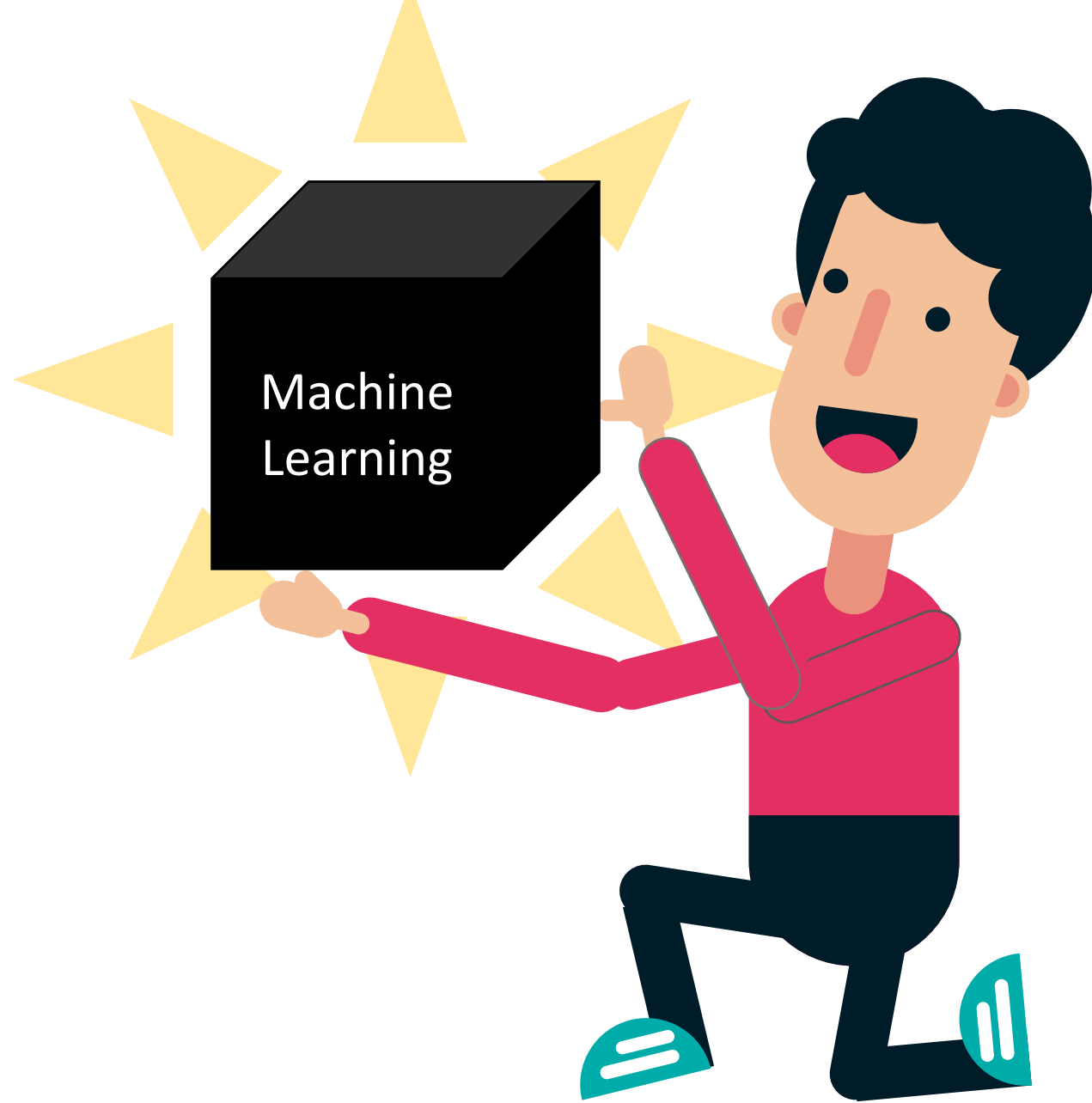


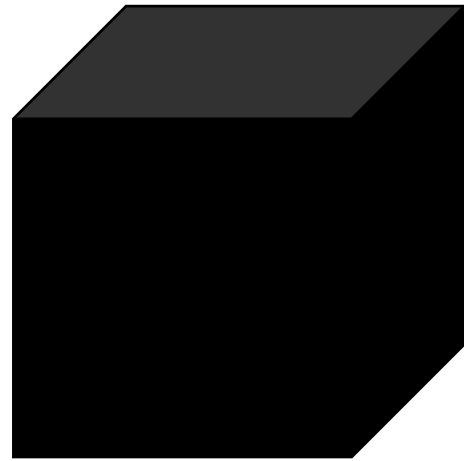
Visualizing Surrogate Decision Trees of Convolutional Neural Networks

