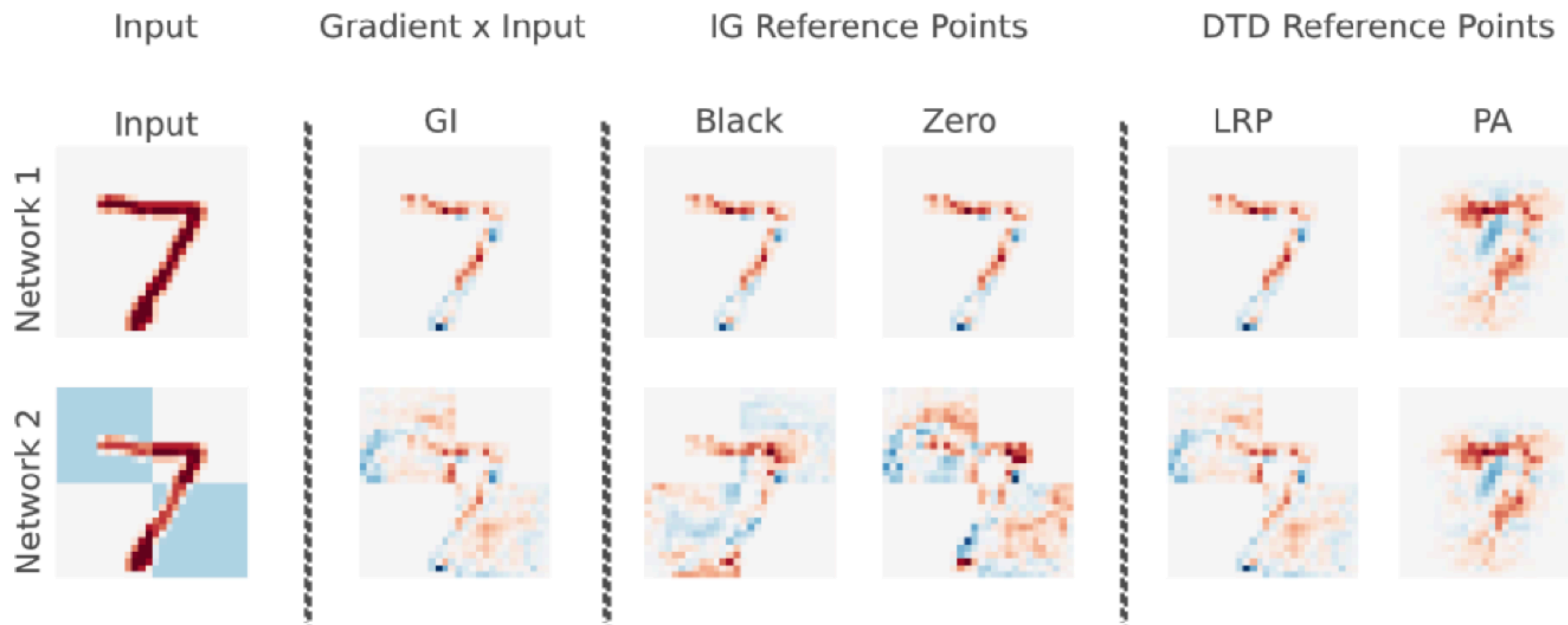


# Add a Constant Vector Shift

MNIST + Constant Shift



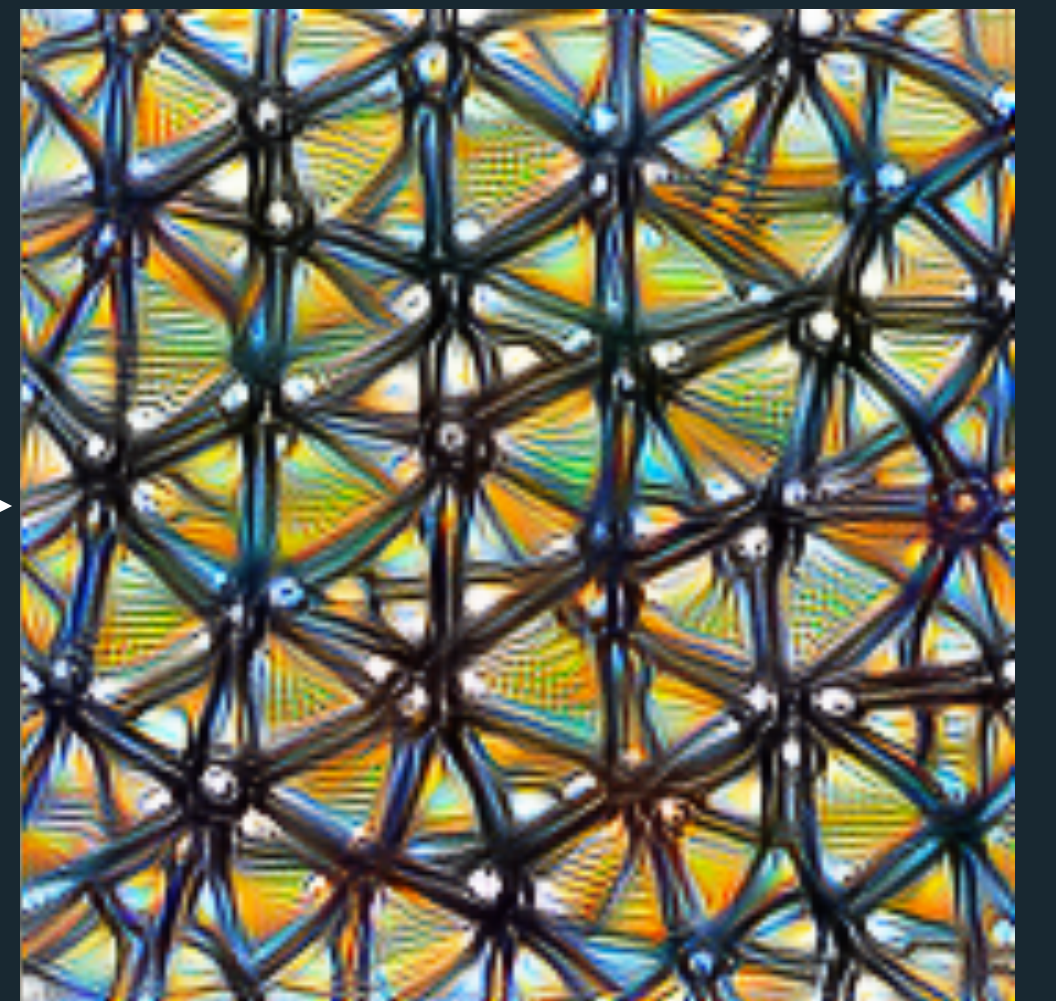
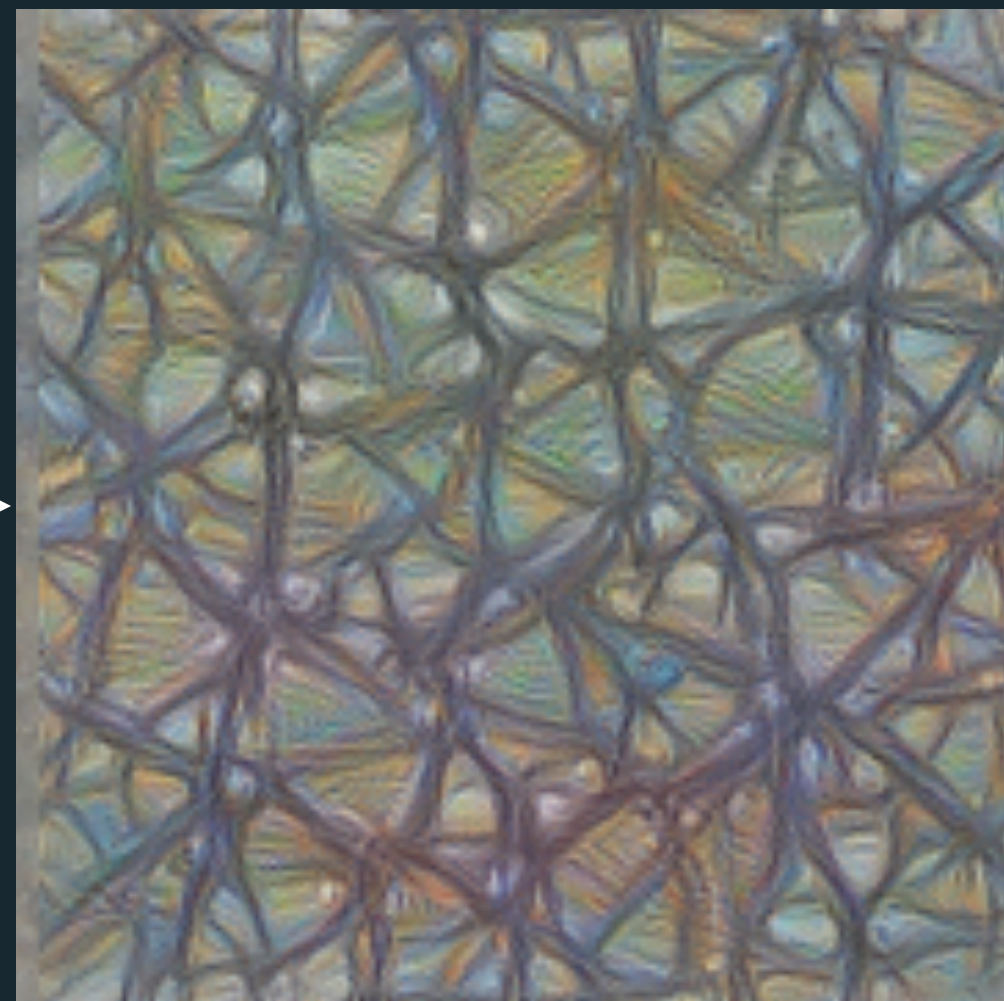
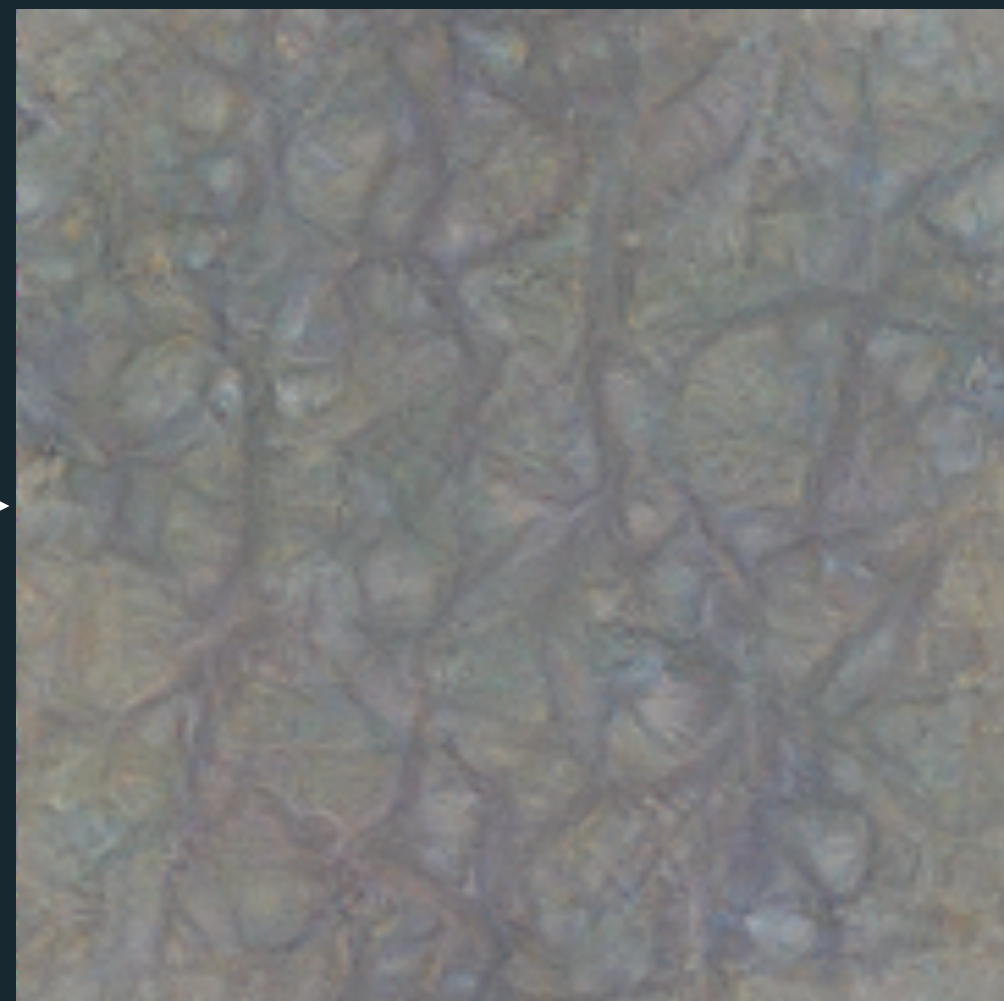
# Attribution Under Constant Vector Shift





