

# Network Dissection: Quantifying Interpretability of Deep Visual Representations

By David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, Antonio Torralba

CS 381V

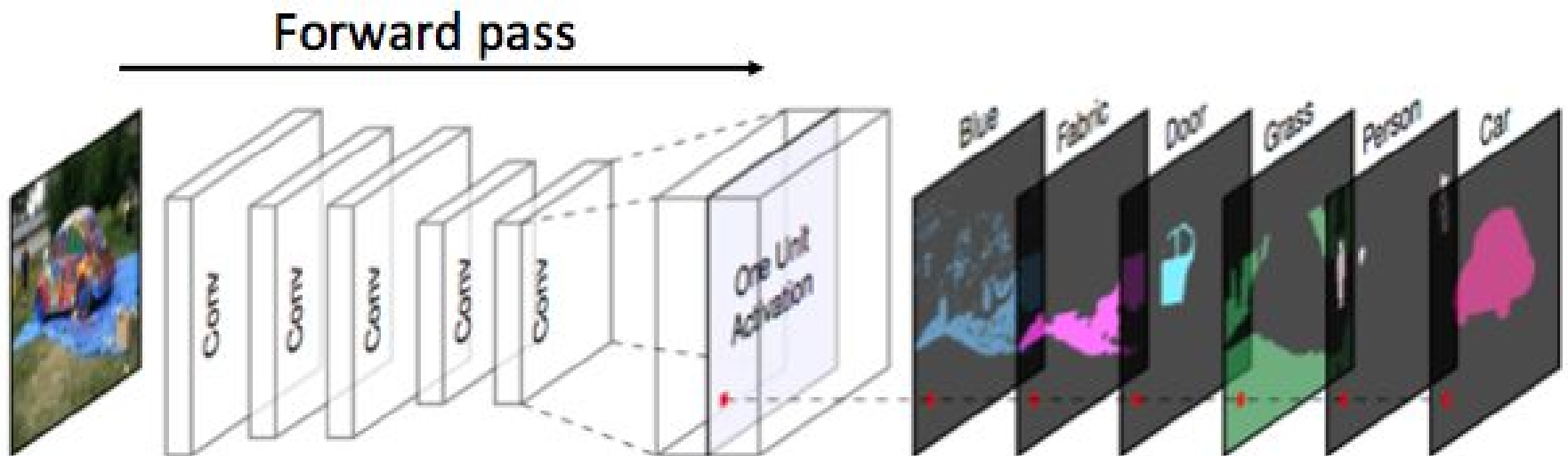
Thomas Crosley and Wonjoon Goo



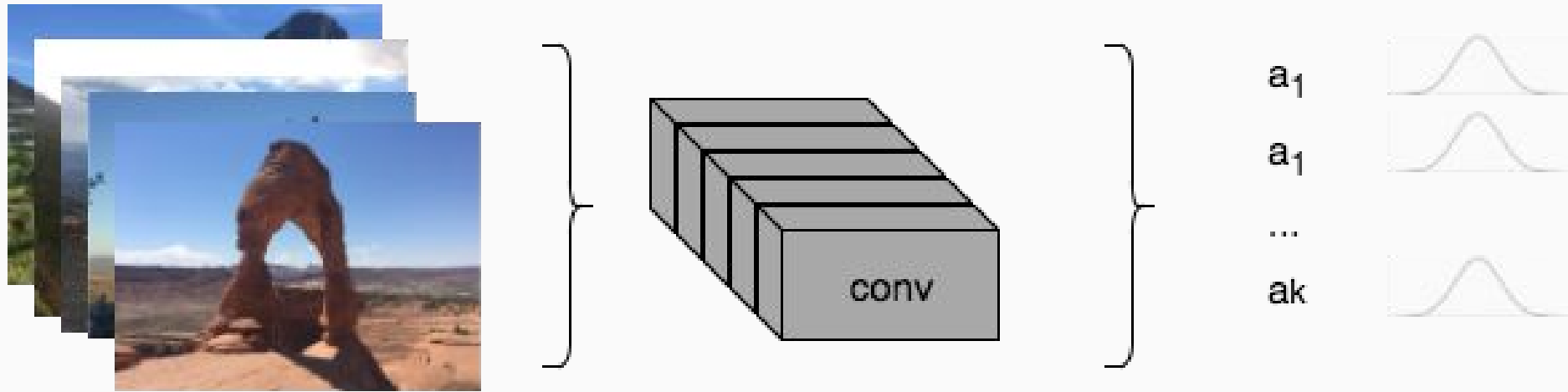
# Detectors

# Network Dissection

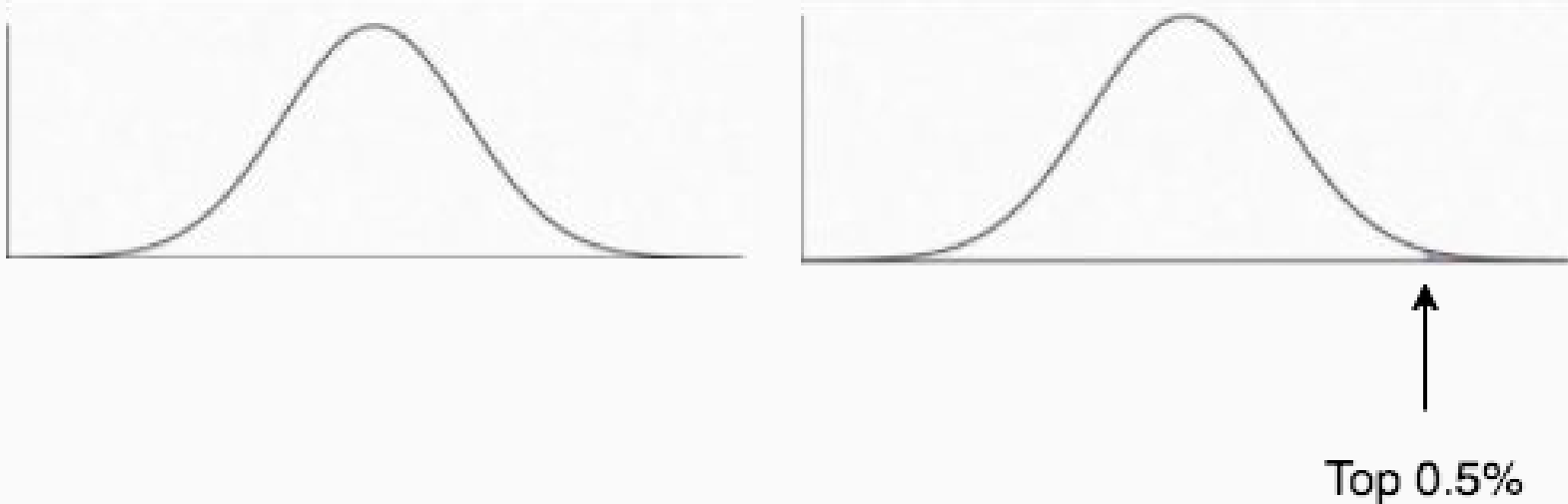
Quantifying the interpretability of units through segmentation



# Unit Distributions



- Compute internal activations for entire dataset
- Gather distribution for each unit across dataset

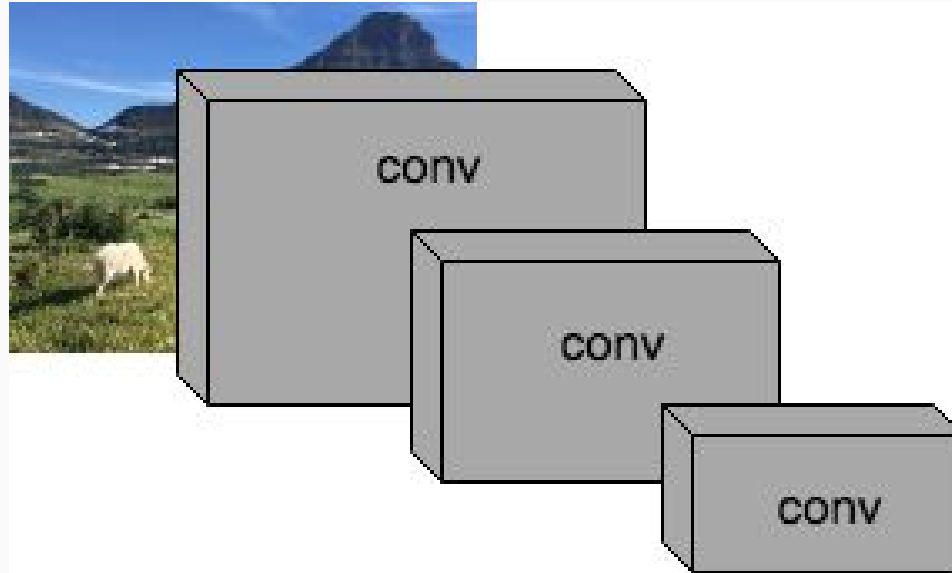


- Compute  $T_k$  such that  $P(a_k > T_k) = 0.005$
- $T_k$  is considered the top-quantile
- Detected regions at test time are those with