Shichao Jia¹, Peiwen Lin¹, Zeyu Li¹, Jiawan Zhang¹ and Shixia Liu²





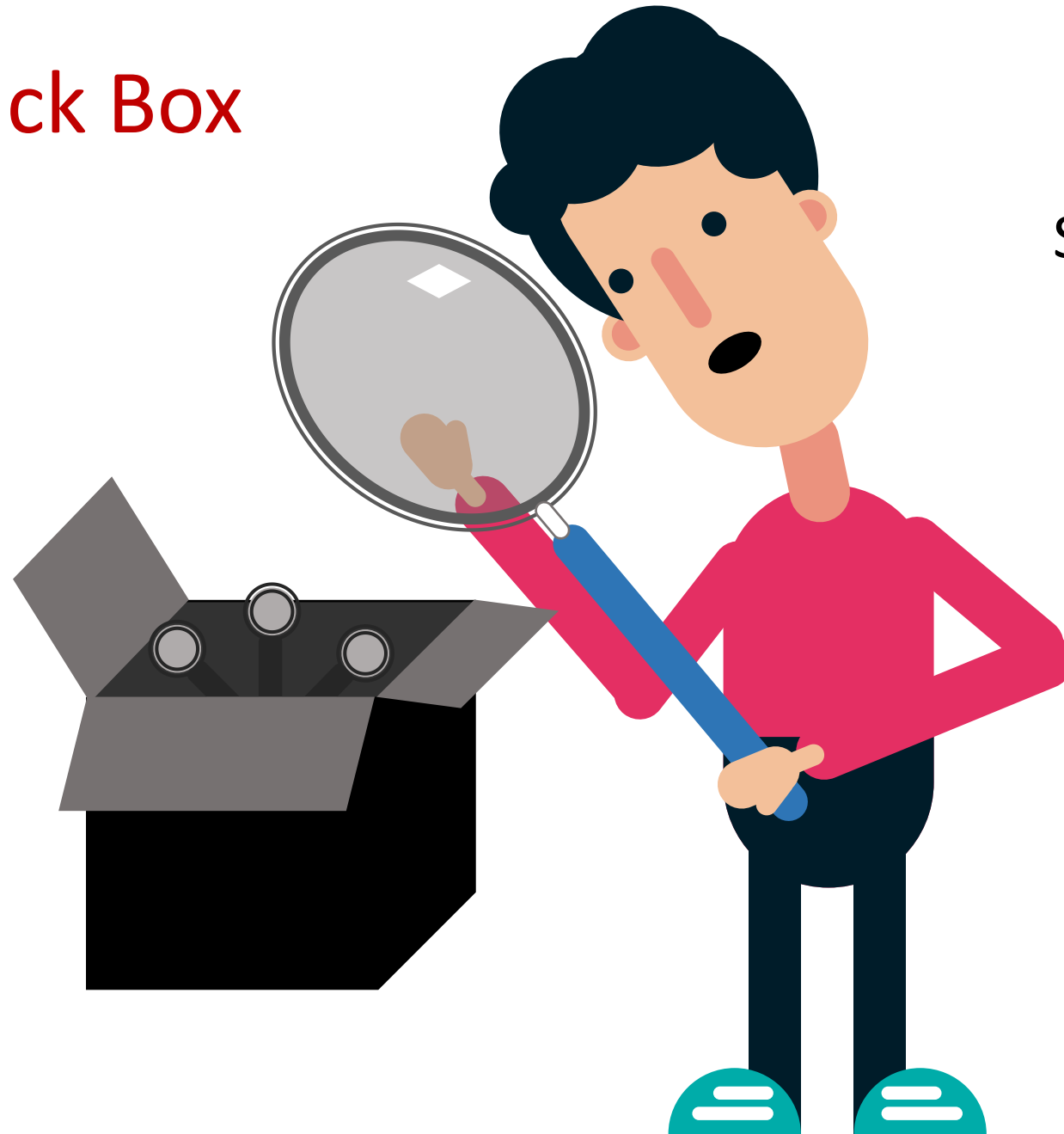
Interpretability Problem



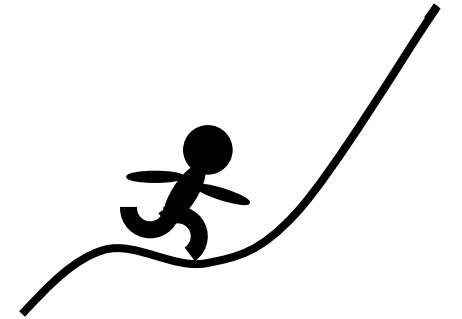
I do not understand !
I can not trust it !

Strategy 1

Open the Black Box



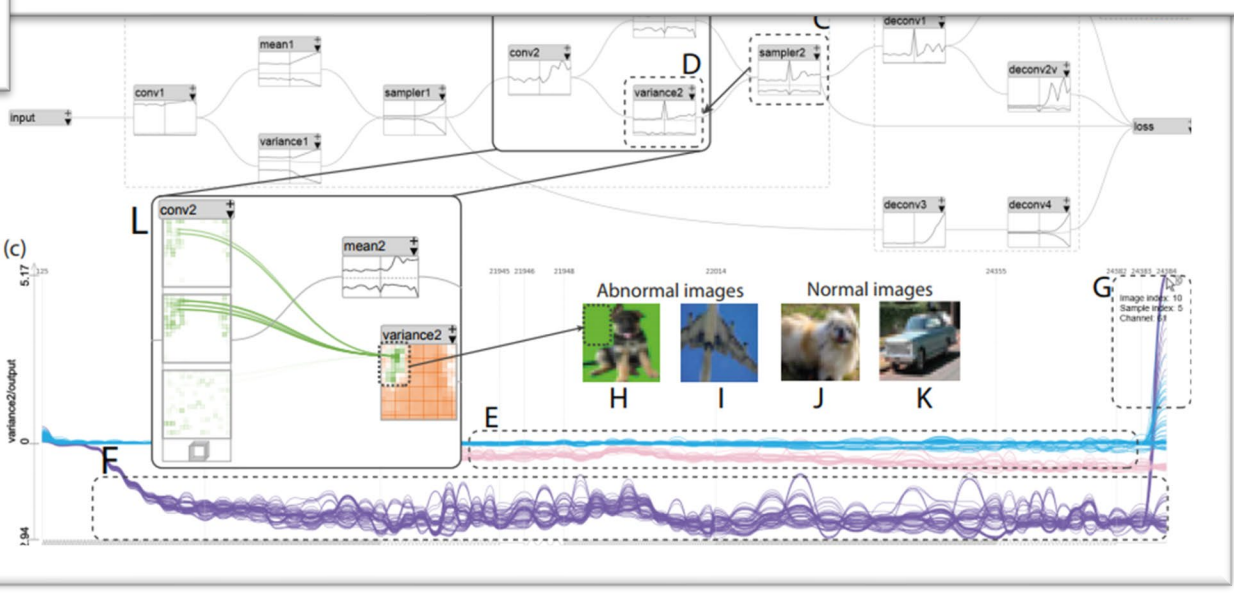
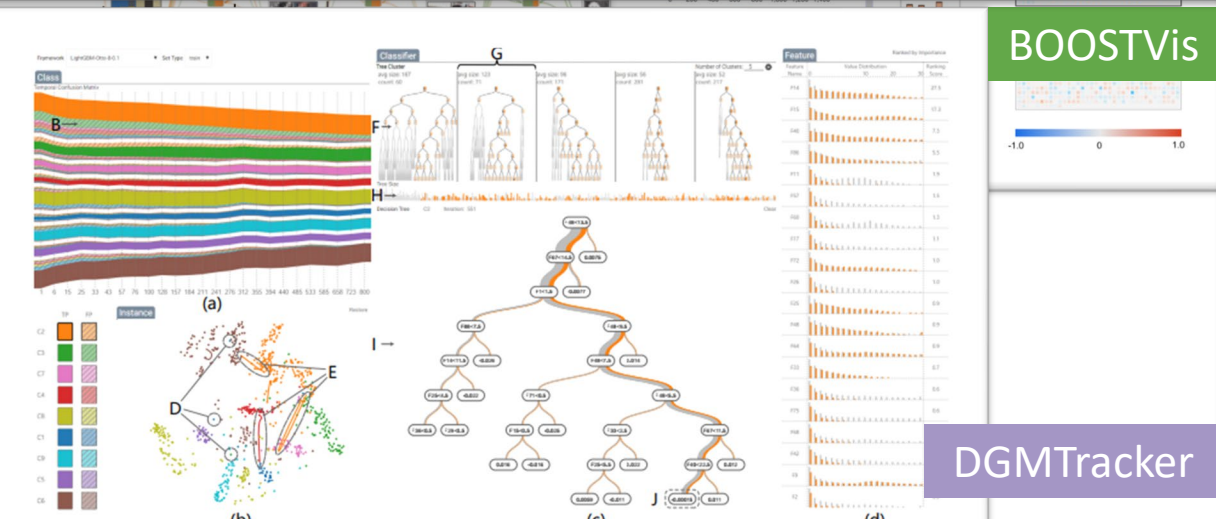
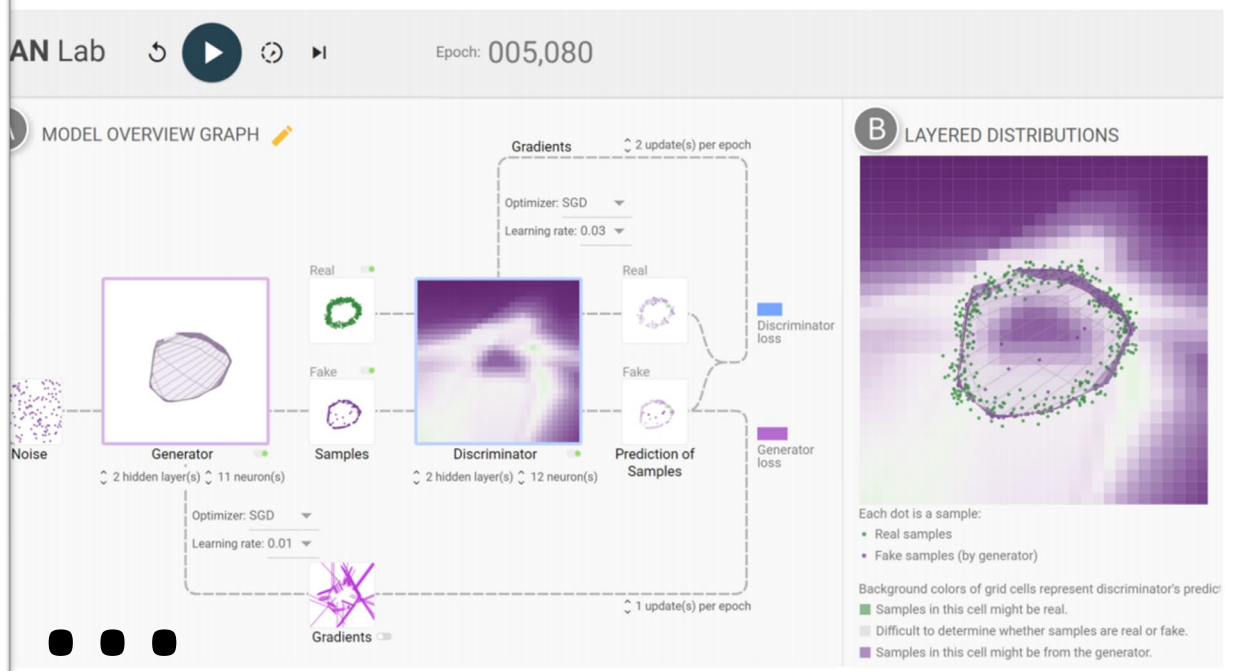
Steep learning curve