# Feature Visualization

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



Step 0

Step 4

Step 48

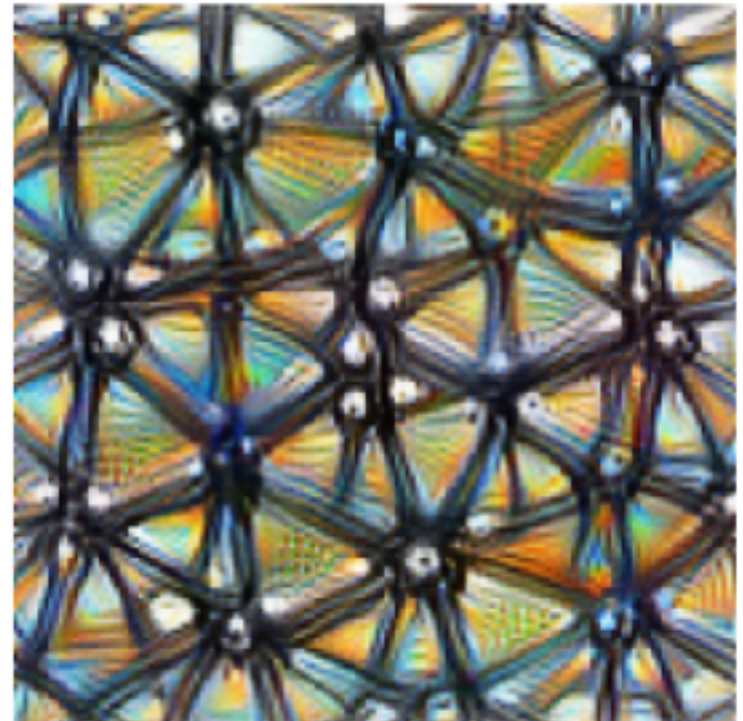
Step 2,048

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



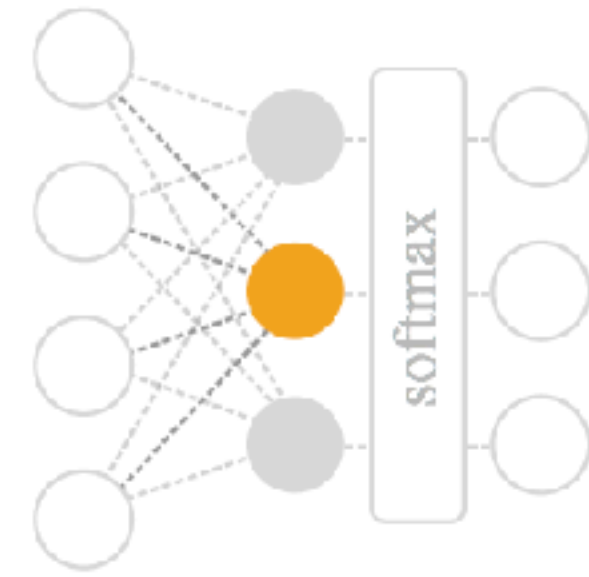
**Neuron**  
`layern[x, y, z]`



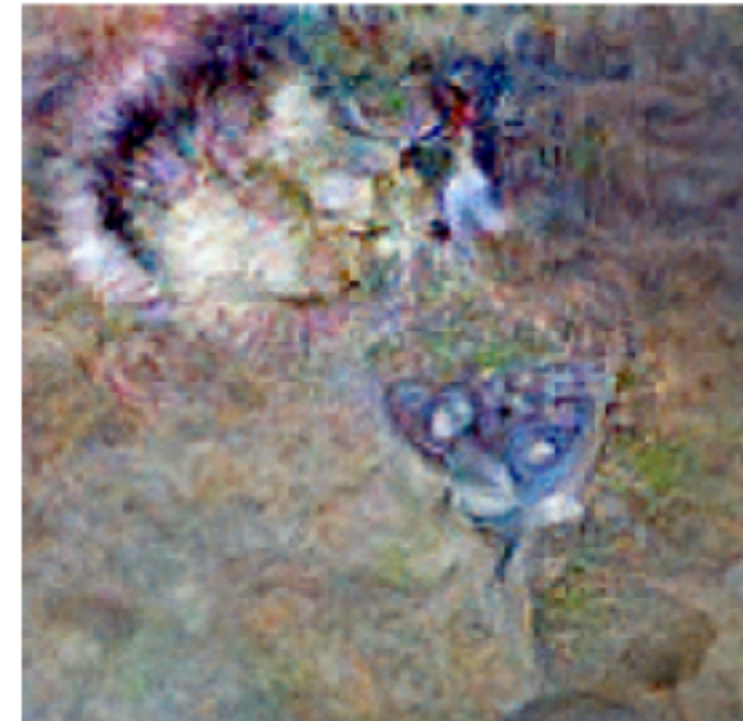
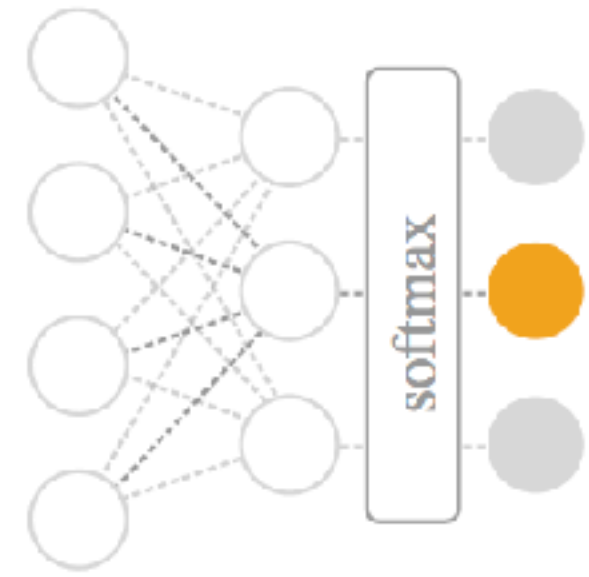
**Channel**  
`layern[:, :, z]`



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



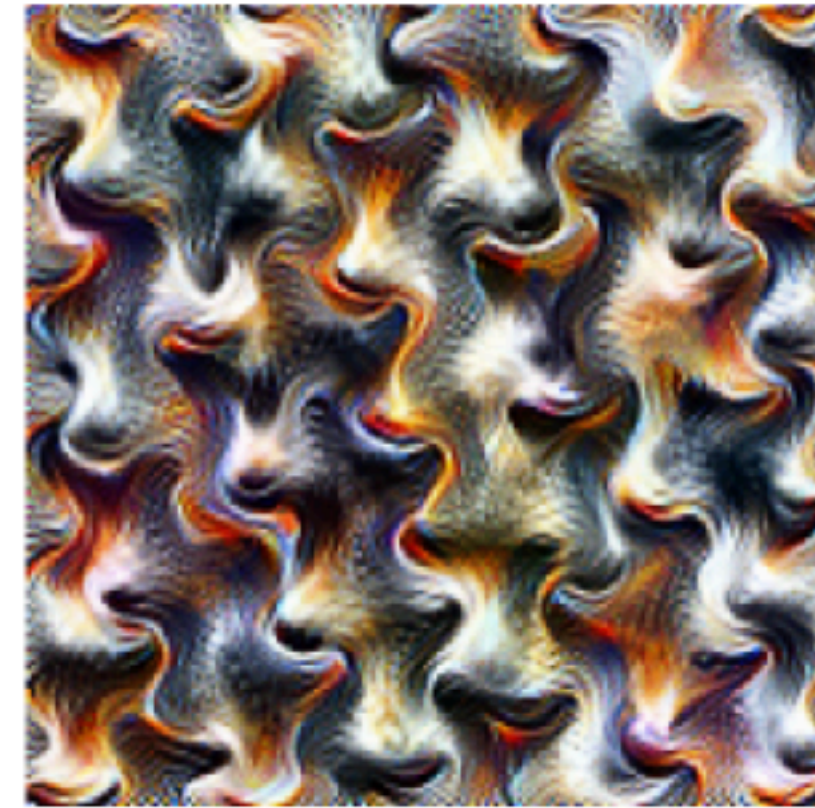
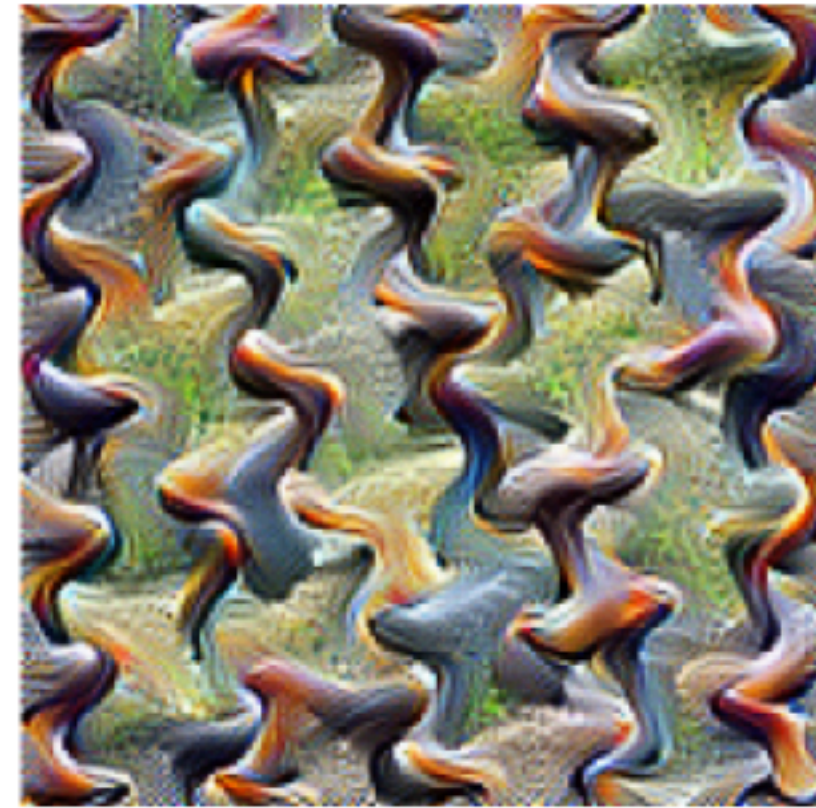
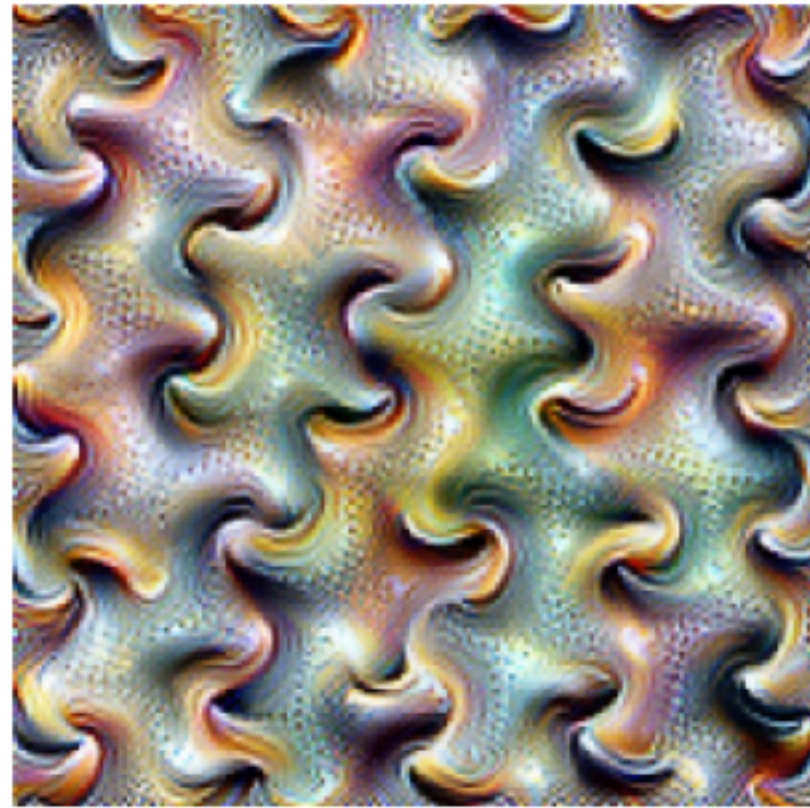
**Class Logits**  
`pre_softmax[k]`



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



Simple Optimization



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



Dataset examples



Simple Optimization



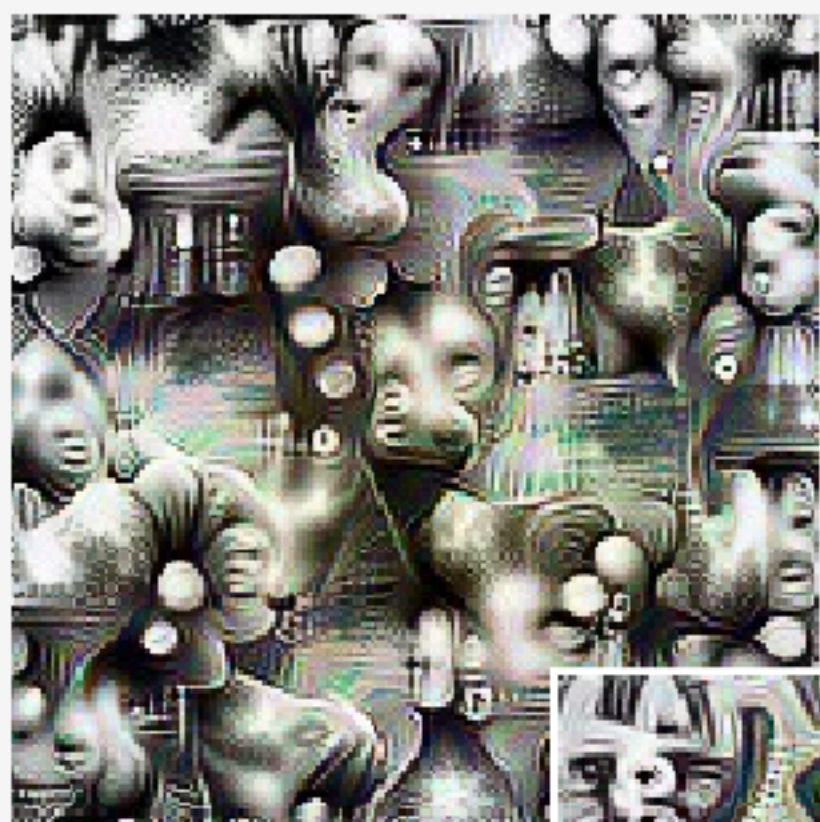
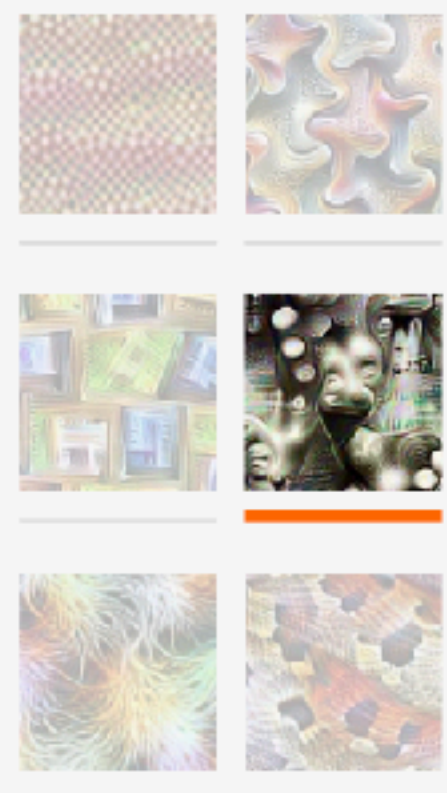
Optimization with diversity reveals multiple types of balls. *Layer mixed5a, Unit 9*