$$a_k > T_k$$

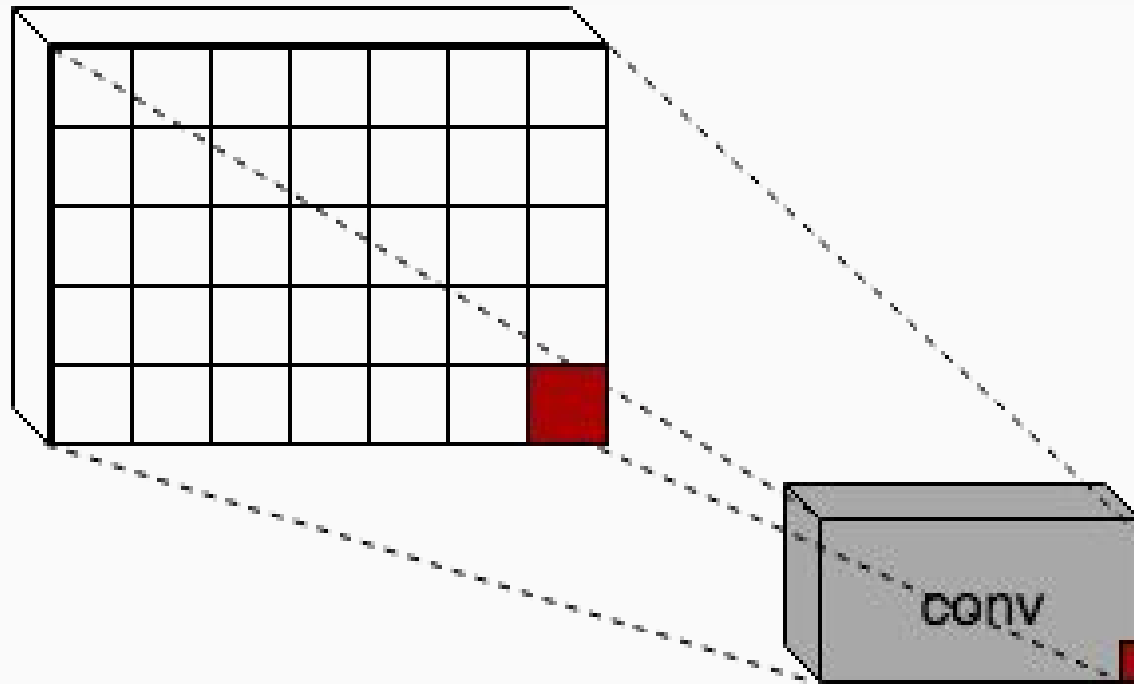
- Score of each unit is its IoU with the label
- Detectors are selected with IoU above a threshold
- Threshold is  $U_{k,c} > 0.04$ .

$$IoU_{k,c} = \frac{\sum |M_k(\mathbf{x}) \cap L_c(\mathbf{x})|}{\sum |M_k(\mathbf{x}) \cup L_c(\mathbf{x})|}$$



- Compute activation map  $a_k$  for all  $k$  neurons in the network

## Scaling Up



- Scale each unit's activation up to the original image size
- Call this the mask-resolution  $S_K$
- Use bi-linear interpolation



# Thresholding

99	100	100	70	68
99	100	101	67	65
98	102	100	98	60
82	97	99	97	52
75	70	72	70	45

$S_k$

$T_k=95$

1	1	1	0	0
1	1	1	0	0
1	1	1	1	0
0	1	1	1	0
0	0	0	0	0

$M_k$

- Now make the binary segmentation mask  $M_k$
- $M_k = S_k > T_k$

# Experiment: Detector Robustness

- Interest in adversarial examples
- Invariance to noise
- Composition by parts or statistics

# Noisy Images



+ Unif[0, 1]



+ 5 \* Unif[0, 1]

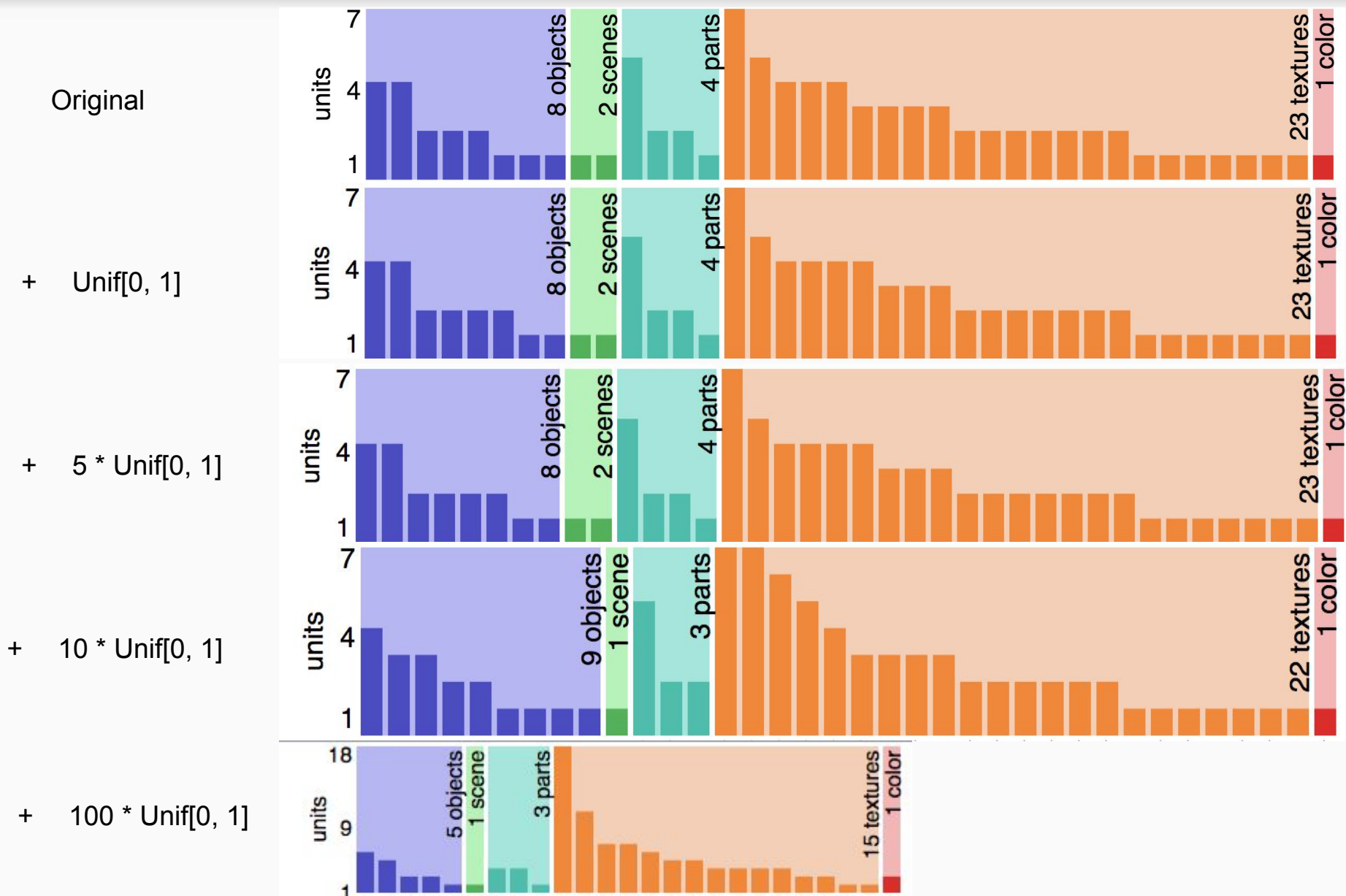


+ 10 \* Unif[0, 1]