Related Work

DQNViz

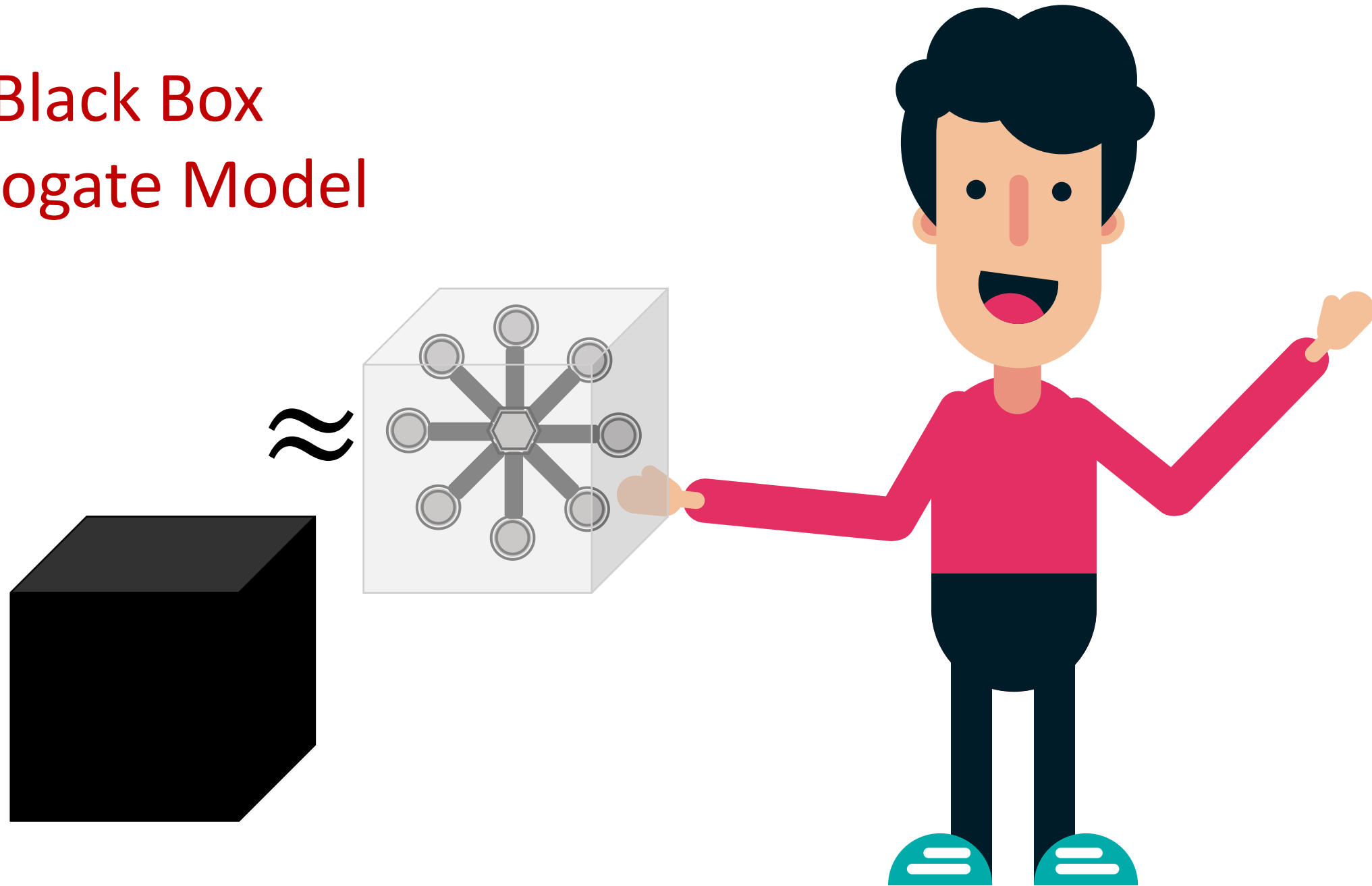
GANLab



Strategy 2

Close the Black Box

Using Surrogate Model

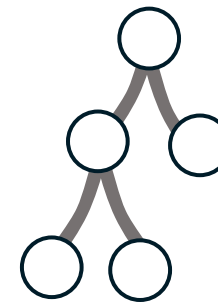


Strategy 2

Close the Black Box

Using Surrogate Model

Such as:

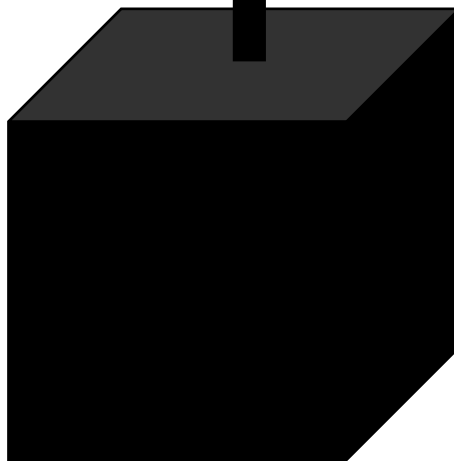
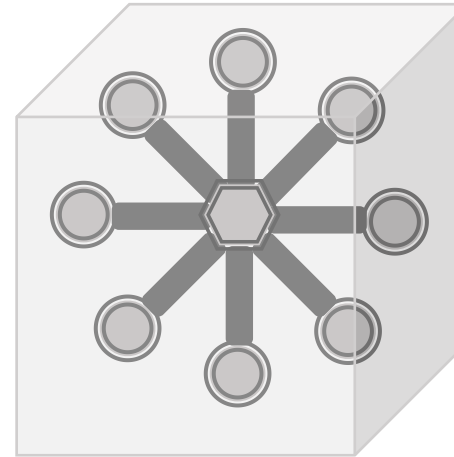


IF $A < B$ Then ...
Else IF $C > B$
Then ...

Decision Tree, Rule List ...

Same Input

Neural Networks, Random Forest...



Similar Output

Related Work

RuleMatrix

A Controls

Model Info:

type: rule-explainer
#rules: 53
model: wine_quality_red-nn-40-40-40-40-40-40-40

Dataset: wine_quality_red

train test sample train
sample test

Styles

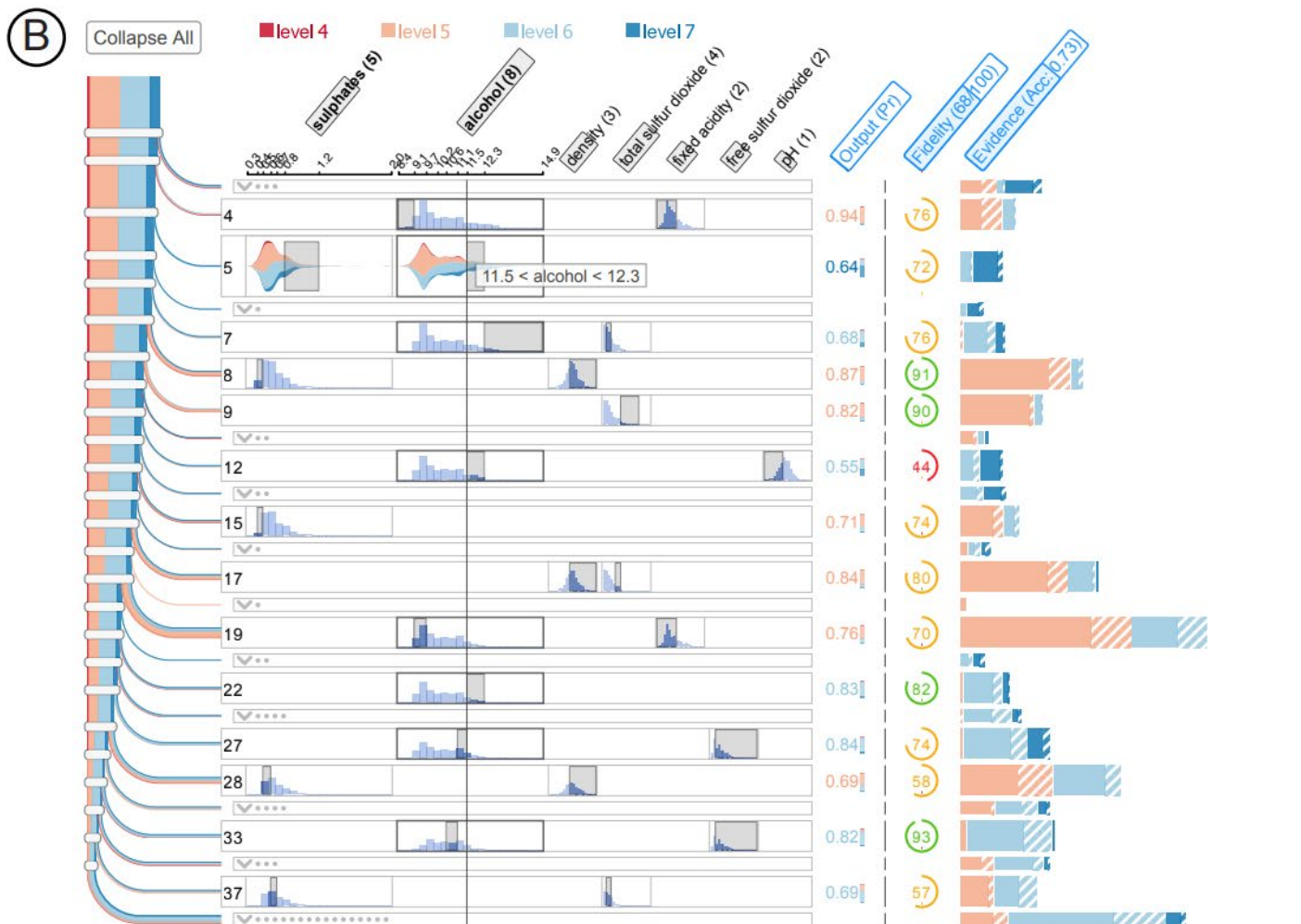
Flow Width:
Rect Width:
Rect Height:
Color Scheme: Seq Div Qual