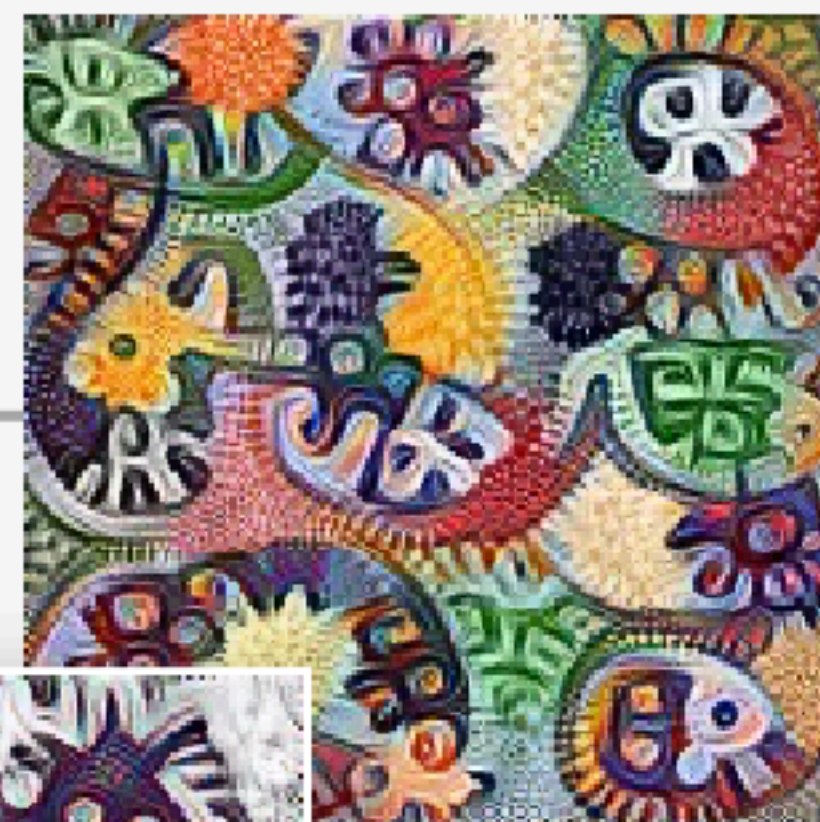
Dataset examples

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

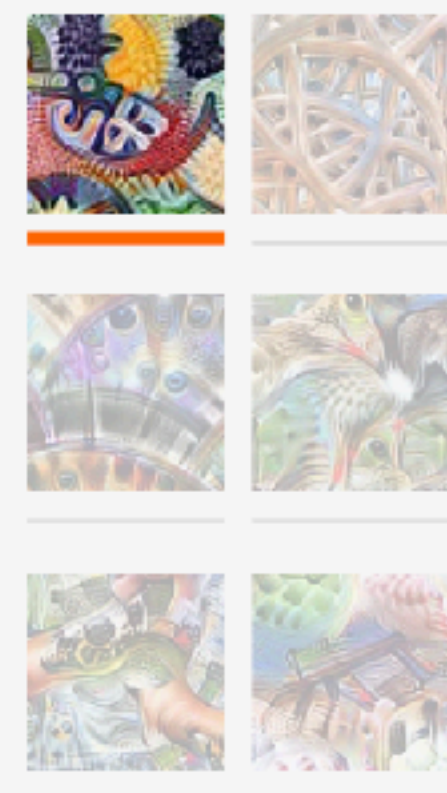
REPRODUCE IN A  
CO NOTEBOOK



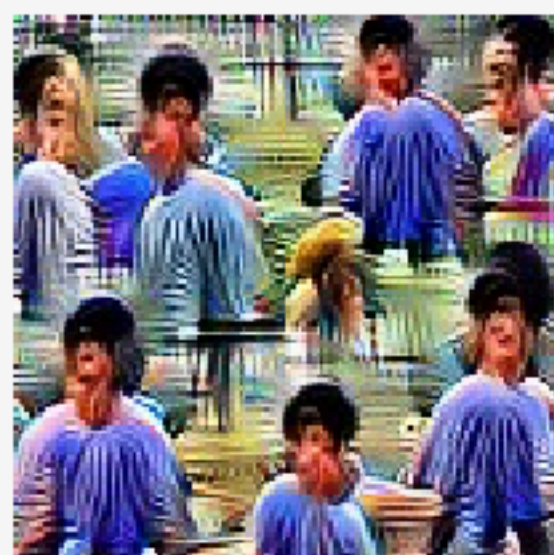
Neuron 1



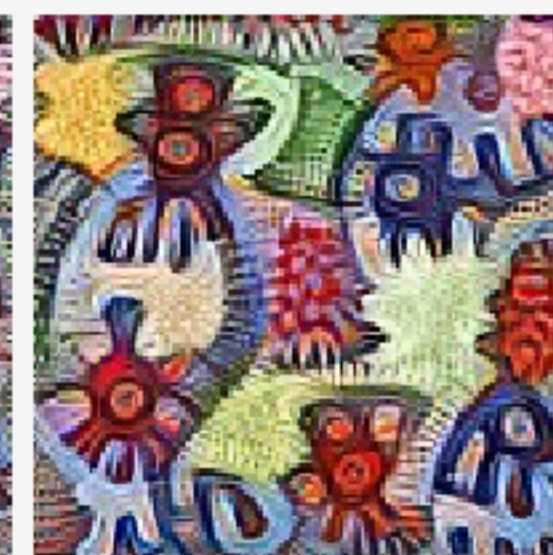
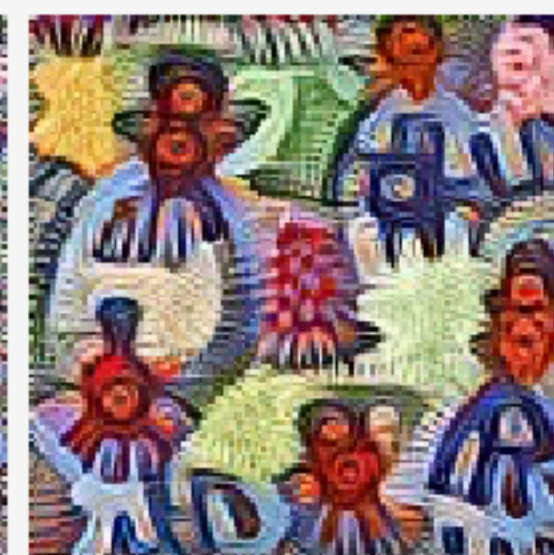
Neuron 2