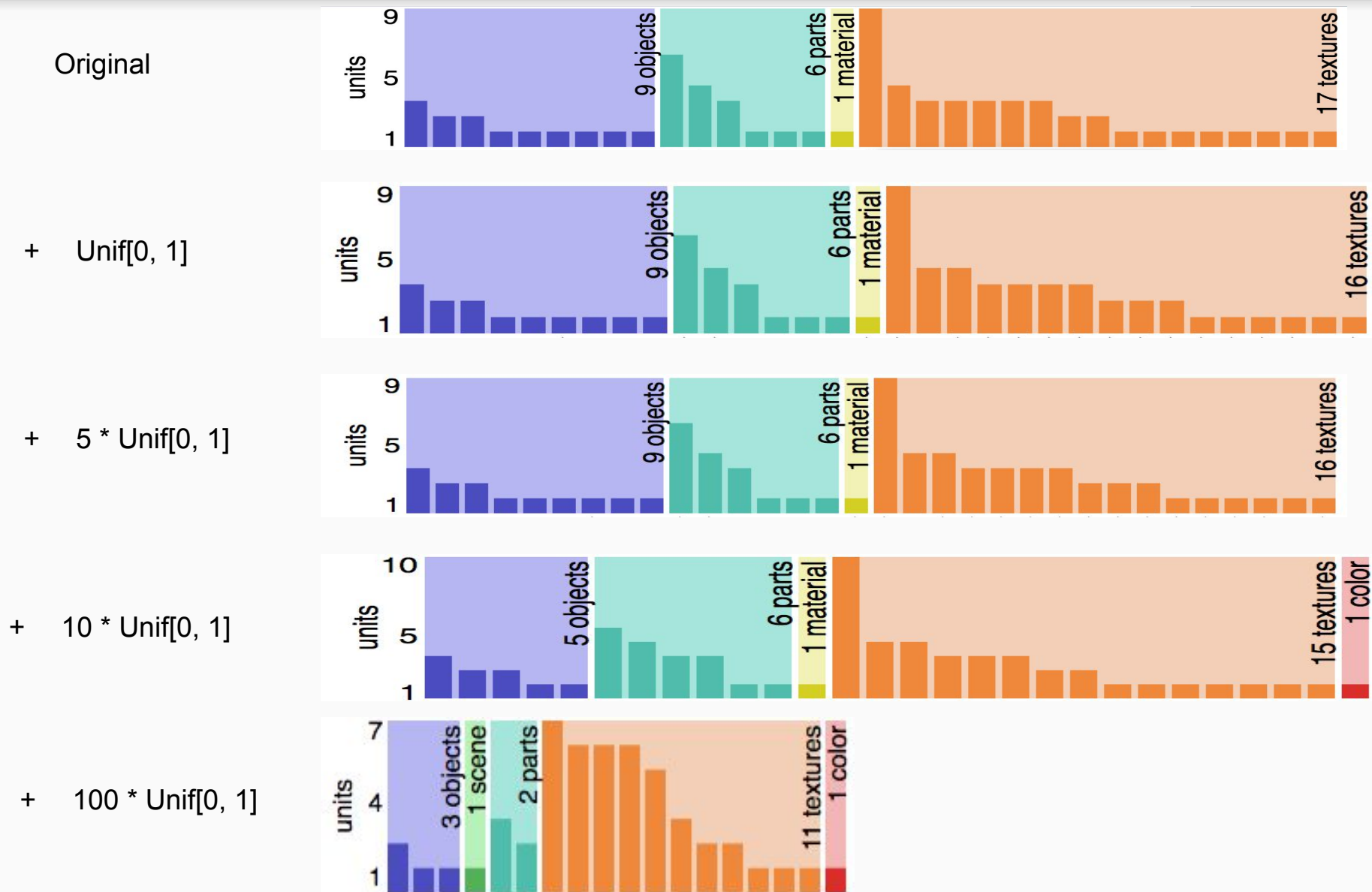


+ 100 \* Unif[0, 1]

# Conv3

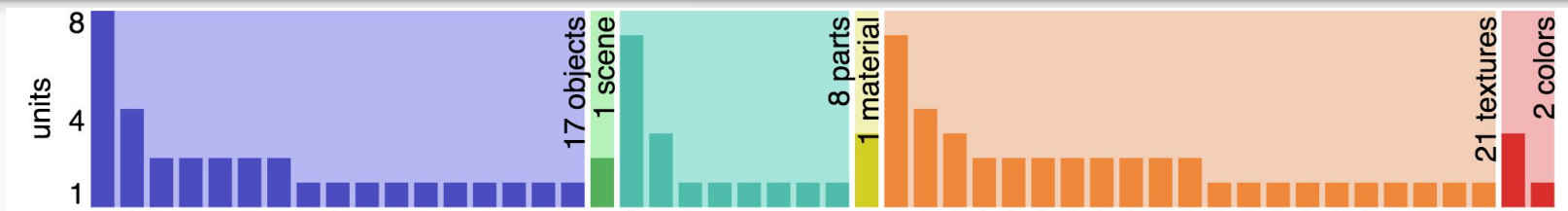


# Conv4

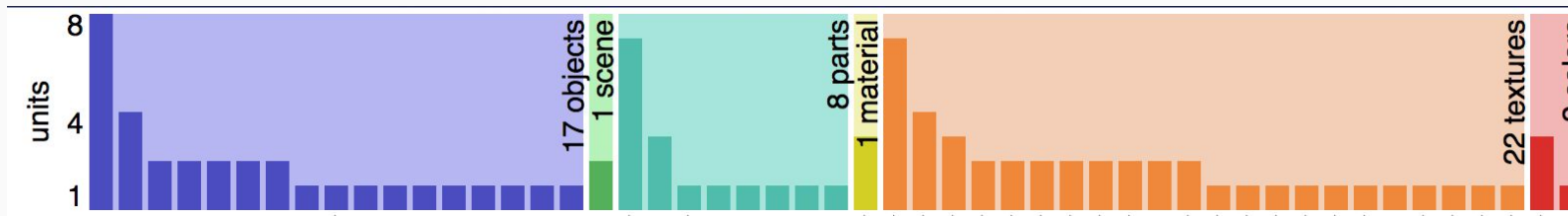


# Conv5

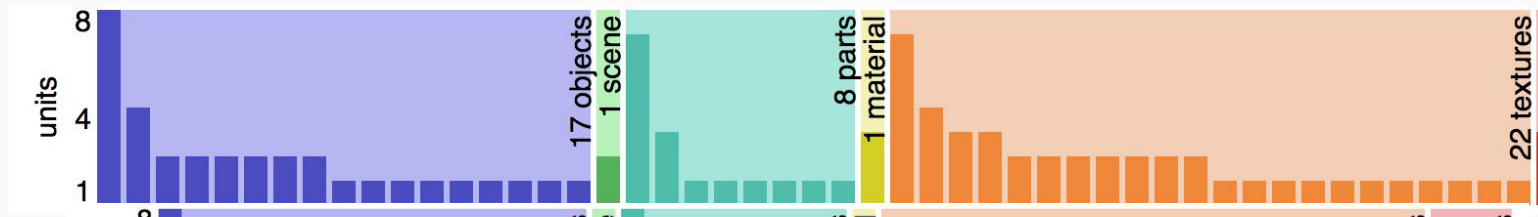
Original



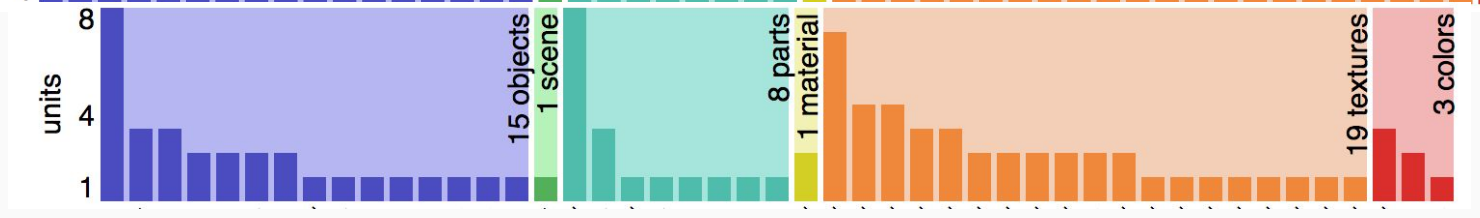
+ Unif[0, 1]



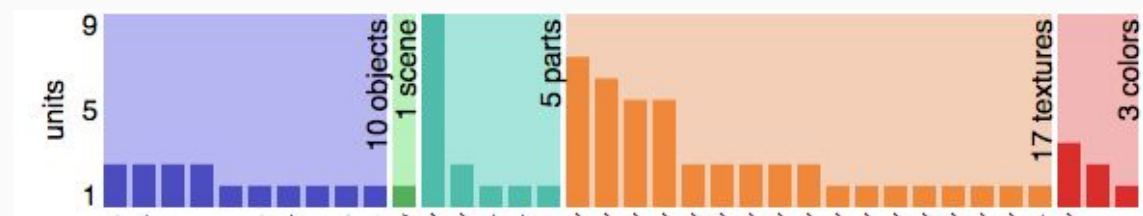
+ 5 \* Unif[0, 1]



+ 10 \* Unif[0, 1]



+ 100 \* Unif[0, 1]