Settings

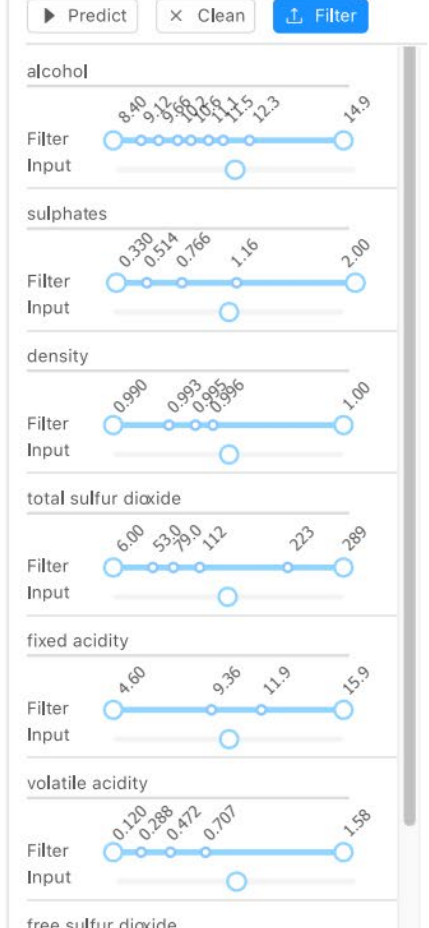
Conditional:
Detail Output:

Rule Filters

Min Evidence:
Fidelity:



C Data Filter

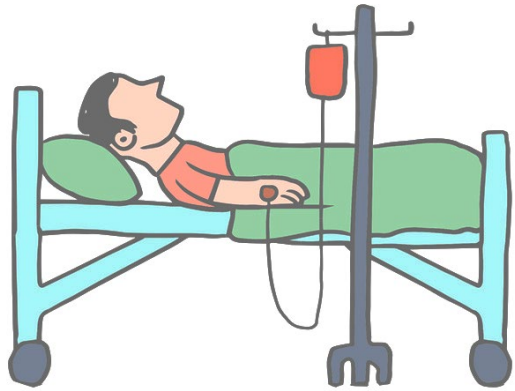


D

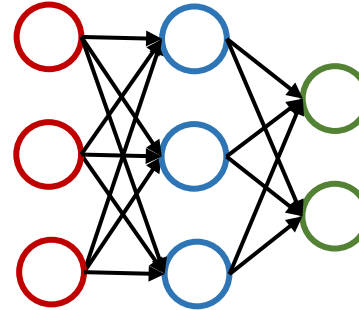
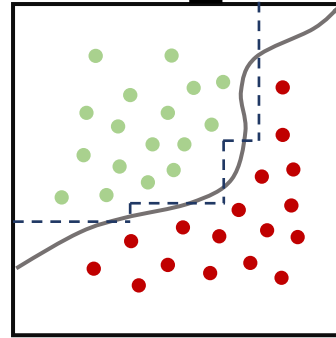
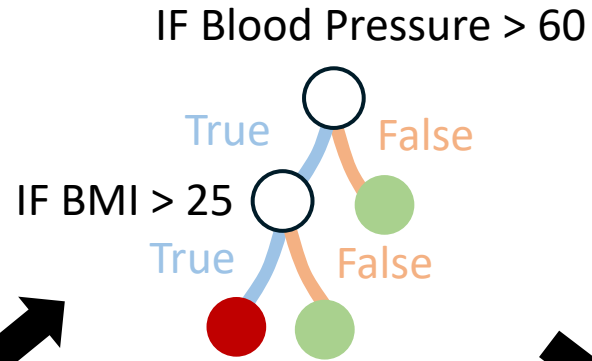
Data Table: train | (1199/1199)

Label	alcohol	sulphates	density	total sulfur dioxide	fixed acidity	volatile acidity	free sulfur dioxide	citric acid	pH	chlorides	residual sugar
level 5	9.500	0.5500	0.9971	22.00	9.300	0.4300	9.000	0.4400	3.280	0.08500	1.900