Jointly optimized

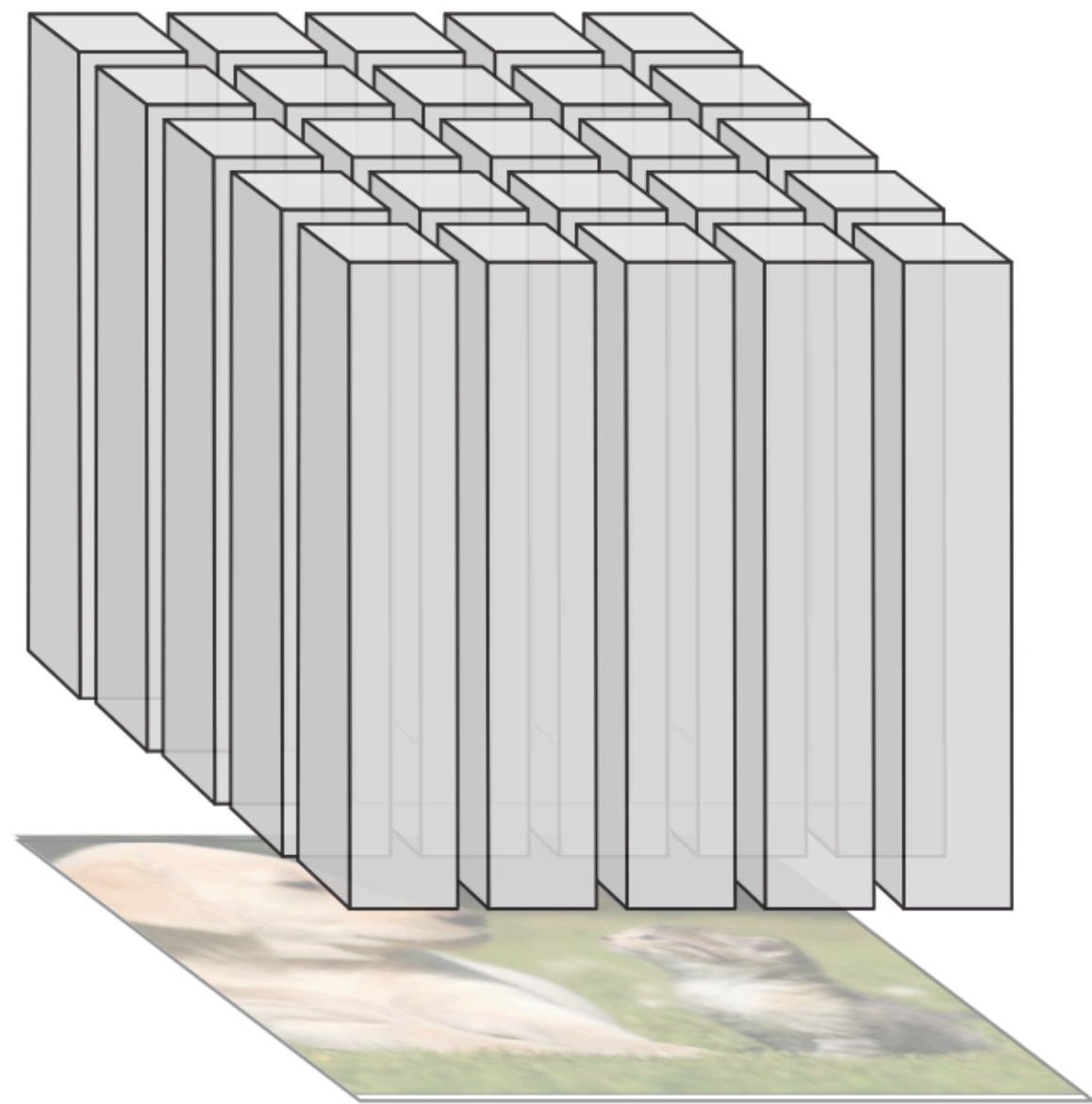


Layer 4b, Unit 475



Layer 4a, Unit 476

## Spatial Activations

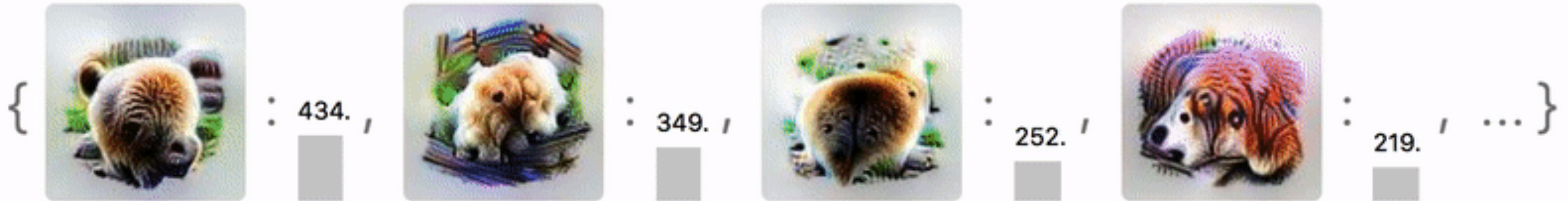


$a_{1,0} = [0, 0, 0, 0, 49.6, 0, 43.6, 30.2, 119.8, 62.7, 0, 51\dots]$

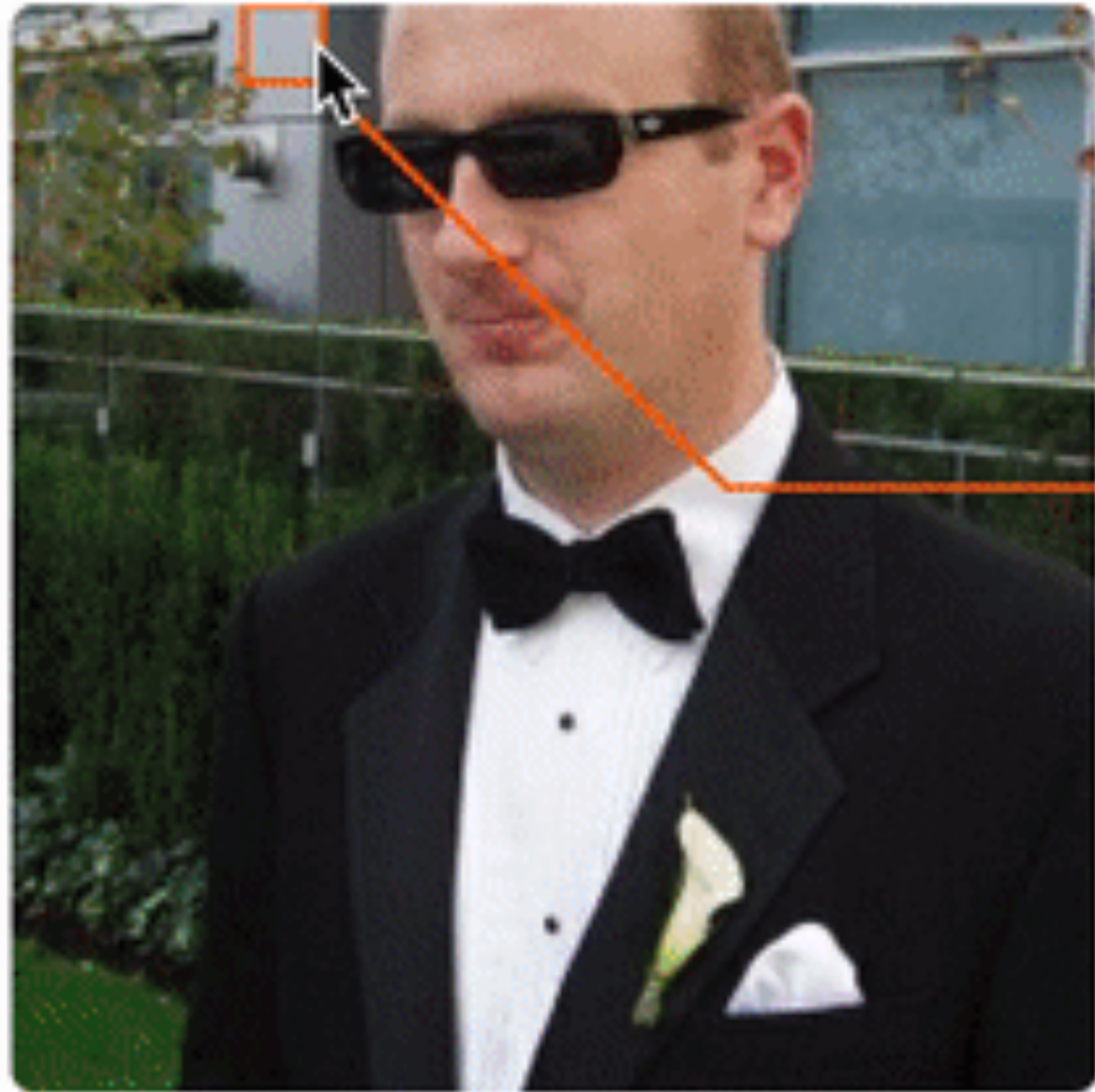
# Semantic Dictionaries



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



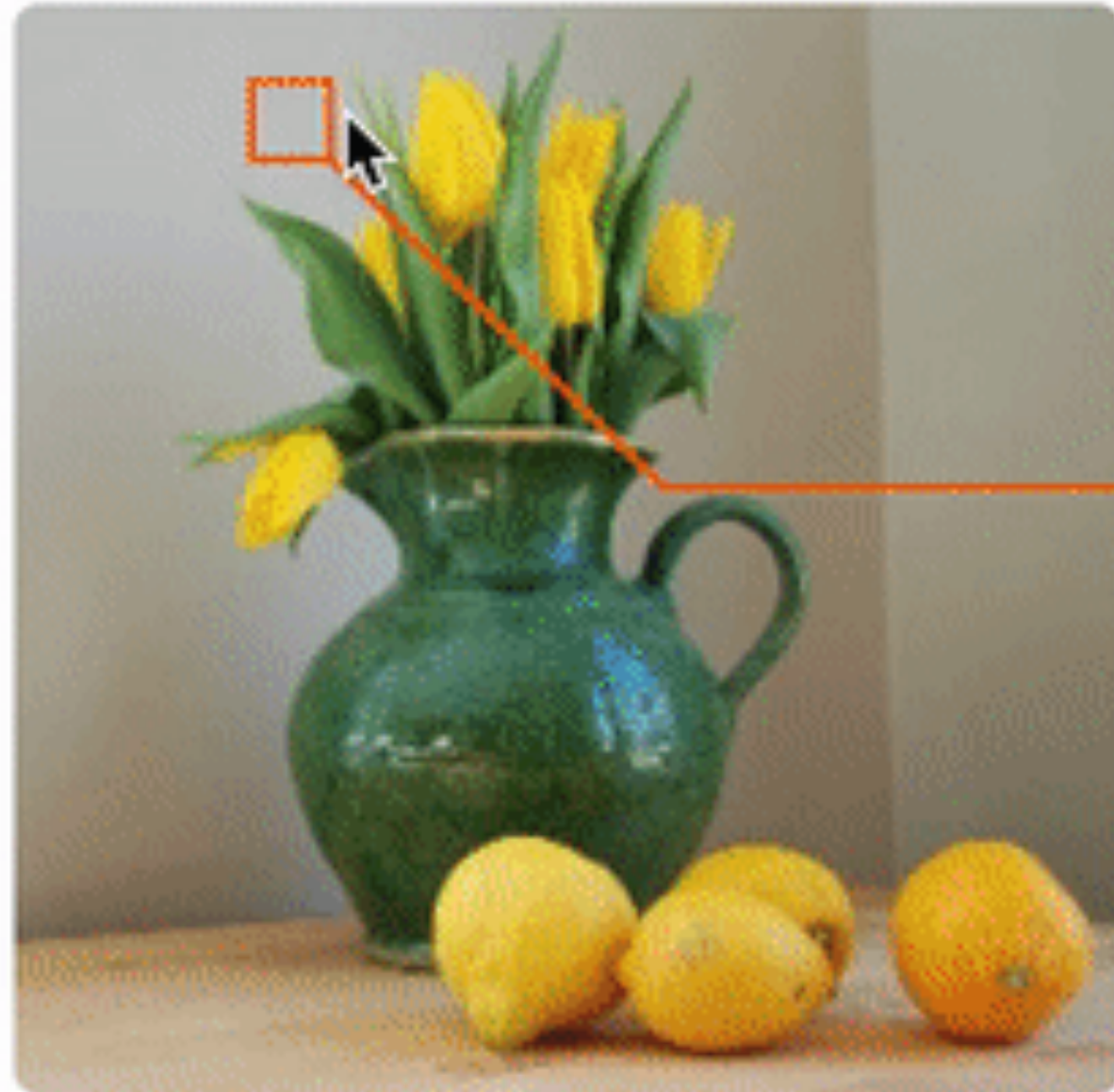
# Semantic Dictionaries



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



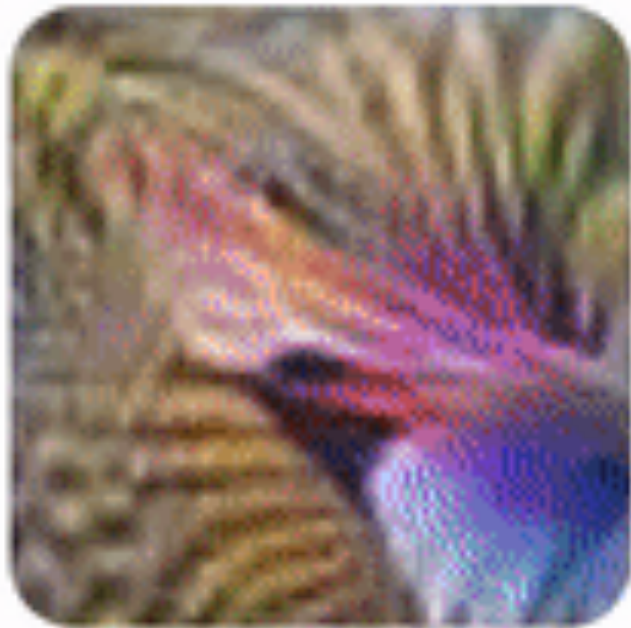
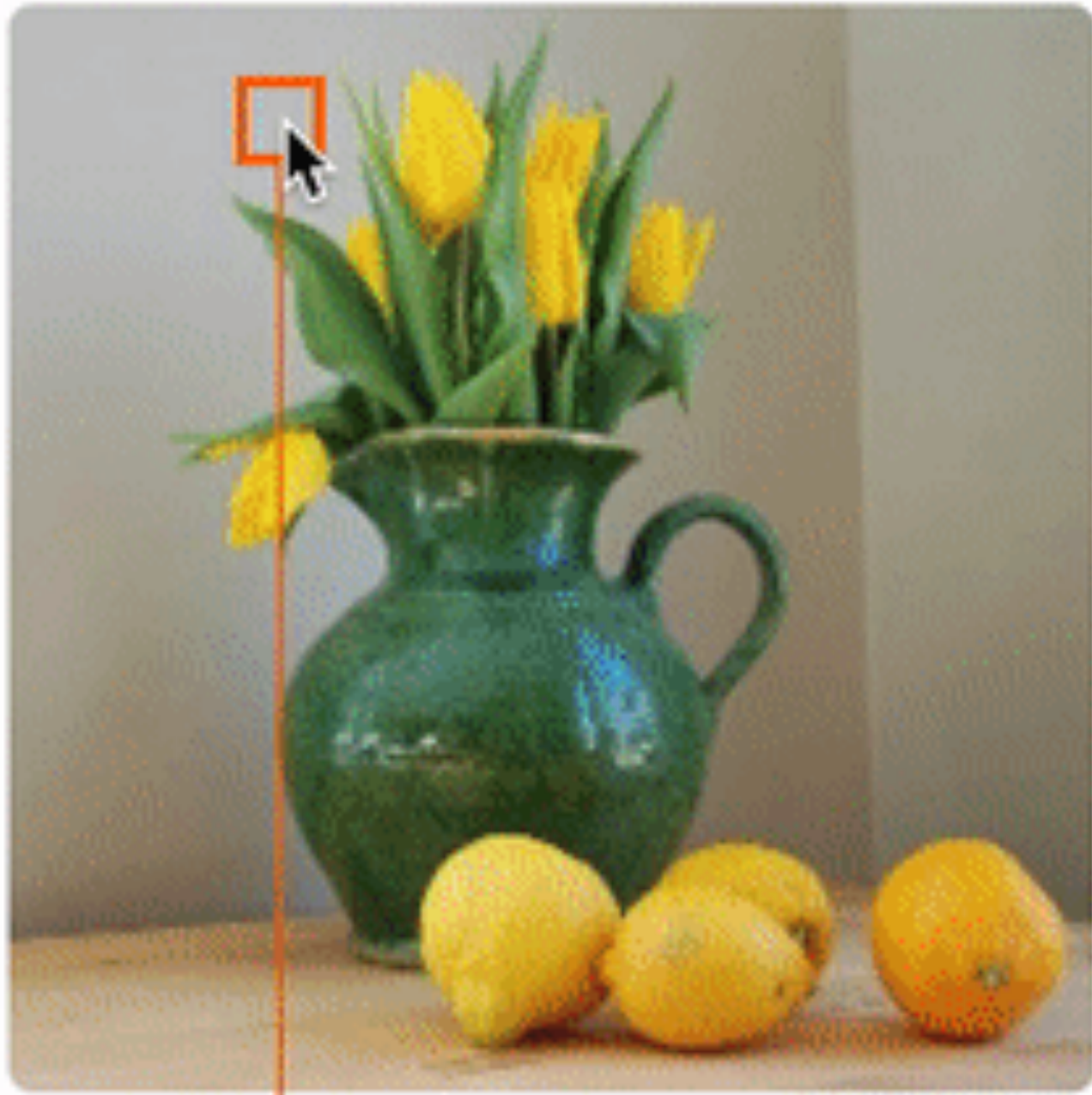
# Semantic Dictionaries



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

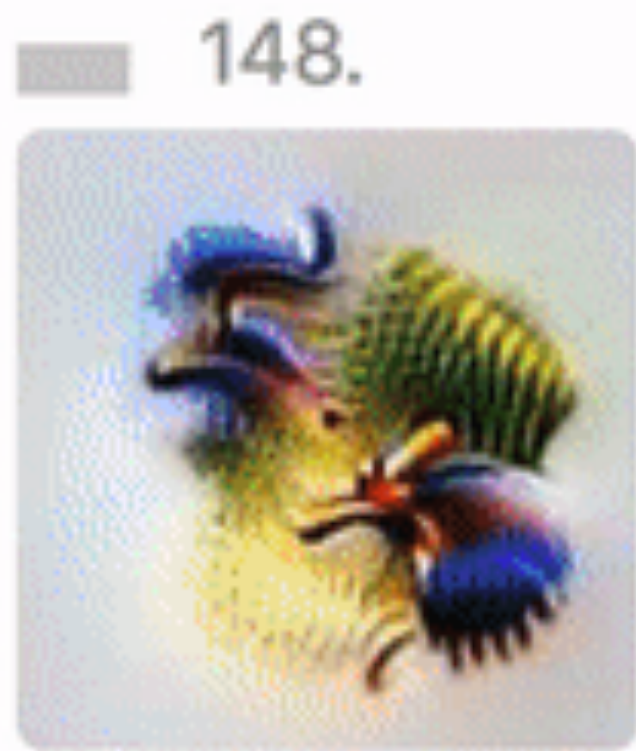




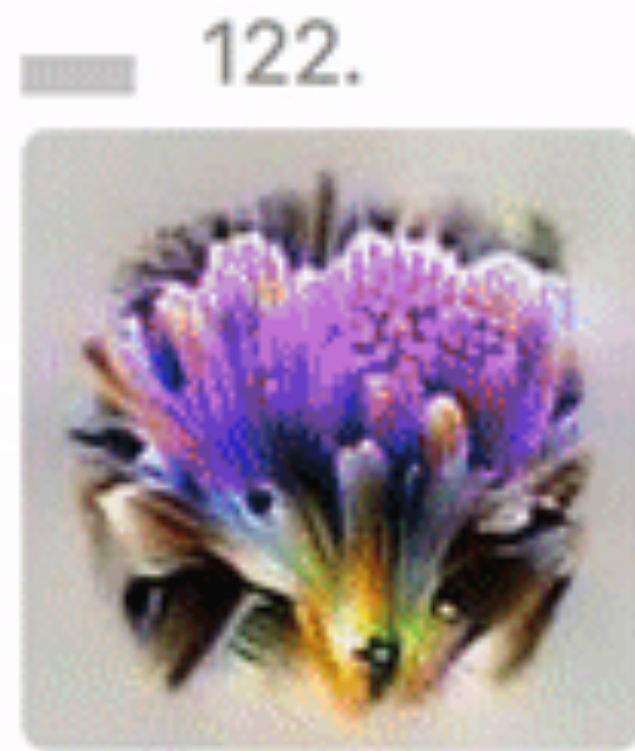


Activation Vector

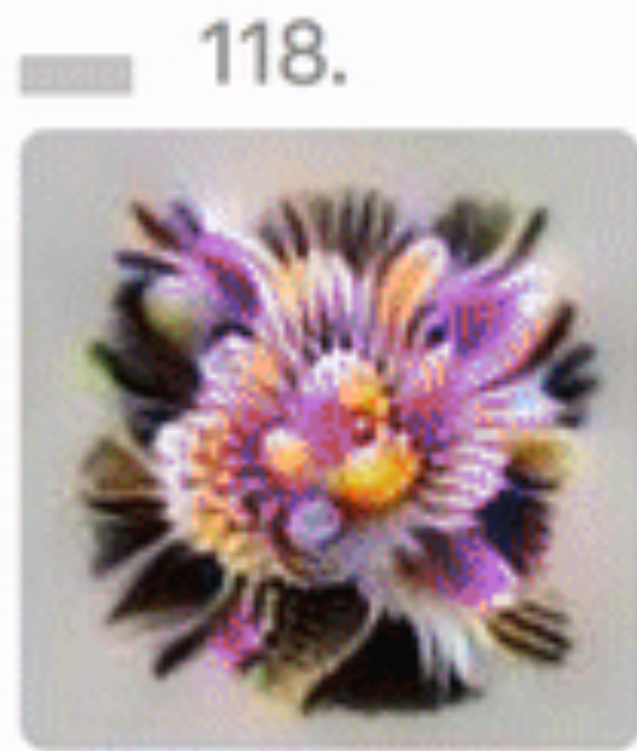
=



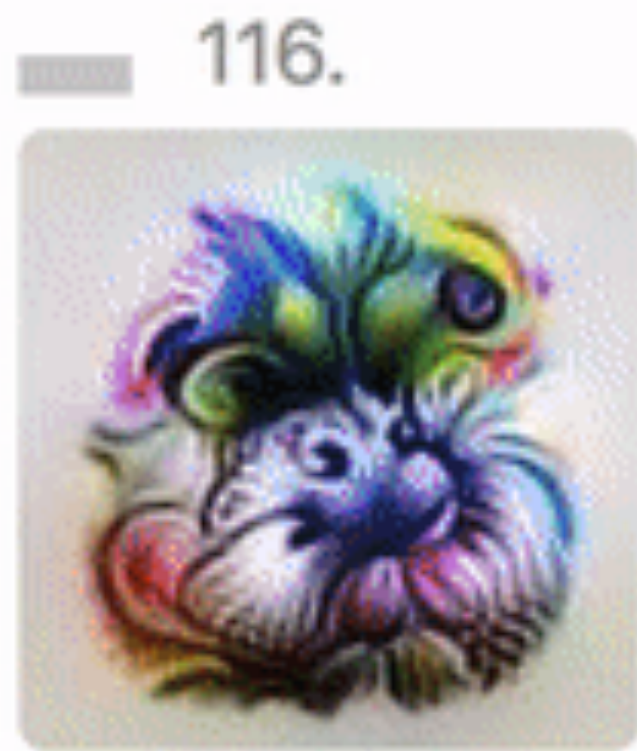
+



+



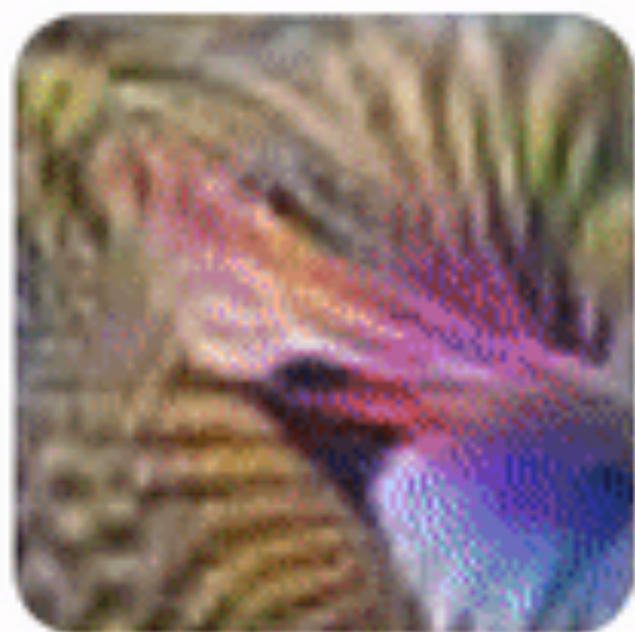
+



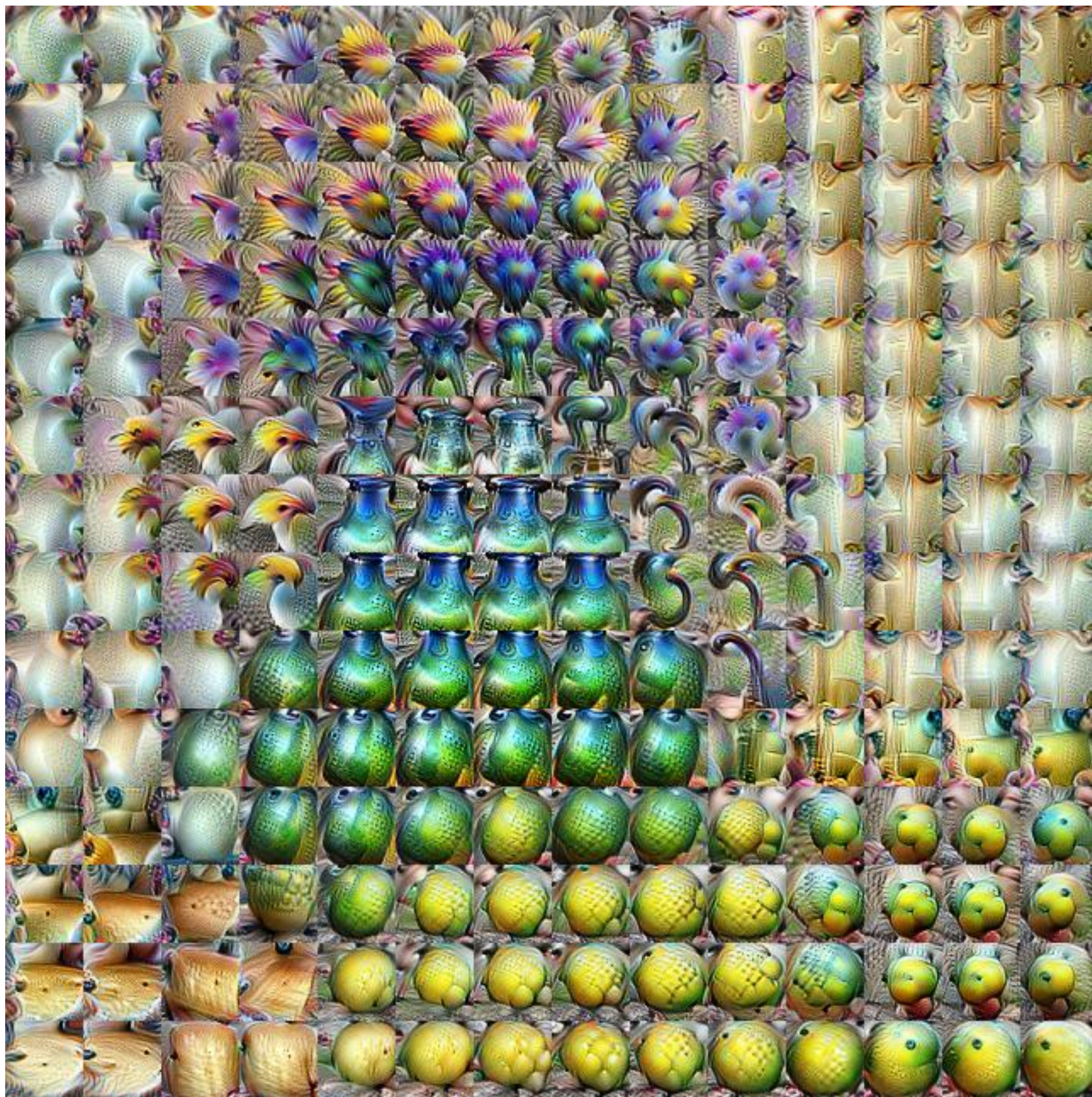
+

...