# Rotated Images

Original



10 degrees



45 degrees



90 degrees

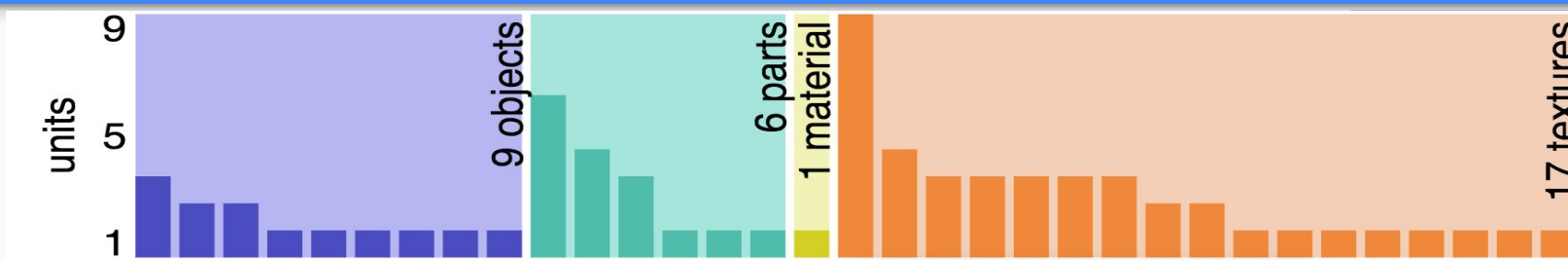




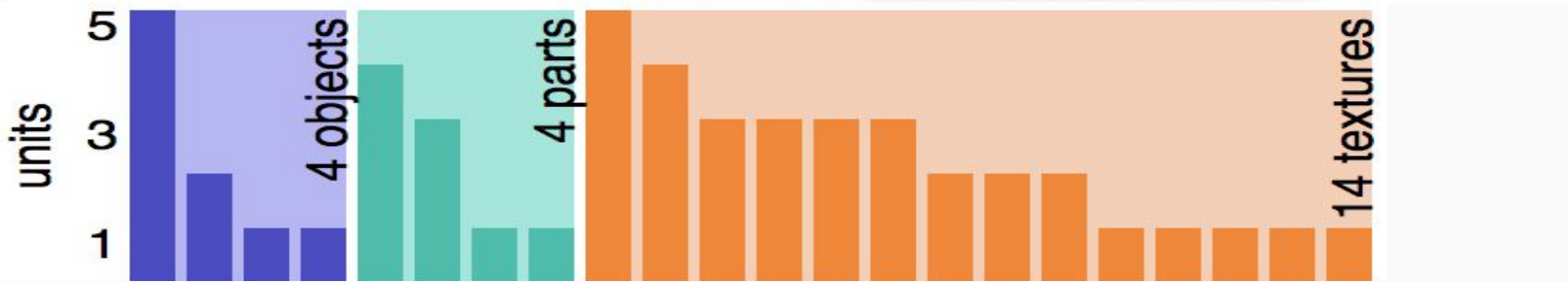


# conv4

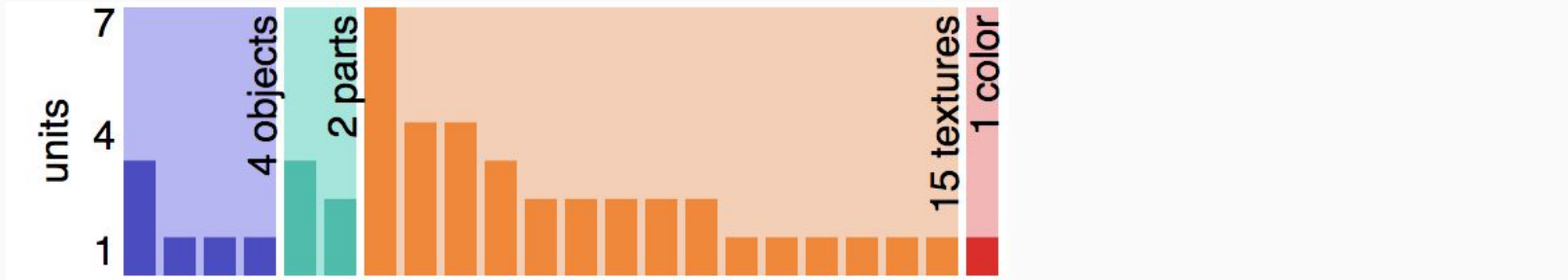
Original



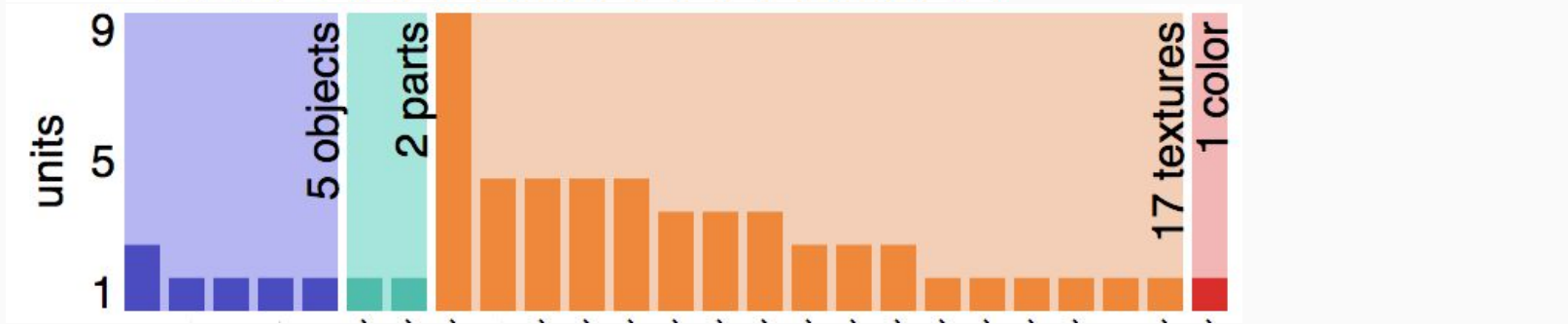
10 degrees



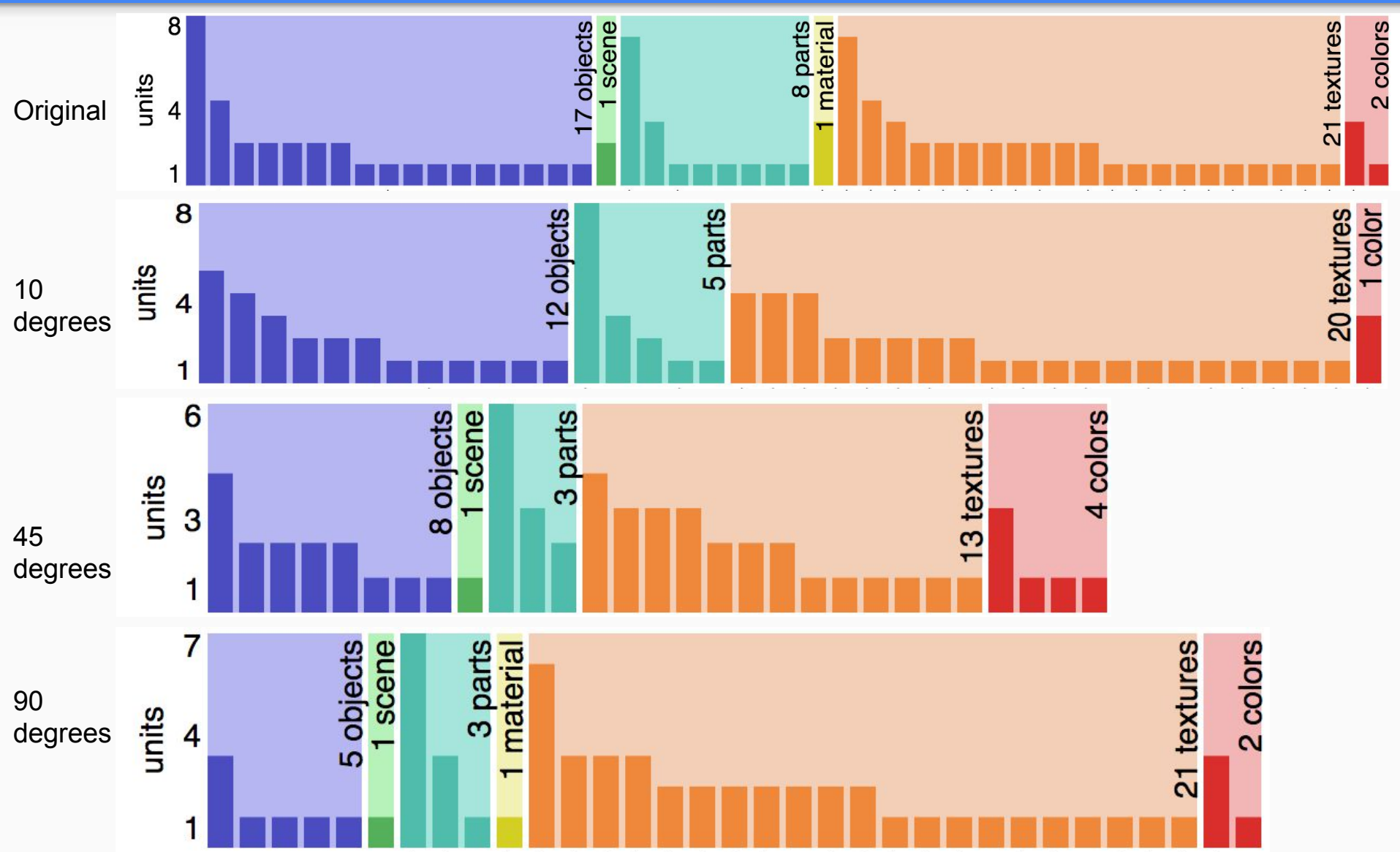
45 degrees



90 degrees



# conv5



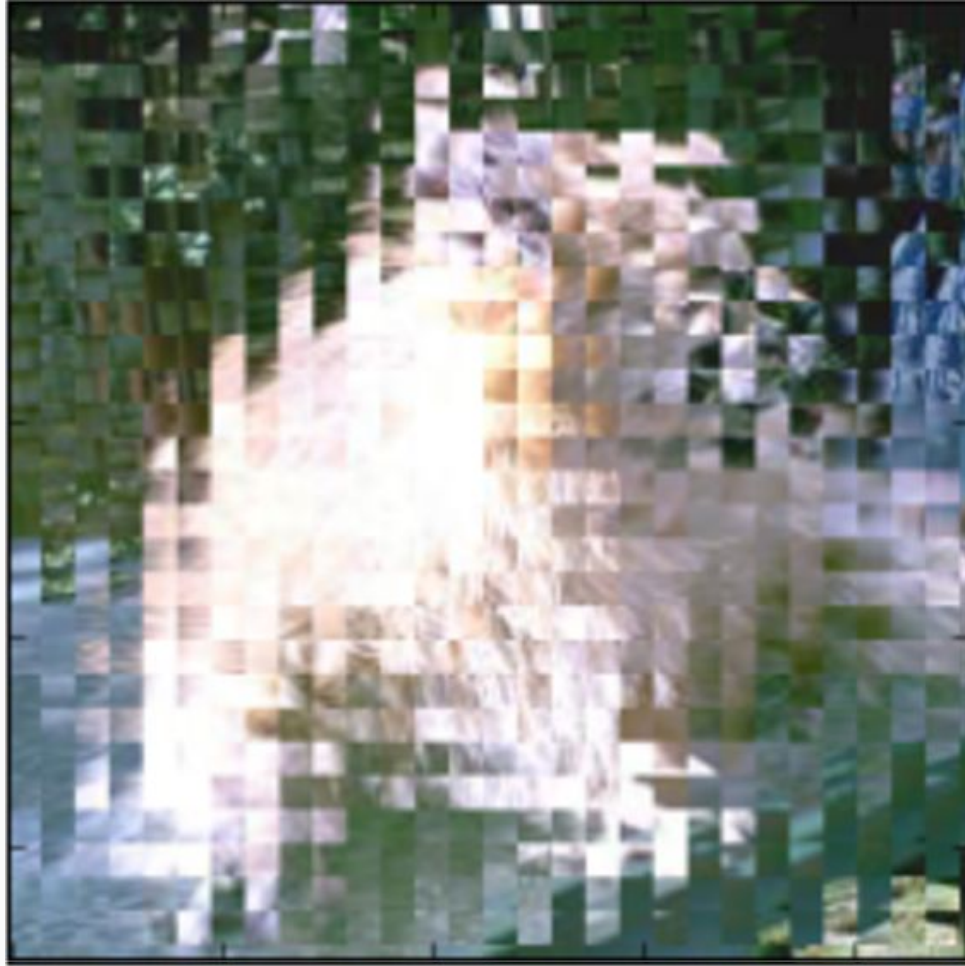
## Rearranged Images



## Rearranged Images

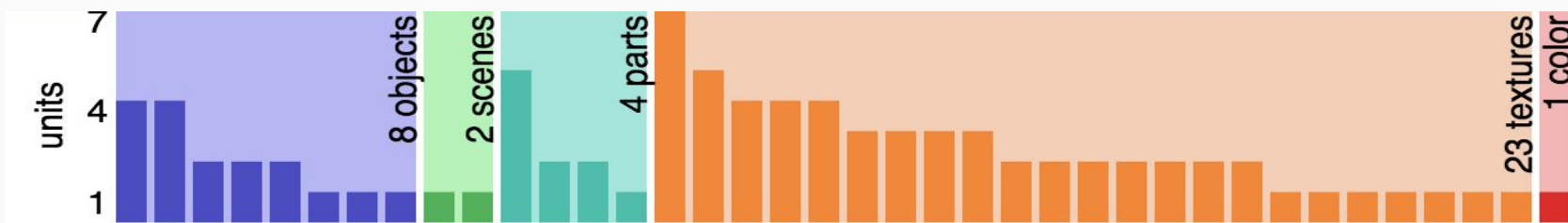