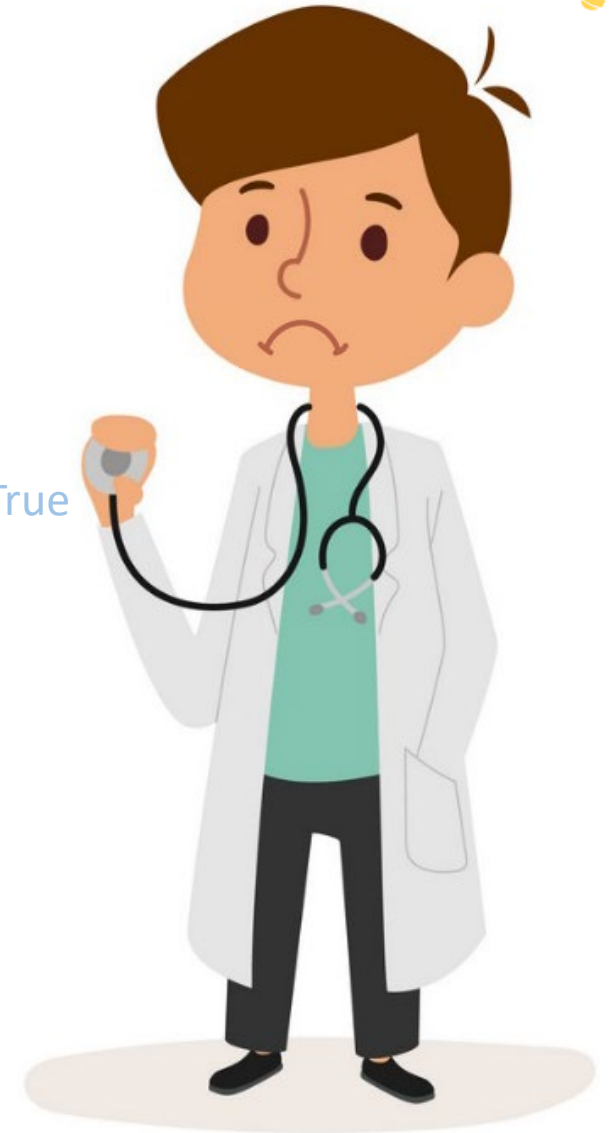
For example



Blood Pressure = 72
Skin Thickness = 35
BMI = 33.6

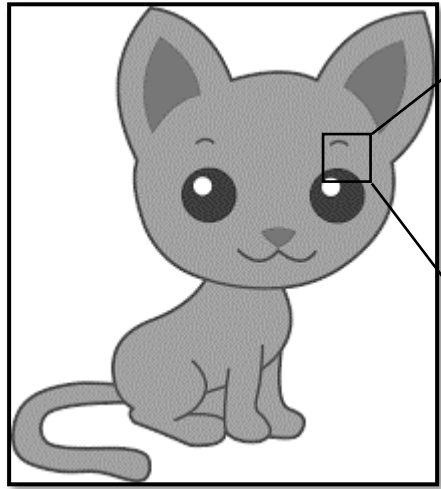


Diabetes? True

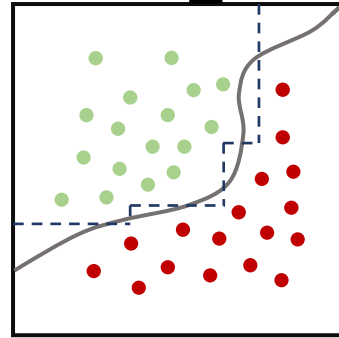
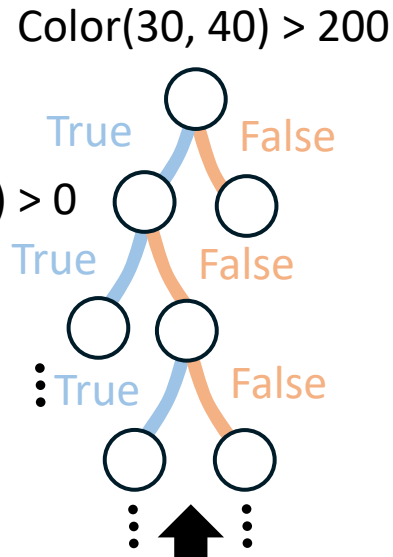
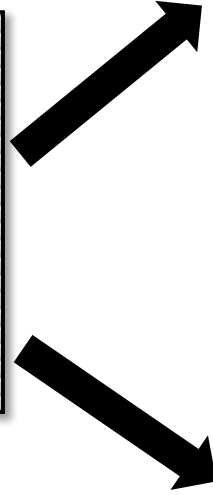
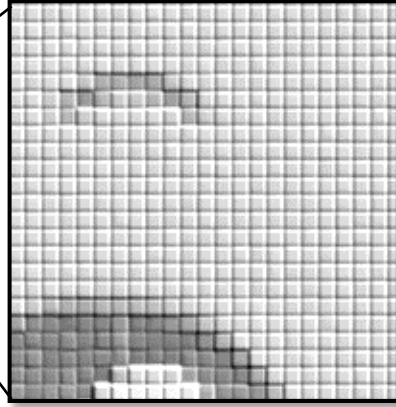


However

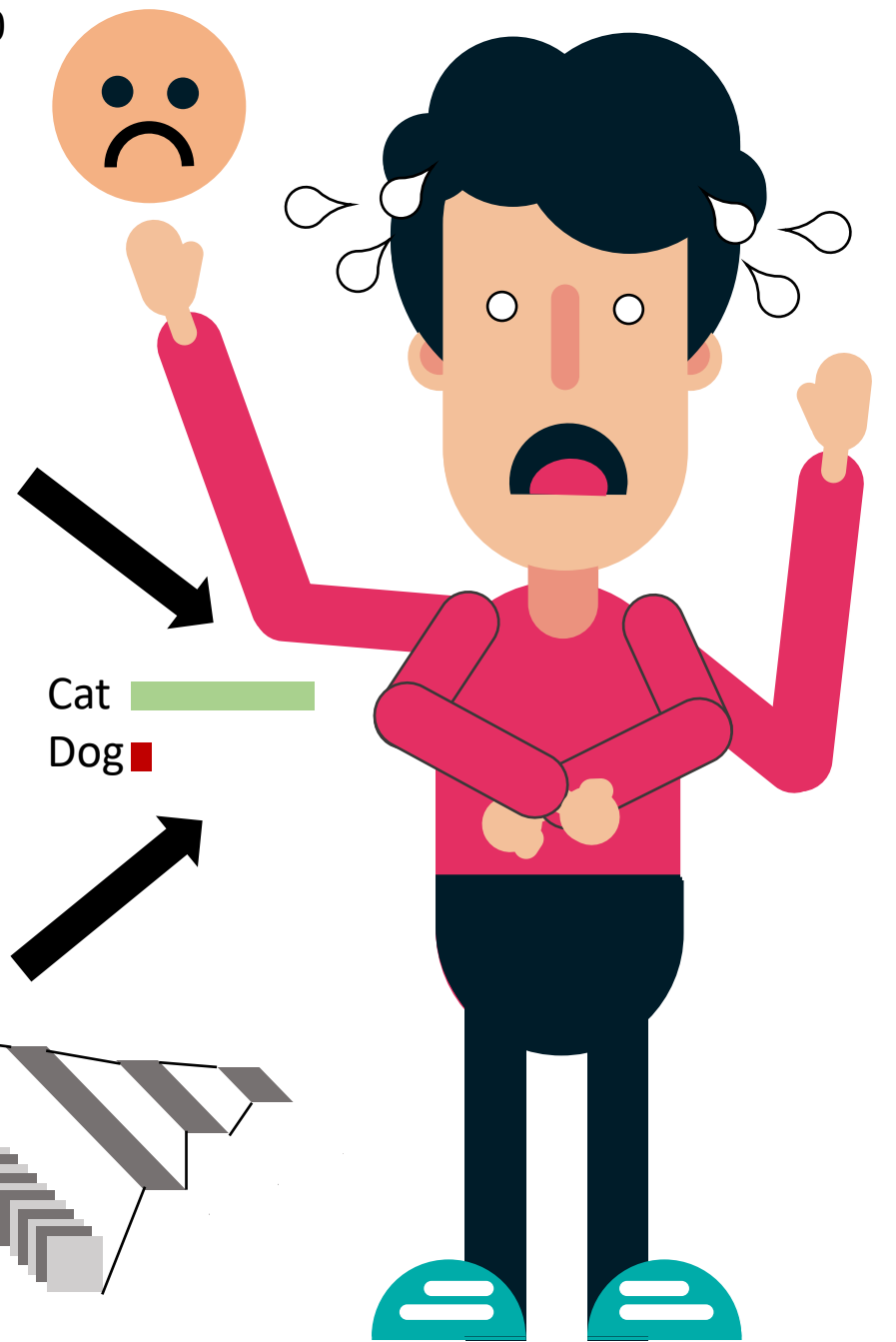
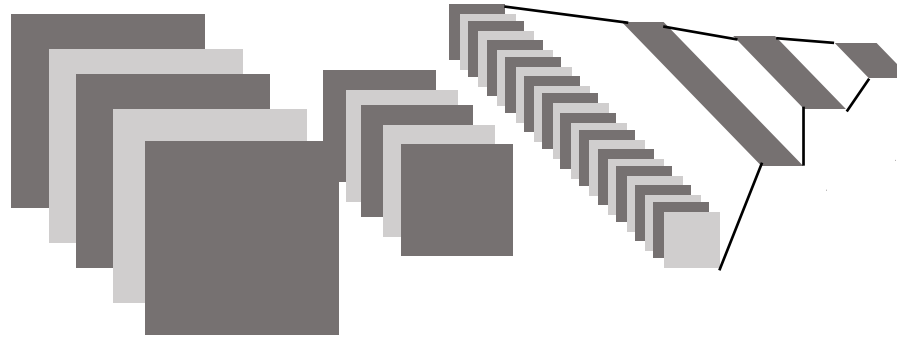
- Low Level Input
- Too Large Tree