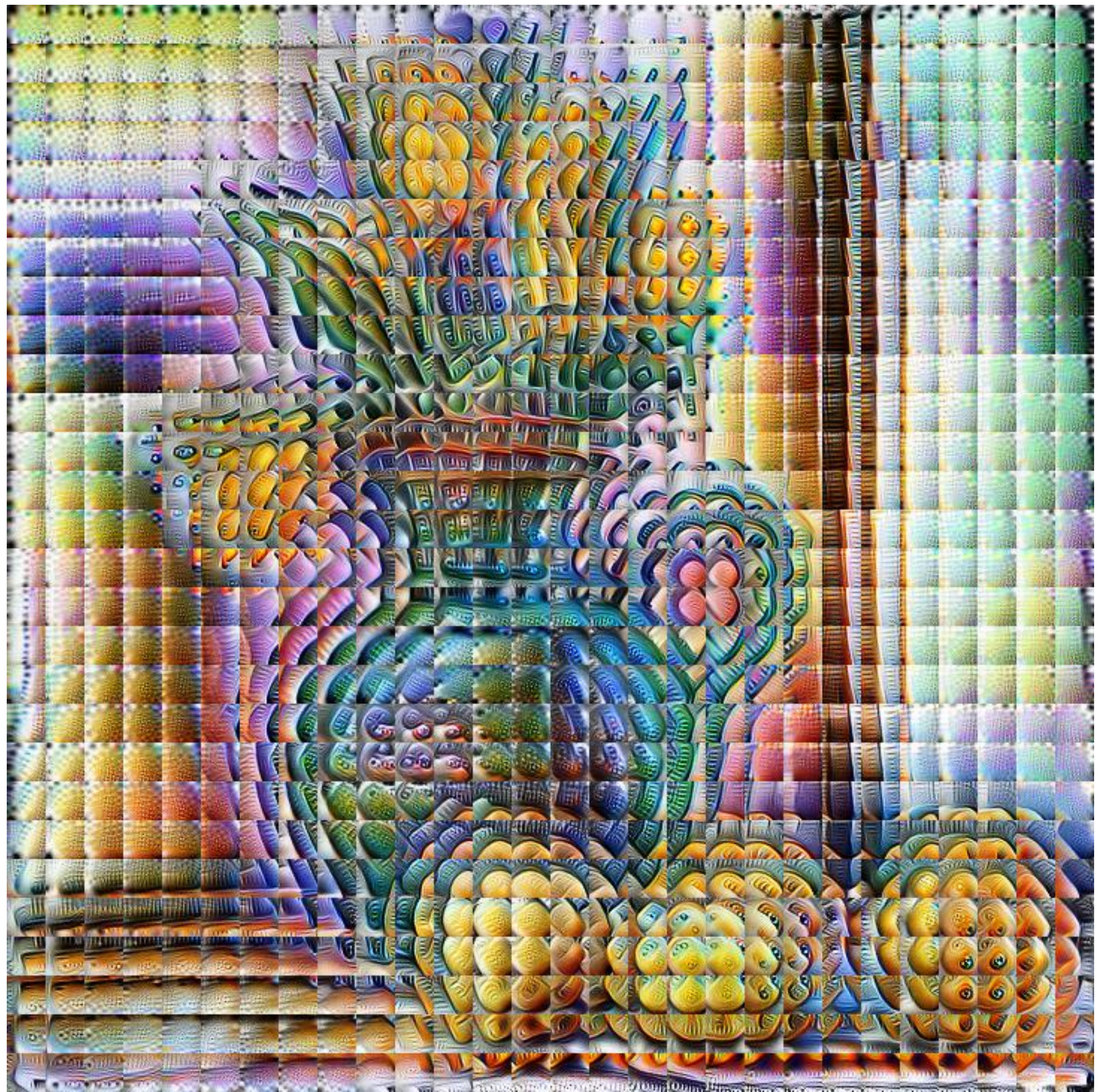
Channels



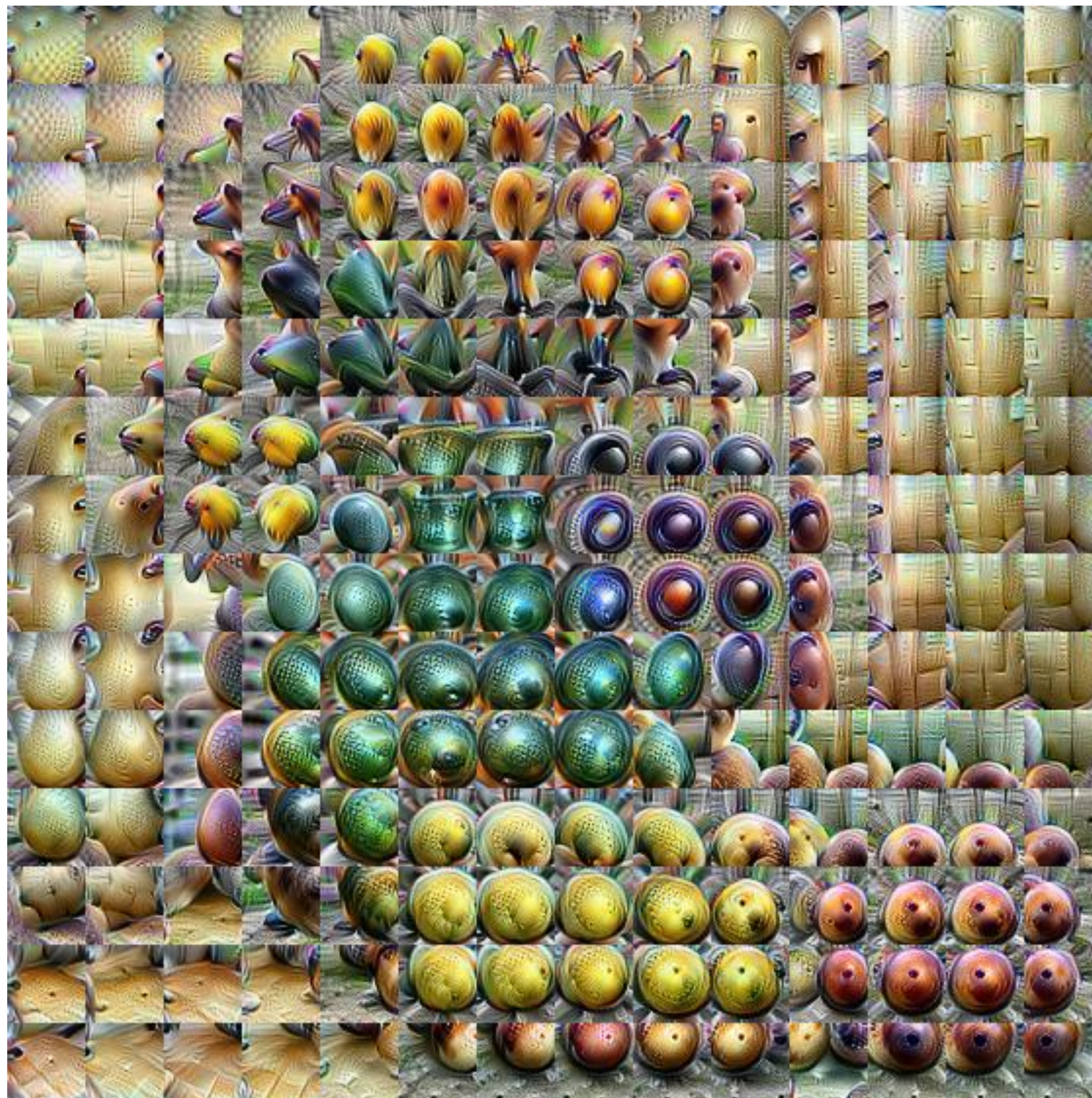
Activation Vector



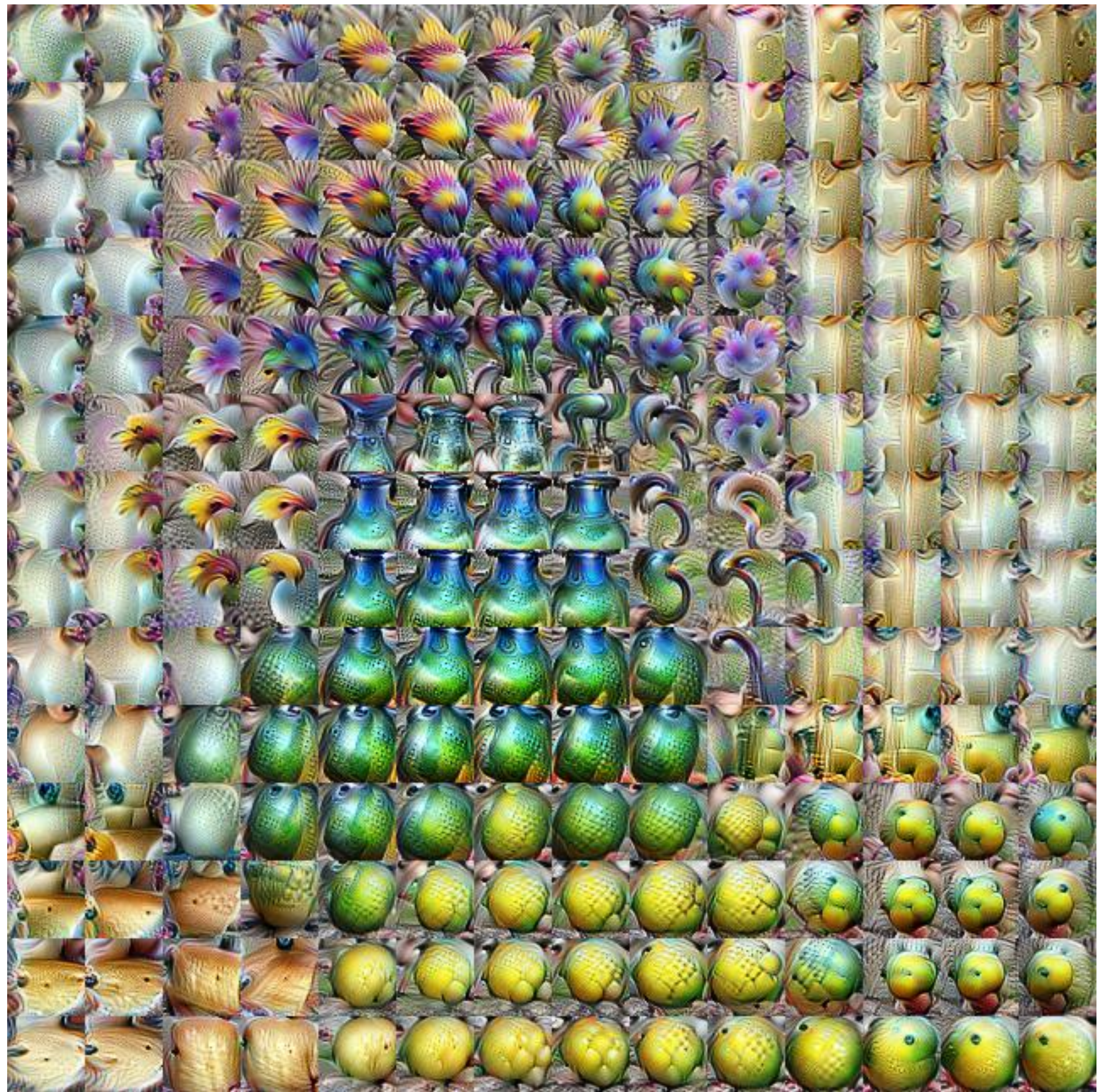
mixed3a: *Edges*

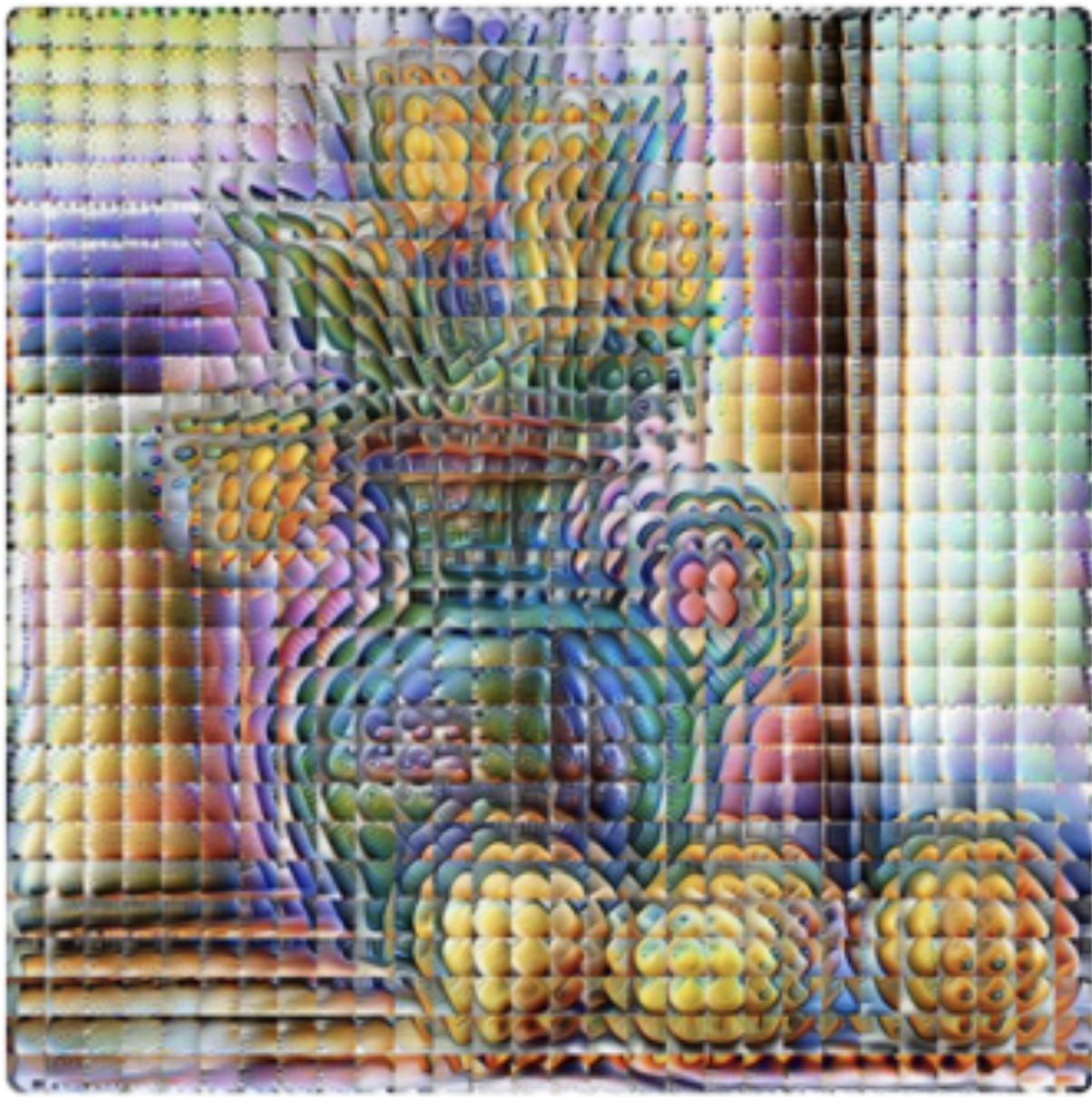


mixed4a: *Geometries*

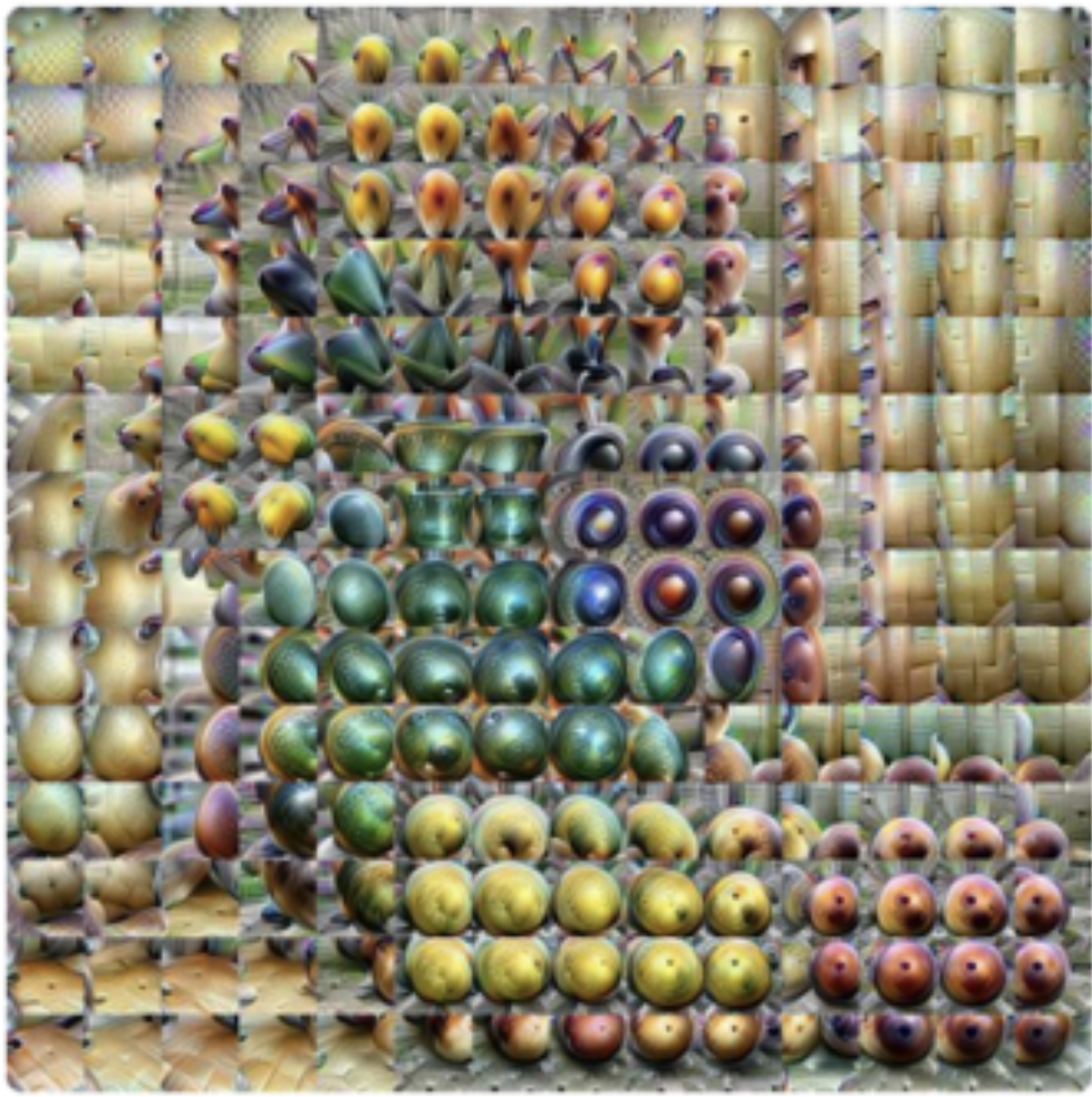


mixed4d: *Objects*