# Rearranged Images

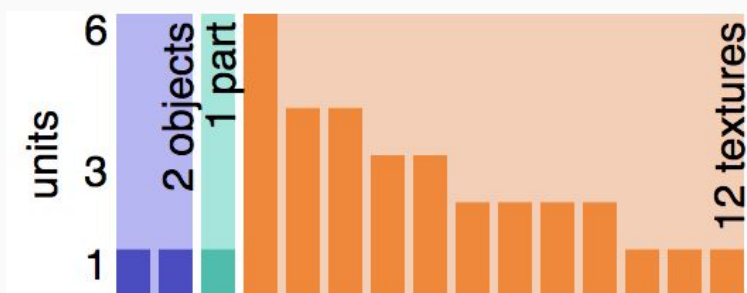


# Conv3

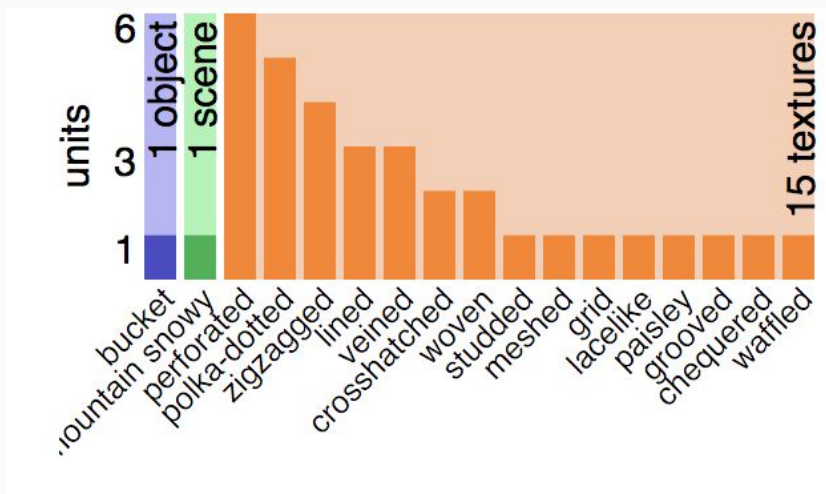
Original



4x4 Patches

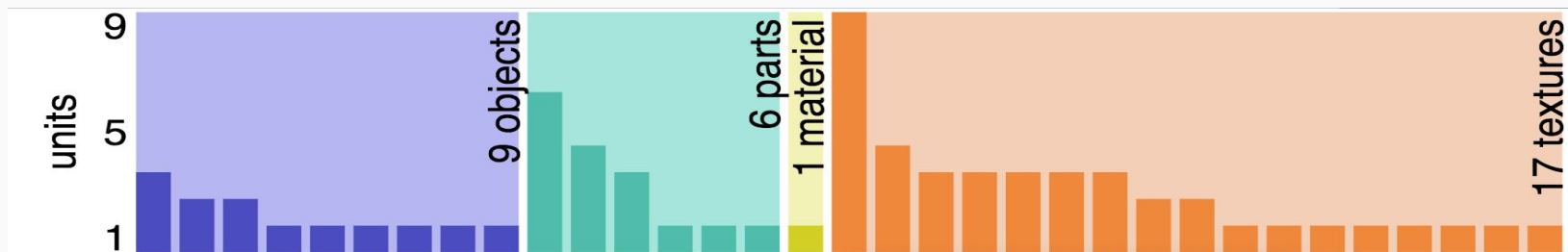


8x8 Patches

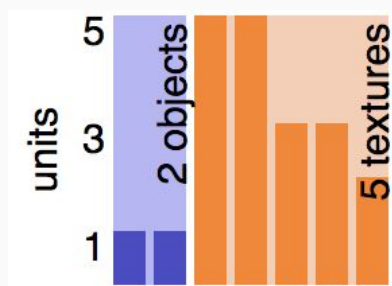


# Conv4

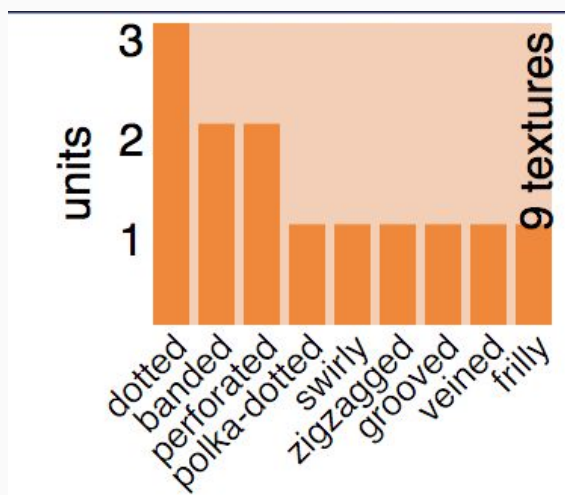
Original



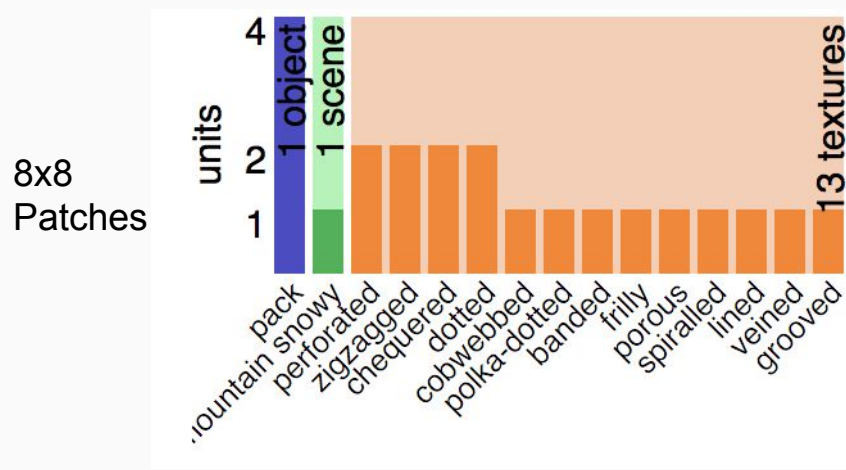
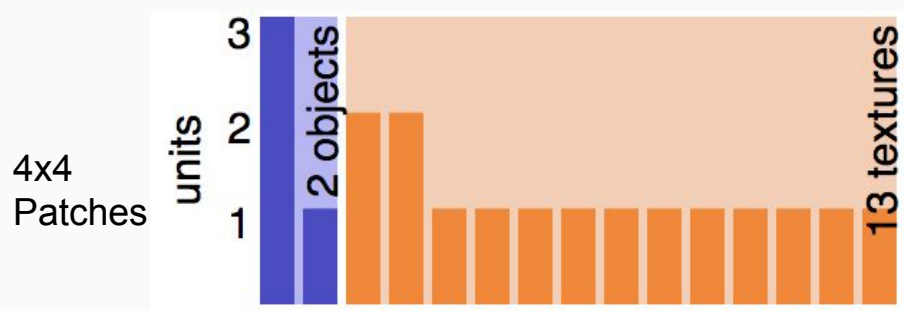
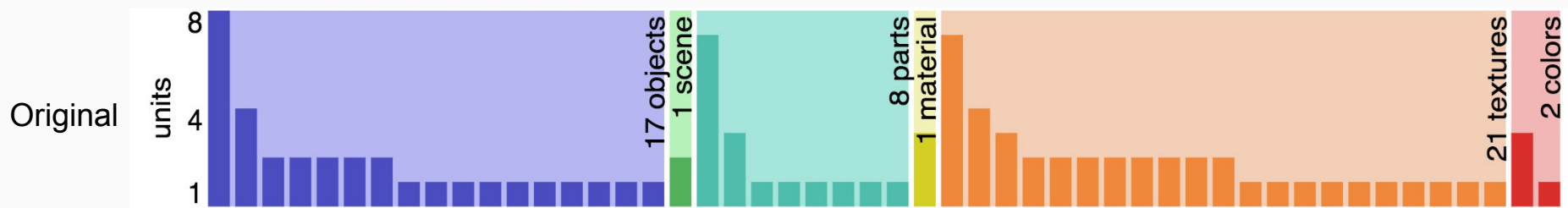
4x4  
Patches



8x8  
Patches



# Conv5



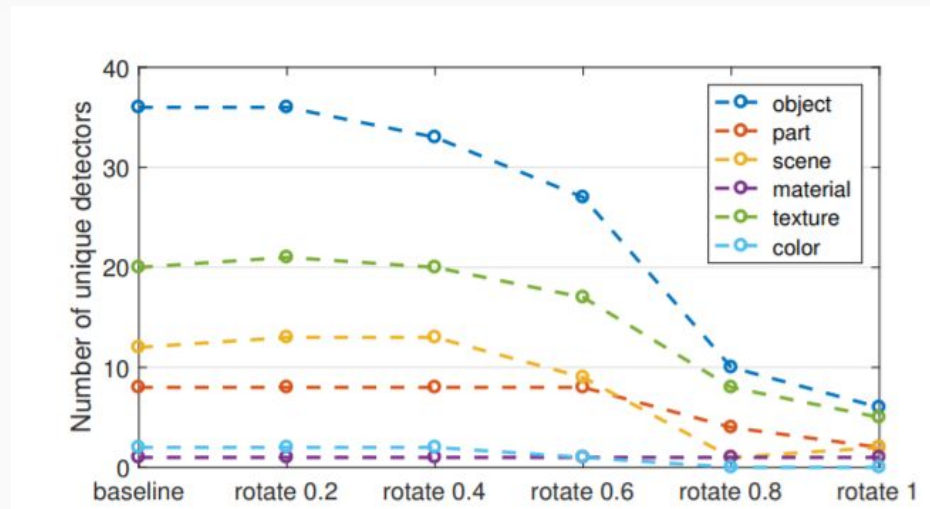


# Axis-Aligned Interpretability

- Hypothesis 1:
  - A linear combination of high level units serves just same or better
  - No specialized interpretation for each unit
- Hypothesis 2: (the authors' argument)
  - A linear combination will degrade the interpretability
  - Each unit serves for unique concept

How similar is the way CNN learns to human?

# Axis-Aligned Interpretability Result from the Authors



- It seems valid argument, but is it the best way to show?
- Problems
  - It depends on a rotation matrix used for test
  - A 90 degree rotation between two axis, does not affect the number of unique detectors
  - The test should be done multiple times and report the means and stds.

# Experiment: Axis-Aligned Interpretability

## Is it really axis aligned?

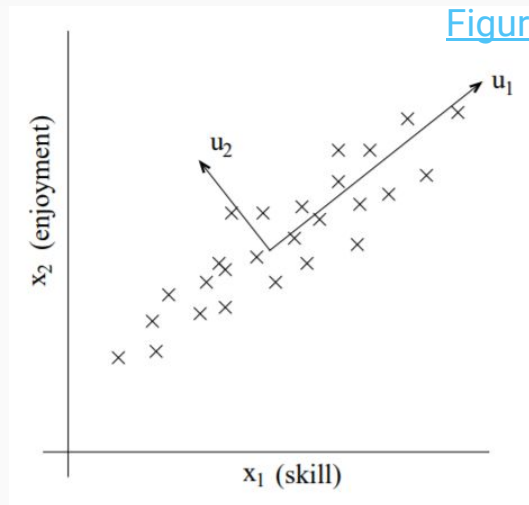


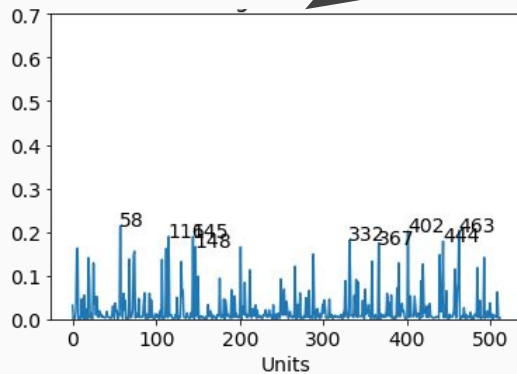
Figure: From Andrew Ng's lecture note on PCA

- Principle Component Analysis (PCA)
  - Find orthonormal vectors explaining samples the most
  - The projections to the vector  $u_1$  have higher variance
- ❖ Argument: a unit itself can explain a concept
  - Projections to unit vectors should have higher variance
  - *Principal axis (Loading)* from PCA should be similar to one of the unit vectors

## Our method

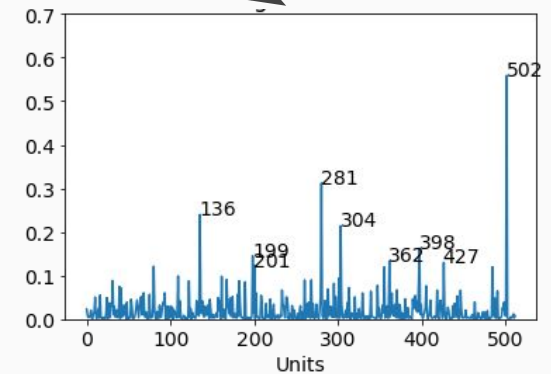
1. Calculate the mean and std. of each unit activation
2. Grab activations for a specific concept
3. Subtract mean and std from activations
4. Perform SVD
5. Print Loading

Hypothesis 1



The concept is interpreted with the combination of elementary basis

Hypothesis 2



The concept can be interpreted with an elementary basis (eg.  $e_{502} := (0, \dots, 0, 1, 0, \dots, 0)$ )

- Optimize target:

$$\max_{u: \|u\|=1} \frac{1}{m} \sum_{i=1}^m (x^{(i)T} u)^2$$

$$\Rightarrow \max_{u: \|u\|=1} \frac{1}{m} \sum_{i=1}^m (x^{(i)T} u)^T (x^{(i)T} u)$$

$$\Rightarrow \max_{u: \|u\|=1} \frac{1}{m} \sum_{i=1}^m (u^T x^{(i)}) (x^{(i)T} u)$$

$$\Rightarrow \max_{u: \|u\|=1} \frac{1}{m} \sum_{i=1}^m u^T x^{(i)} x^{(i)T} u$$

$$\Rightarrow \max_{u: \|u\|=1} u^T \left( \frac{1}{m} \sum_{i=1}^m x^{(i)} x^{(i)T} \right) u$$

$$\Rightarrow \max_{u: \|u\|=1} u^T \Sigma u, \quad \text{where } \Sigma = \frac{1}{m} \sum_{i=1}^m x^{(i)} x^{(i)T}$$

- With Lagrange multiplier:

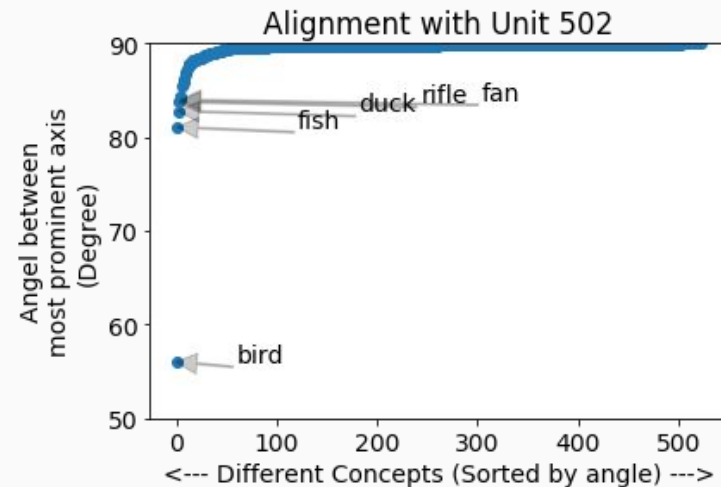
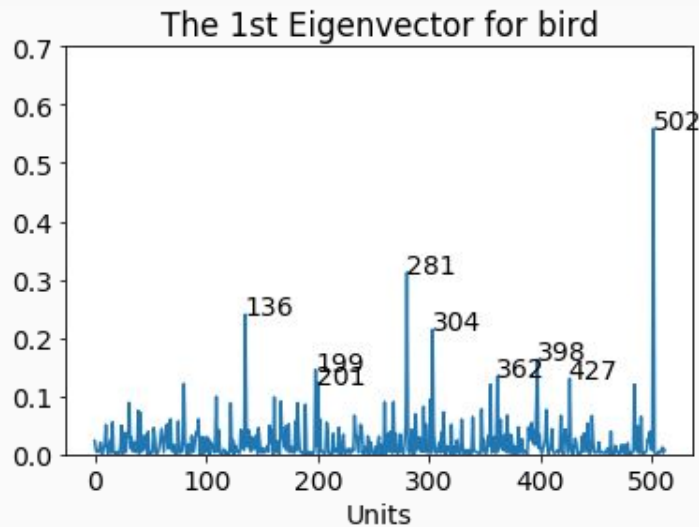
$$\Rightarrow \max u^T \Sigma u \quad \text{subject } u^T u = 1 \text{ (Because } \|u\| = 1)$$

$$\Rightarrow \mathcal{L}(u, \lambda) = u^T \Sigma u + \lambda(u^T u - 1)$$

$$\begin{aligned} \Rightarrow \frac{\partial \mathcal{L}(u, \lambda)}{\partial u} &= \frac{\partial (u^T \Sigma u + \lambda u^T u)}{\partial u} = \frac{\partial (u^T \Sigma u)}{\partial u} + \frac{\partial (\lambda u^T u)}{\partial u} \\ &= 2\Sigma u + 2\lambda u \stackrel{\text{Set}}{=} 0 \end{aligned}$$