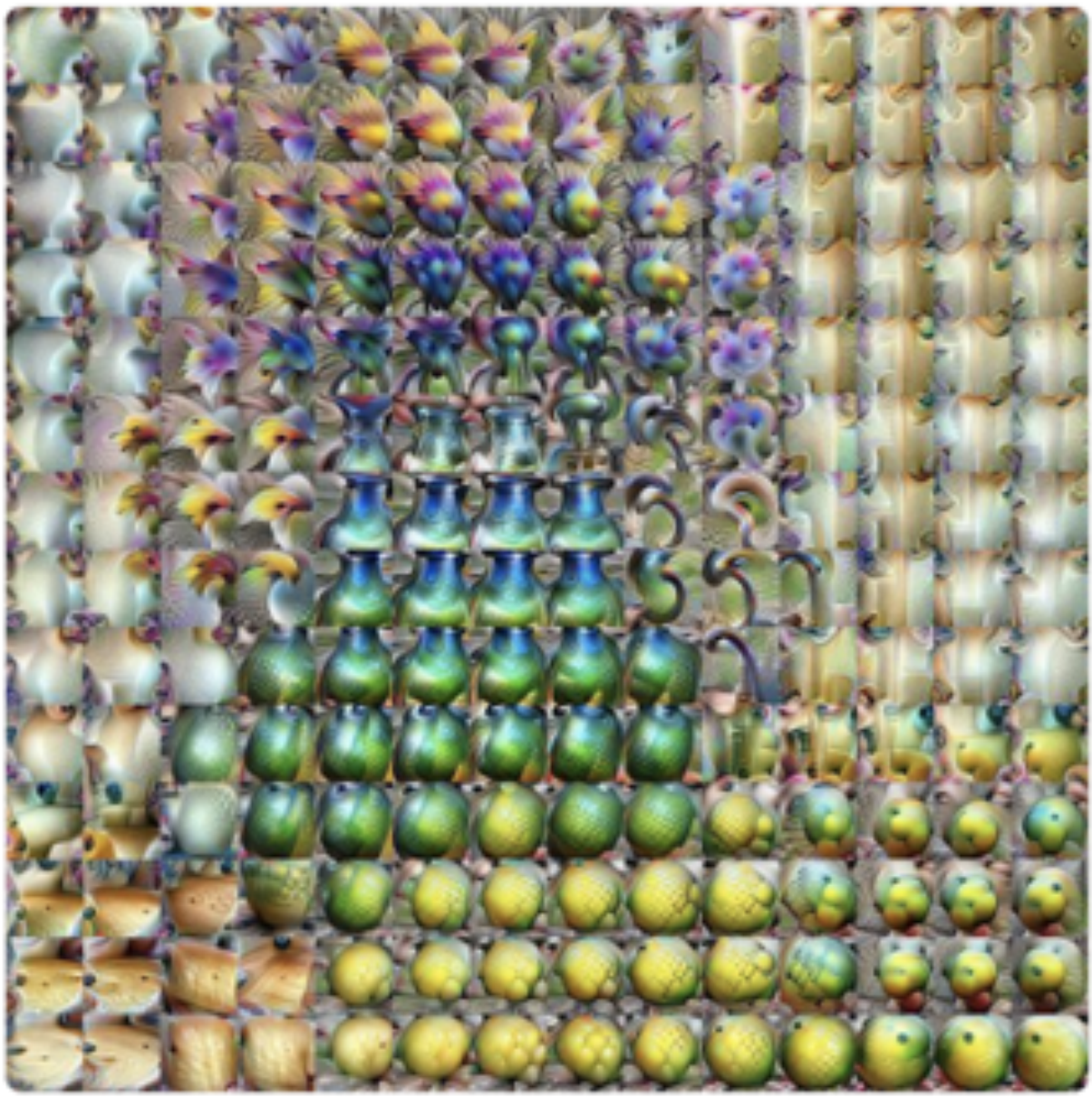




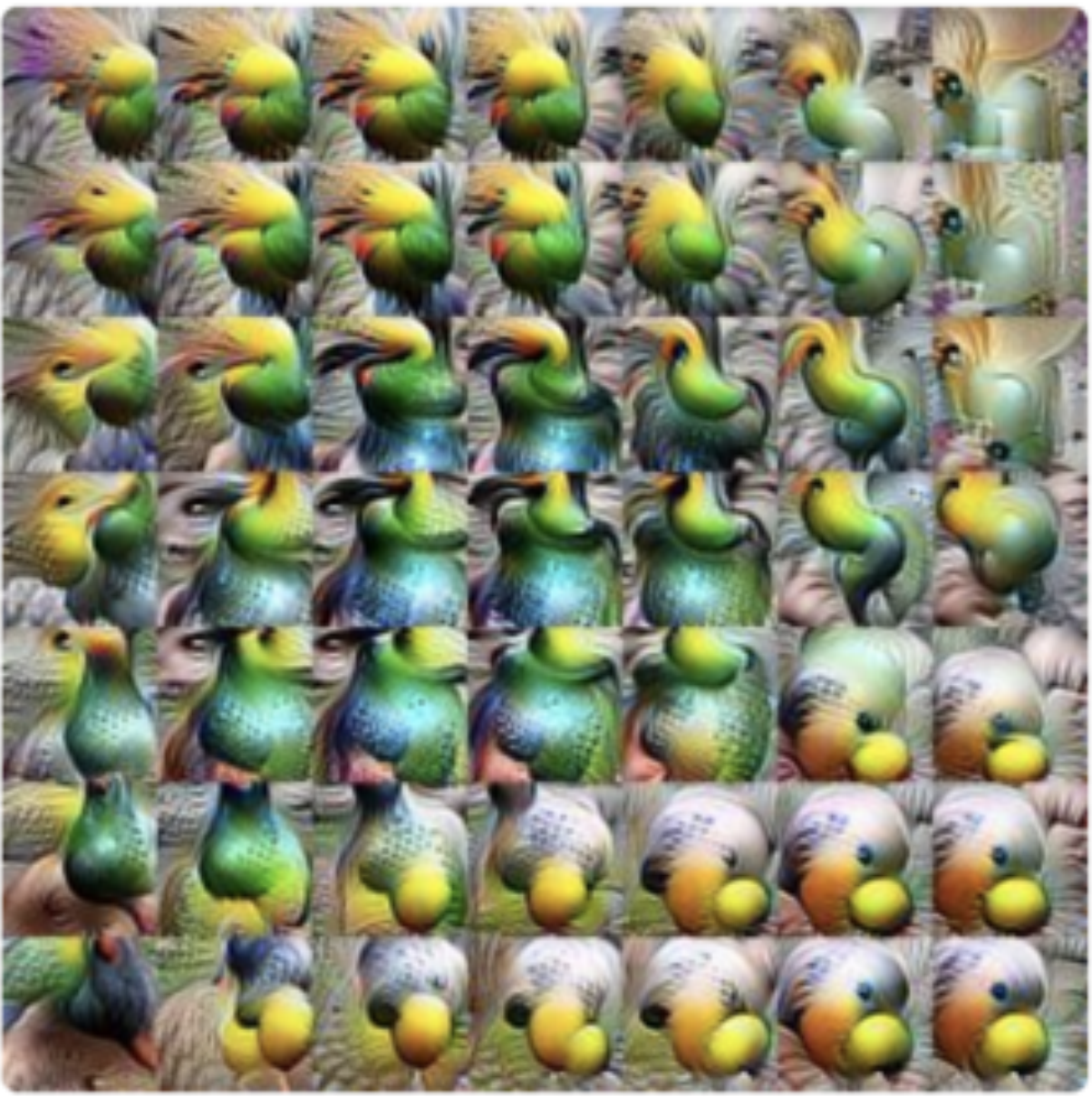
MIXED3A



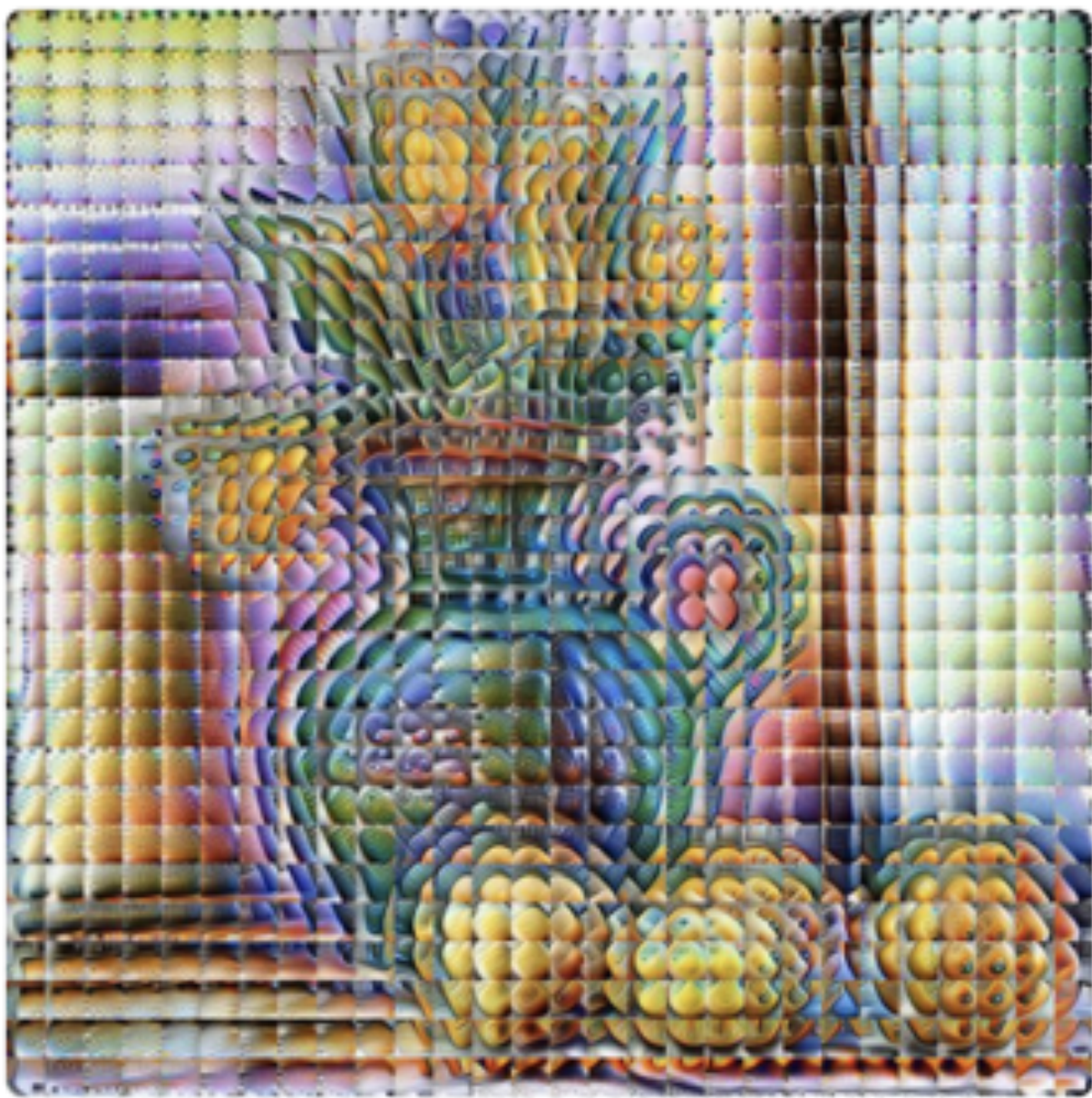
MIXED4A



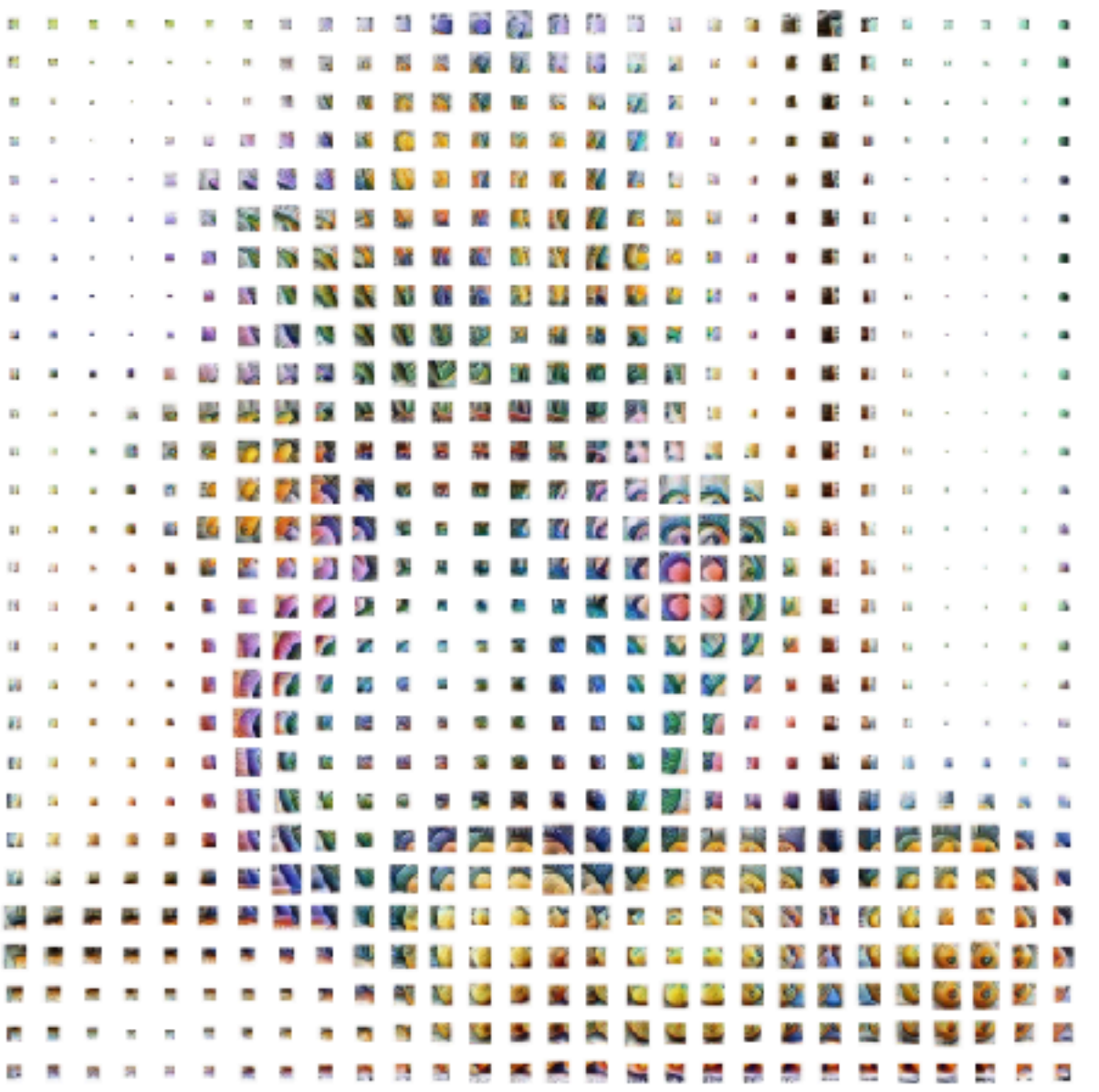
MIXED4D



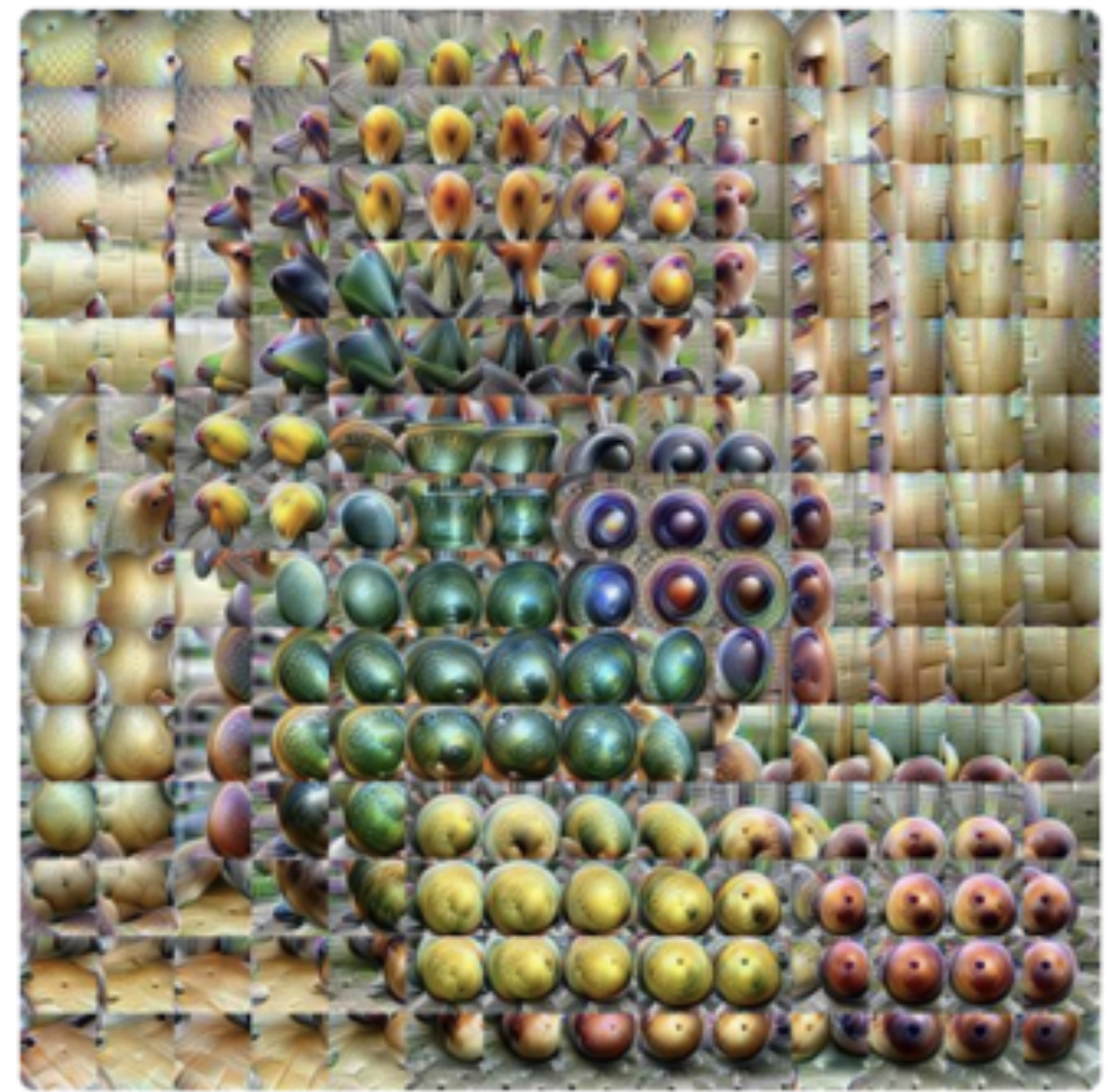
MIXED5A



MIXED3A



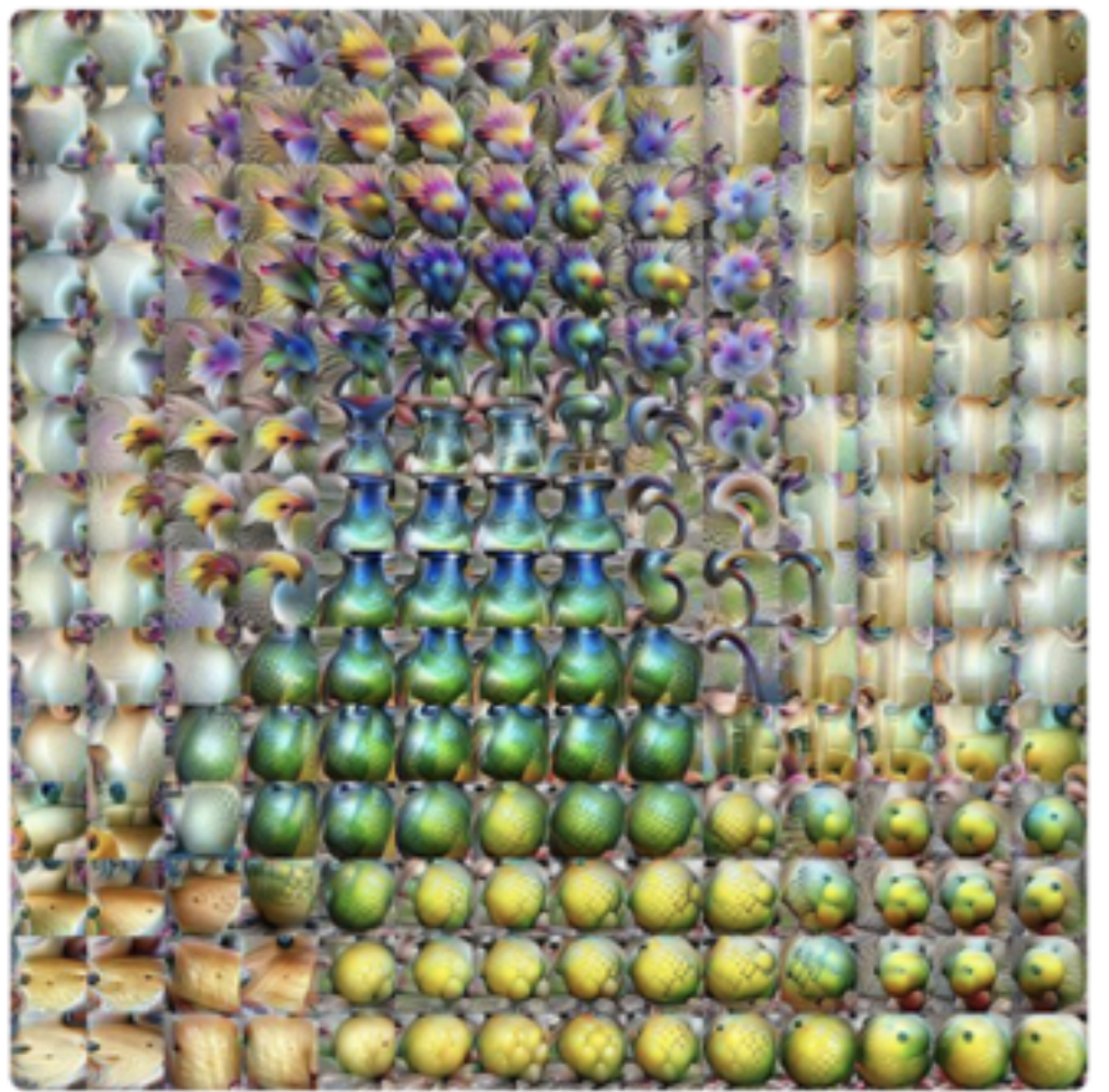
MIXED3A