- The eigenvector for the highest eigenvalue becomes principal axis (loading)

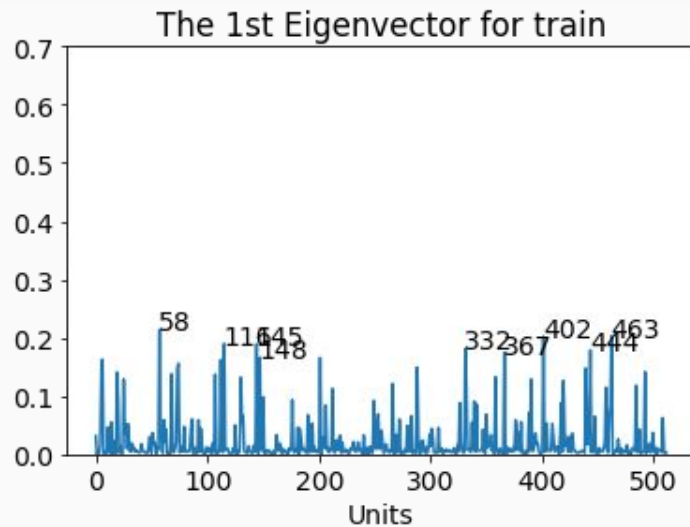
# PCA Results - Activations for Bird Concept



- Unit 502 stands high; concept bird is aligned to the unit
- Does Unit 502 only serve for concept Bird?
  - Yes
  - It does not stand for other concepts except bird
- Support Hypothesis 2

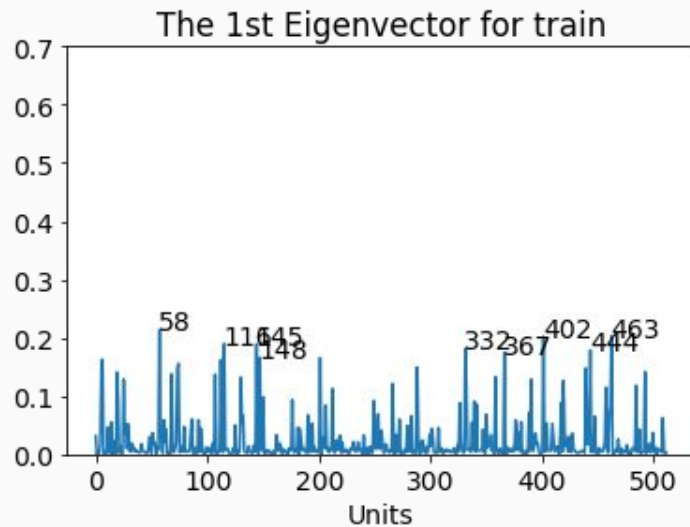


## PCA Results - Activations for Train Concept



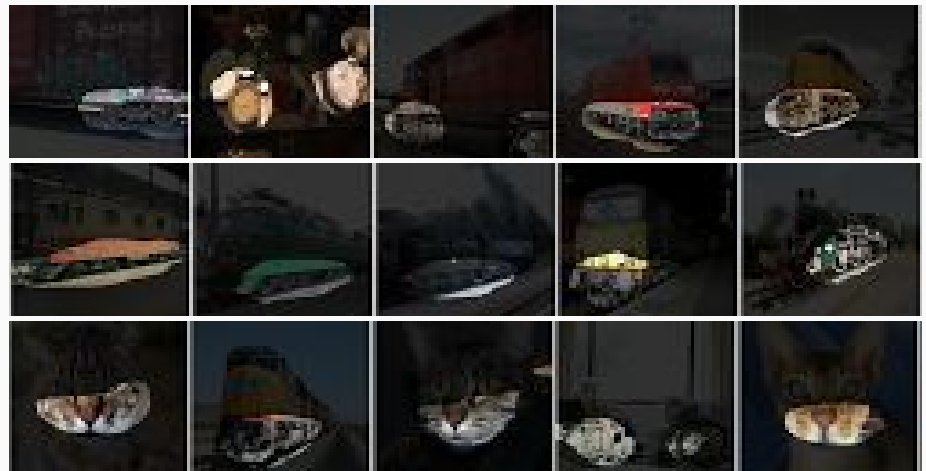
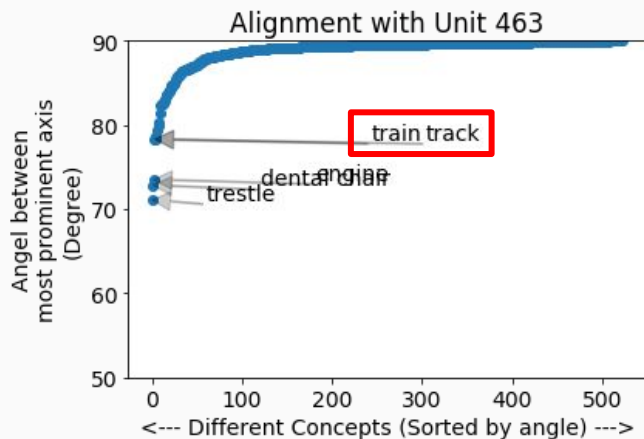
- No units stands out for concept train
  - Linear combination of them have better interpretability
  - Support Hypothesis 1

# PCA Results - Activations for Train Concept

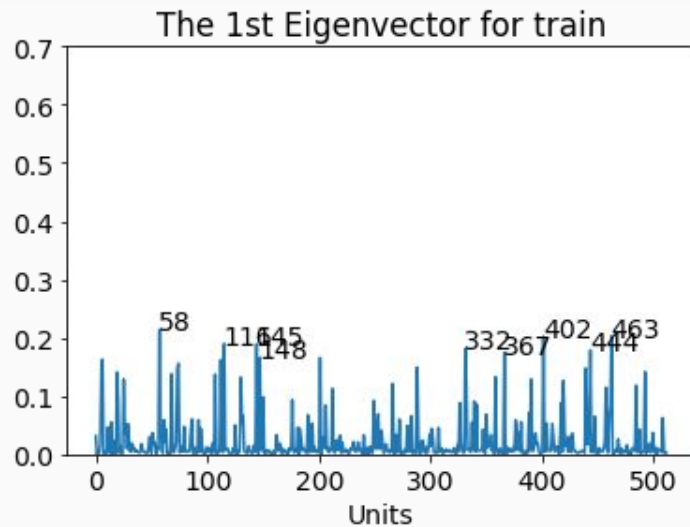


- No units stands out for concept train
  - Linear combination of them have interpretability

Some objects with circle and trestle?



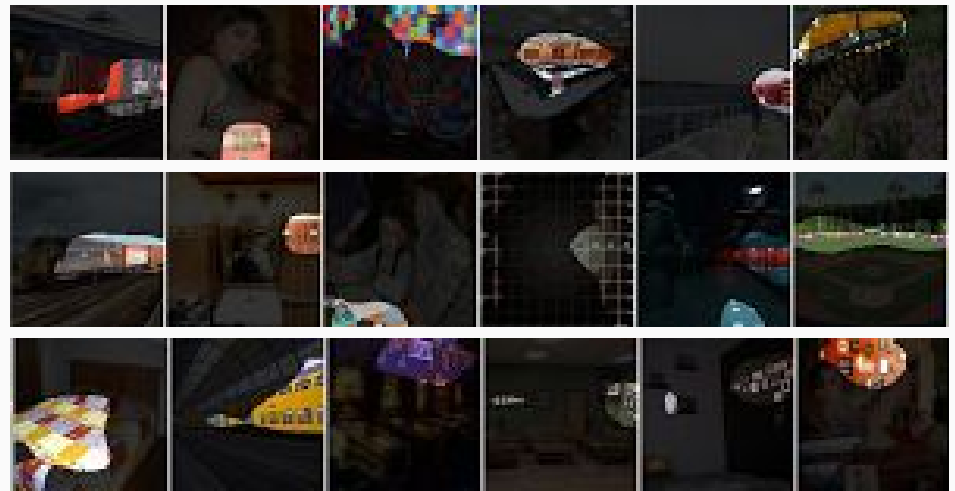
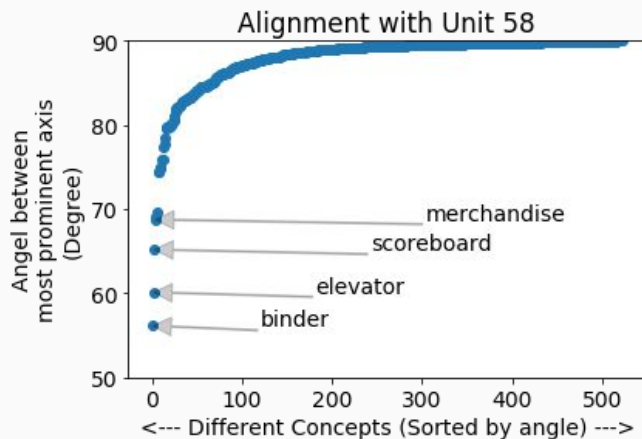
# PCA Results - Activations for Train Concept



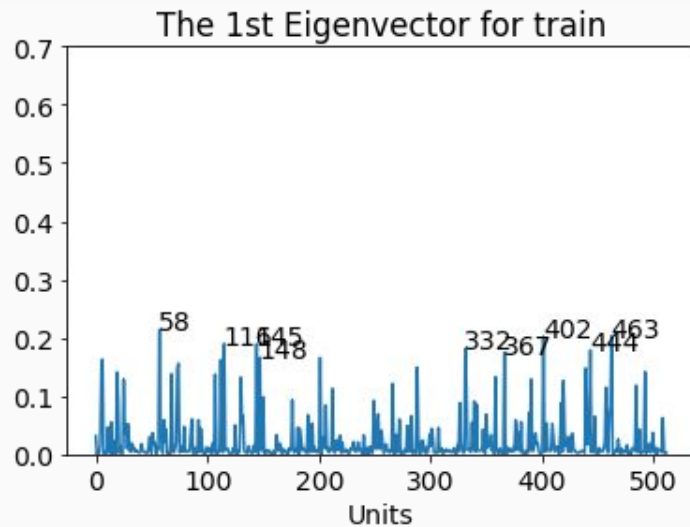
- No units stands out for concept train

The sequence of square boxes?

- Linear combination of them have interpretability

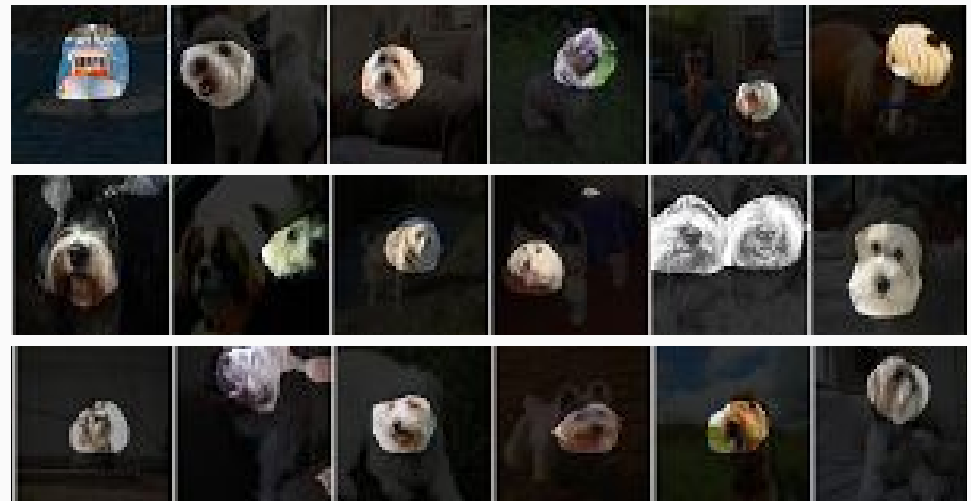
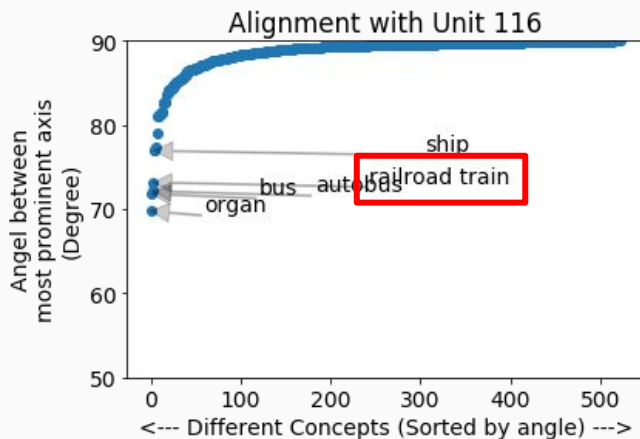


# PCA Results - Activations for Train Concept



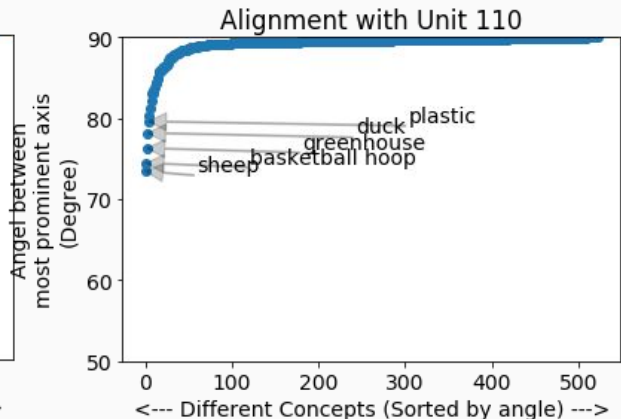
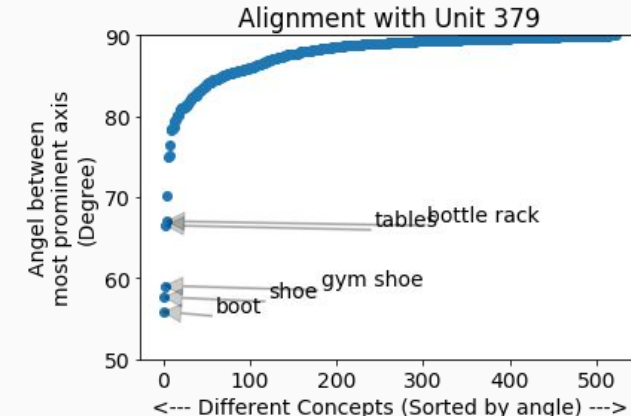
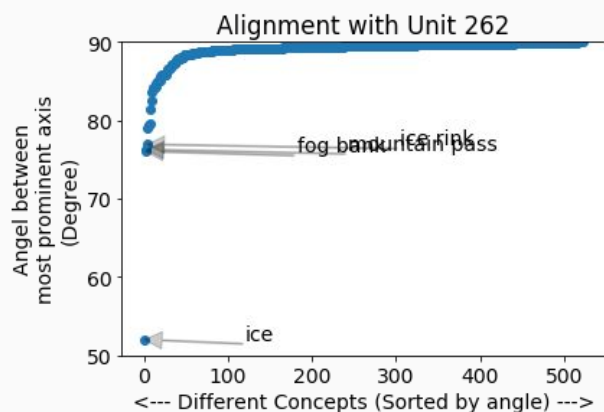
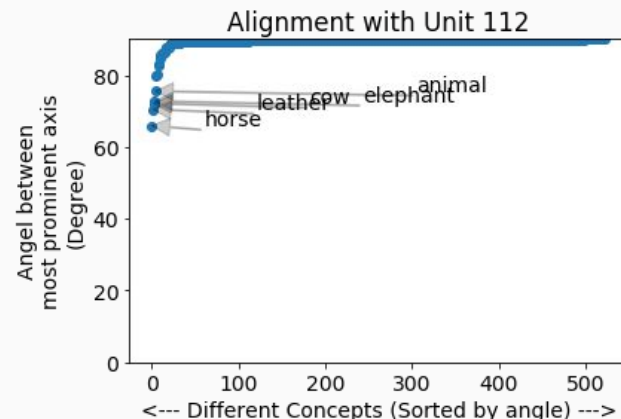
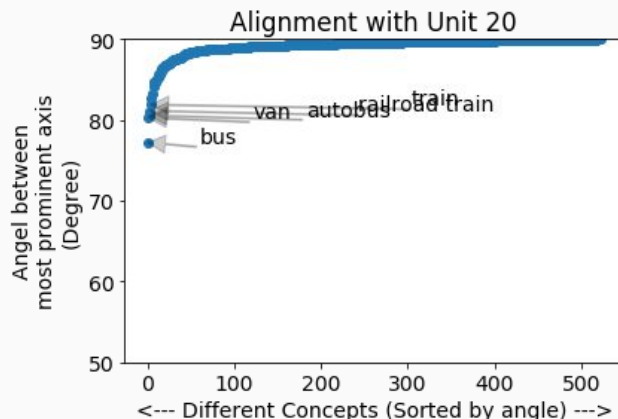
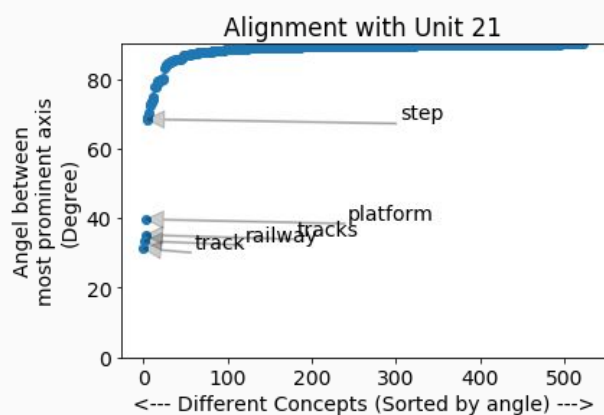
- No units stands out for concept train
  - Linear combination of them have interpretability

Dog face!



# Conclusion...?

- Actually, it seems mixed!
- CNN learns some human concepts naturally, but not always
  - It might highly correlated with the label we give



- What if we regularize the network to encourage its interpretability?
  - [Taxonomy-Regularized Semantic Deep Convolutional Neural Networks](#),  
Wonjoon Goo, Juyong Kim, Gunhee Kim, and Sung Ju Hwang, ECCV 2016

Thanks!

Any questions?