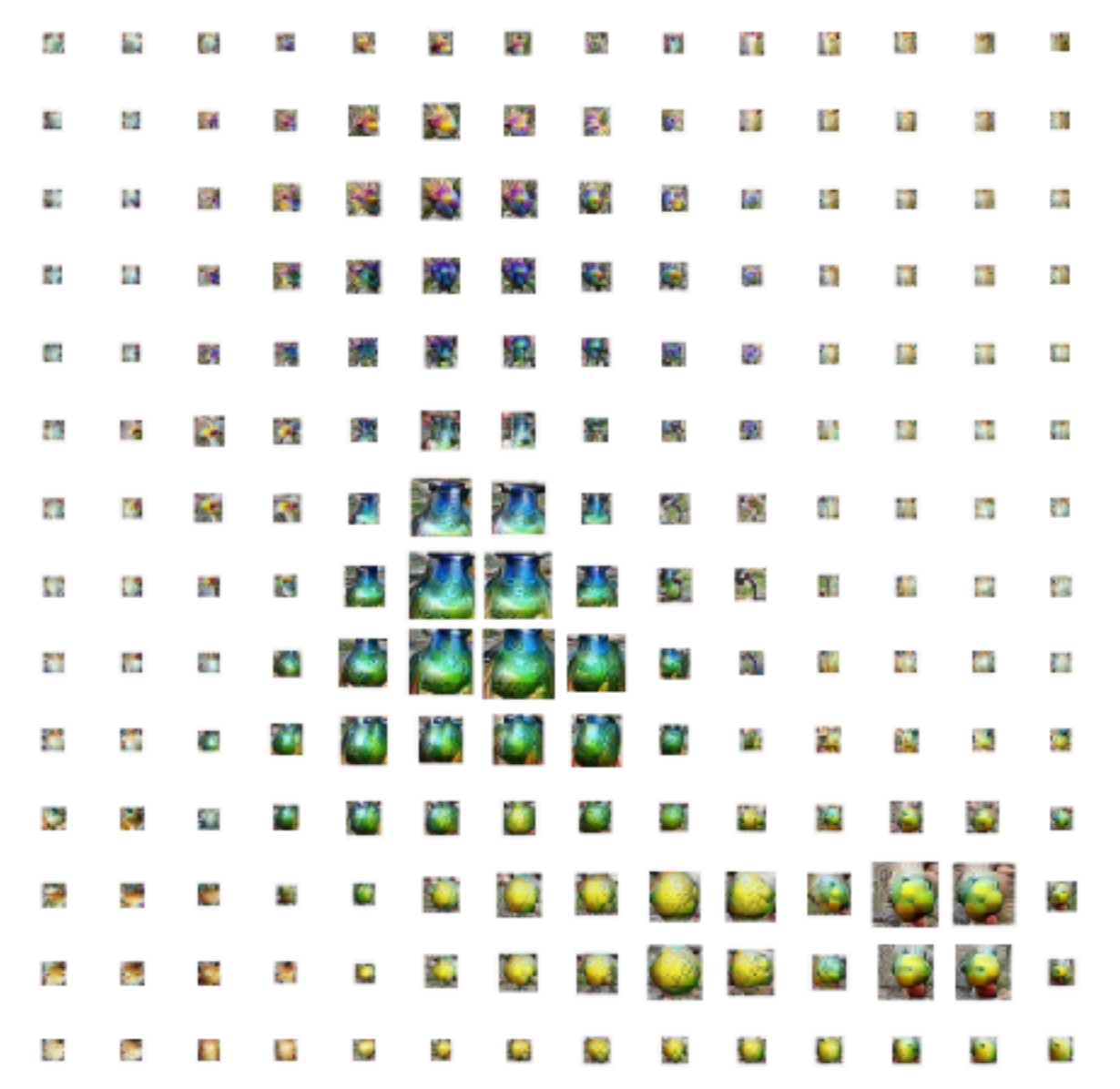
MIXED4A



MIXED4A



MIXED4D



MIXED4D

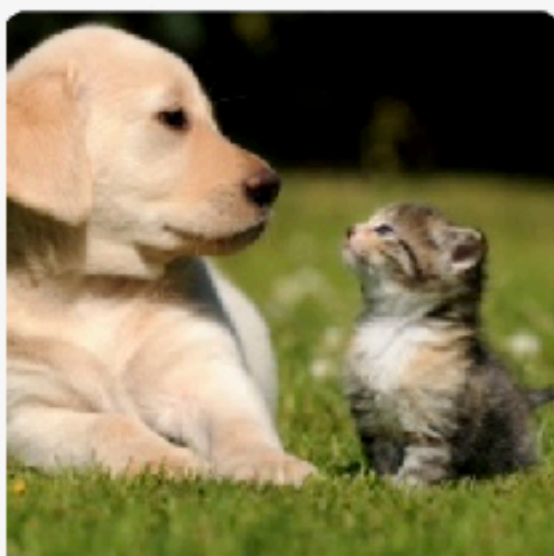


MIXED5A



MIXED5A

### INPUT IMAGE



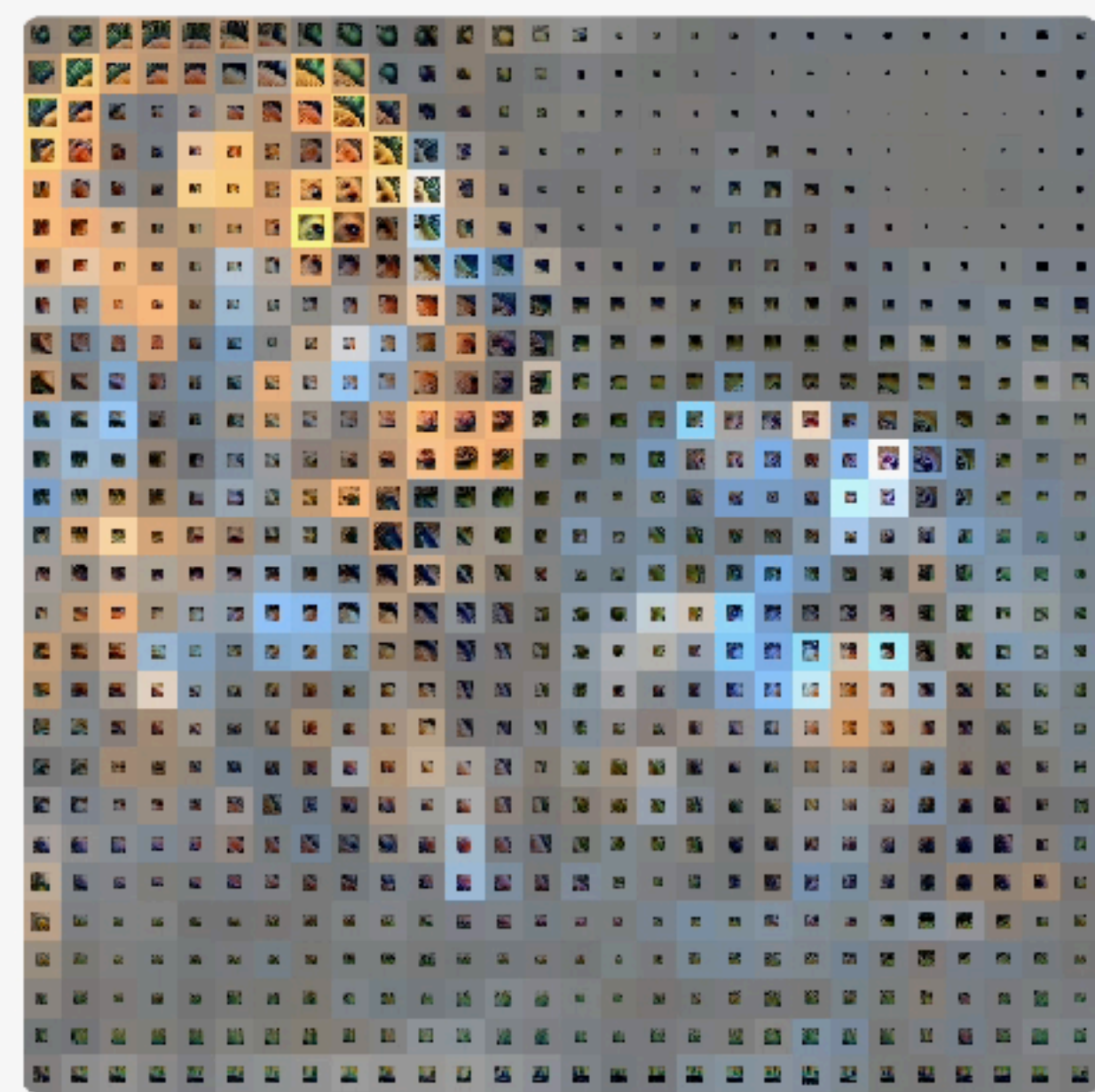
### OUTPUT CLASSES

Labrador Retriever	██████████
Golden Retriever	█
Tennis Ball	
Rhodesian Ridgeback	
Appenzeller	

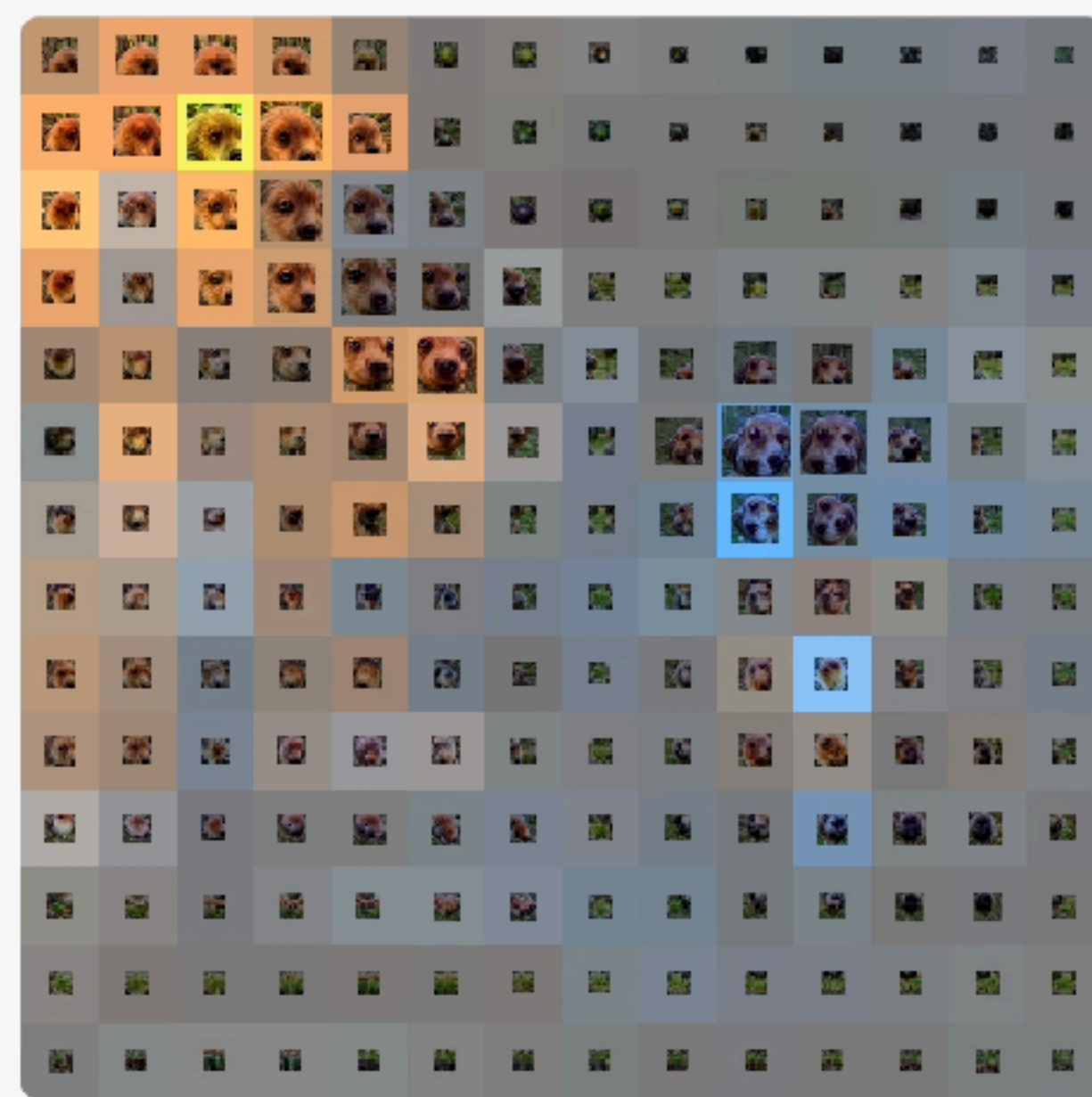
### OUTPUT FACTORS

Labrador Retriever	██████████	Tiger	██████████
Golden Retriever	██████████	Tiger Cat	██████████
Beagle	██████████	Lynx	██████████
Kuvasz	██████████	Collie	██████████
Redbone	██████████	Border Collie	██████████

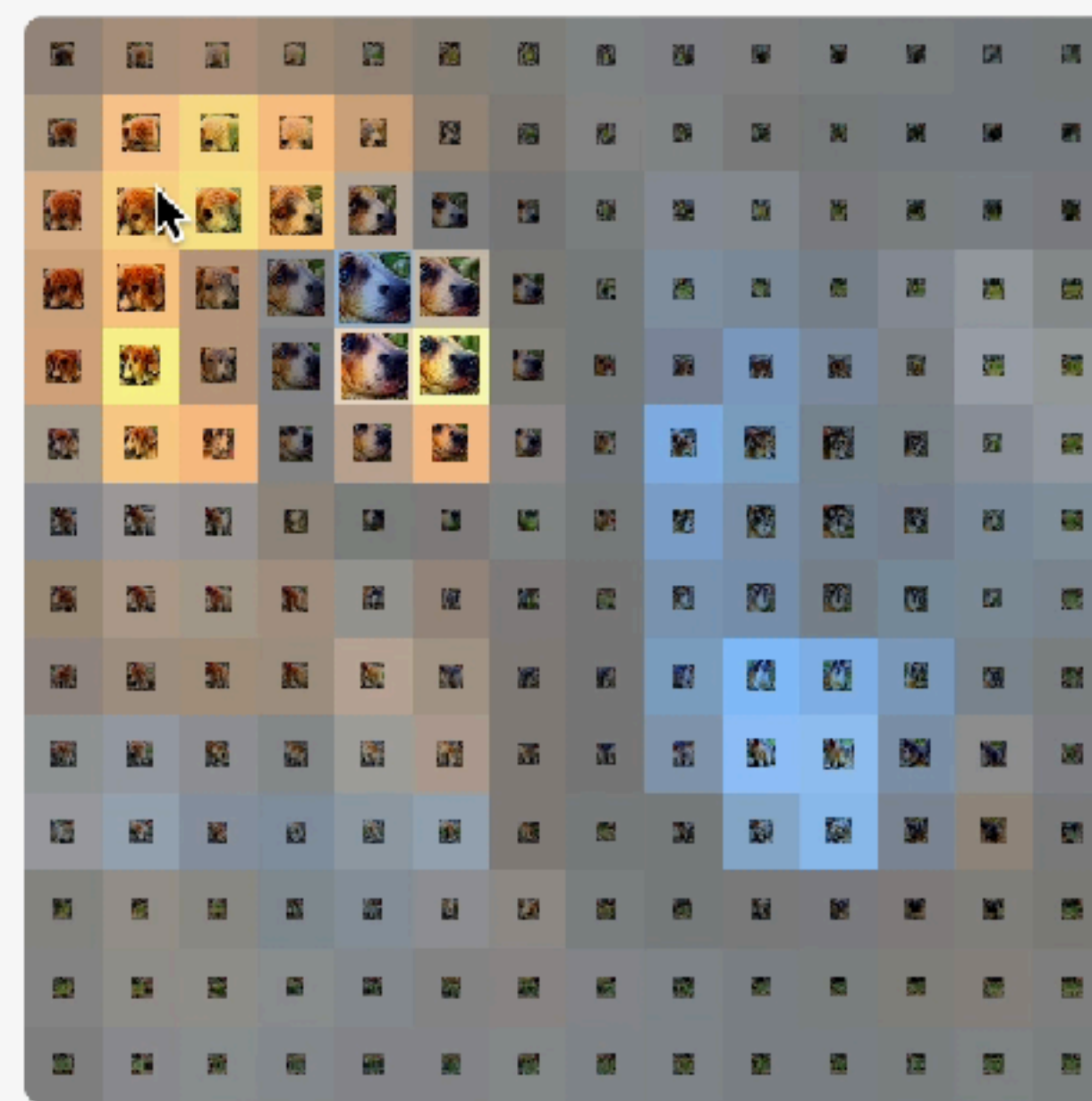
mixed3a



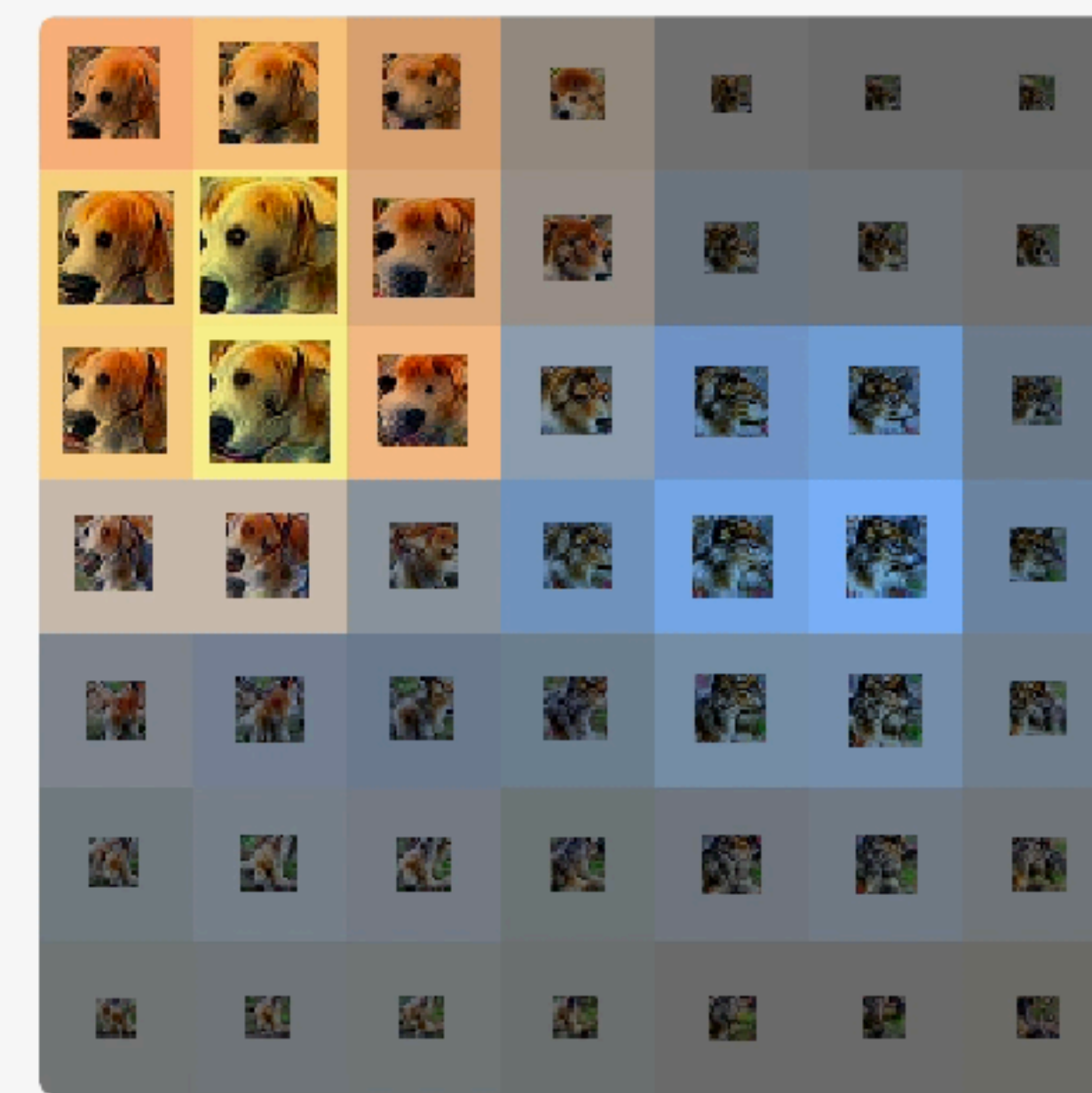
mixed4a



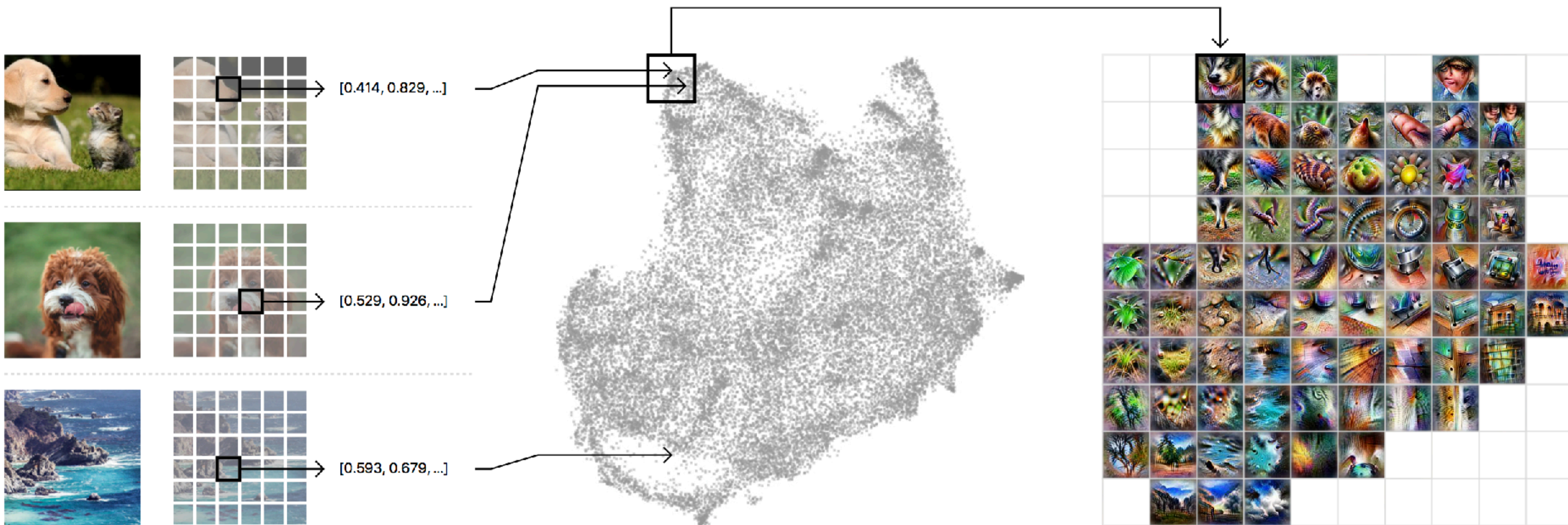
mixed4d



mixed5a







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](#)