image

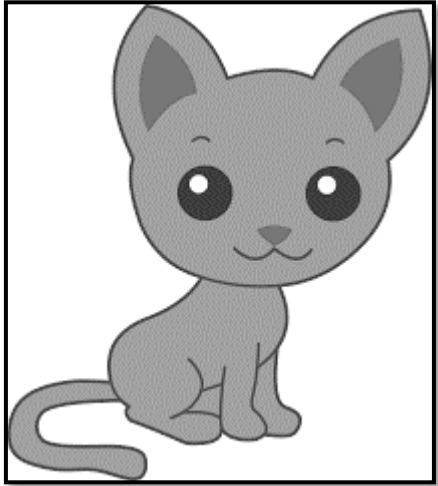


Cat █
Dog █

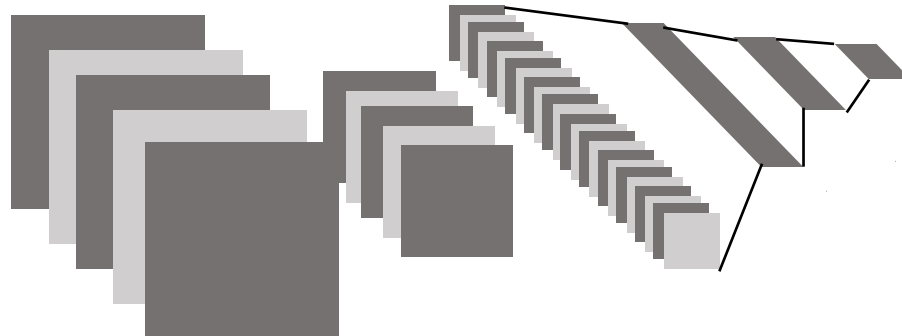


Our Strategy:

- 1 Slice the Black Box

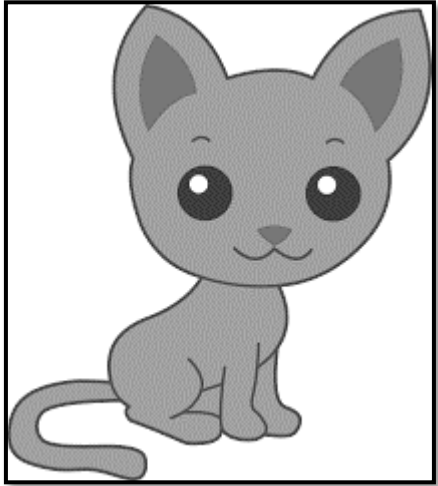


image

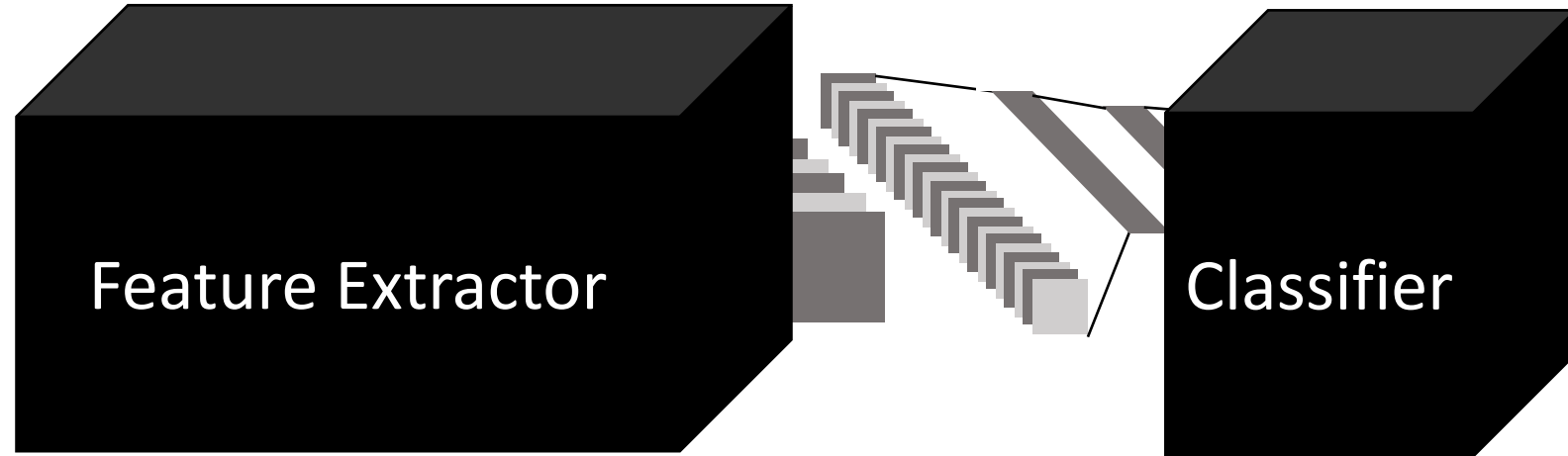


Our Strategy:

1 Slice the Black Box

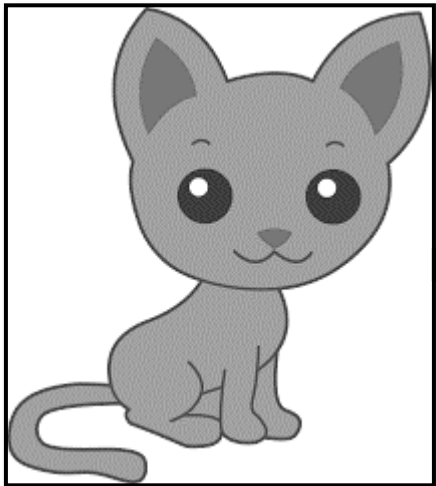


image

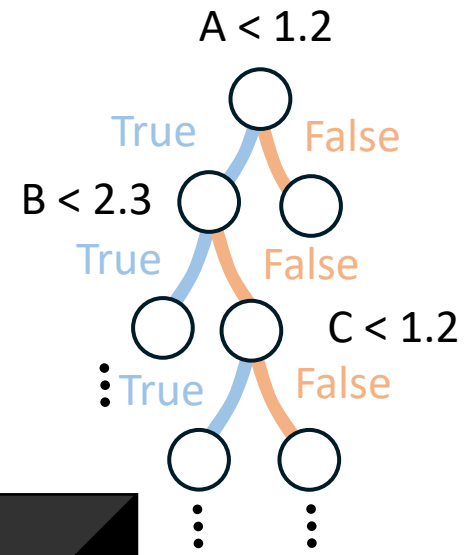
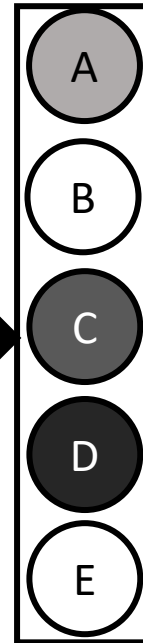
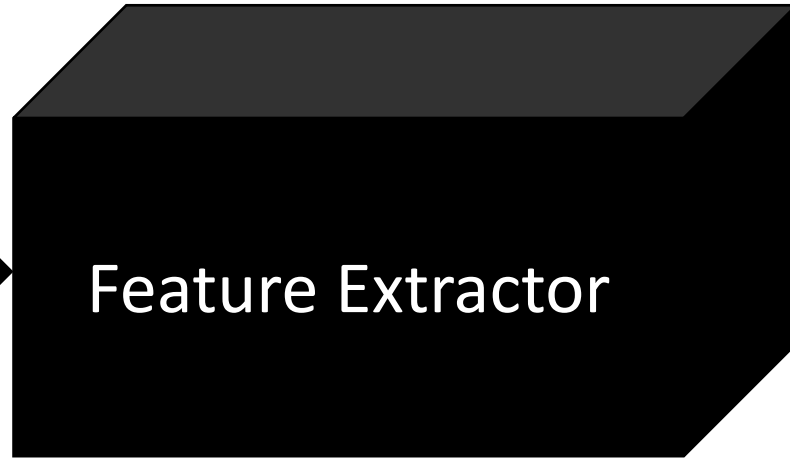



Our Strategy:

2 Extract the Decision Tree from Classifier



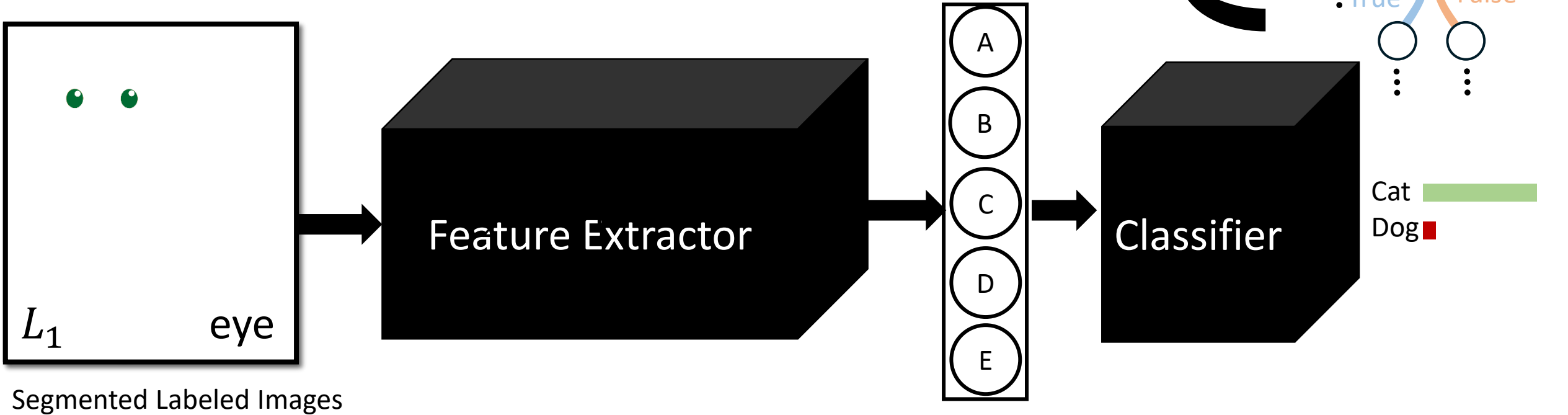
image



Cat 
Dog 

Our Strategy:

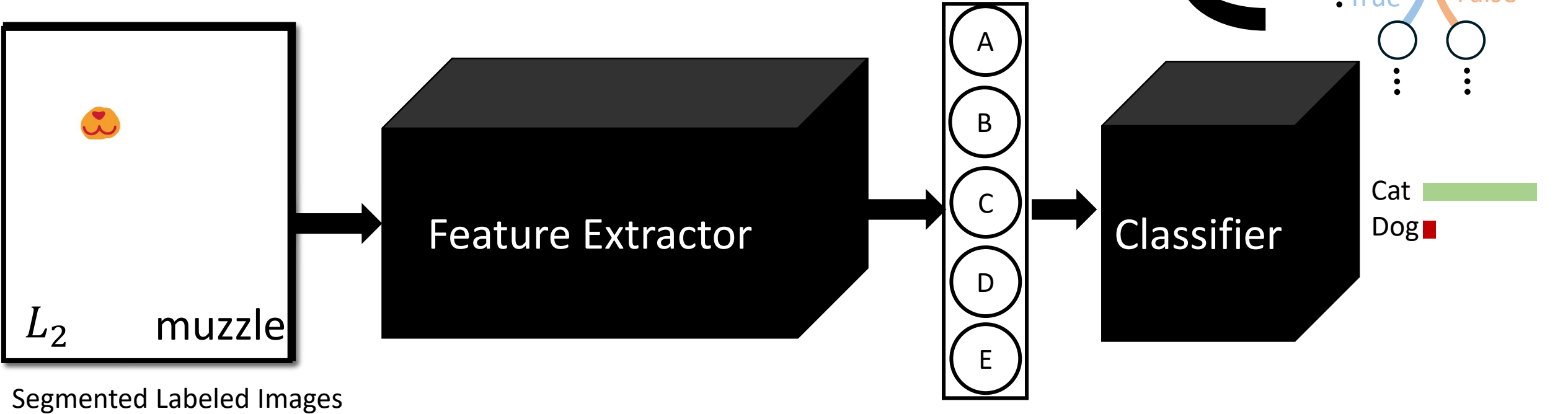
3 Match Features to Concepts¹



1. Network Dissection: Quantifying Interpretability of Deep Visual Representations (CVPR 2017), David Bau, Bolei Zhou, et.al

Our Strategy:

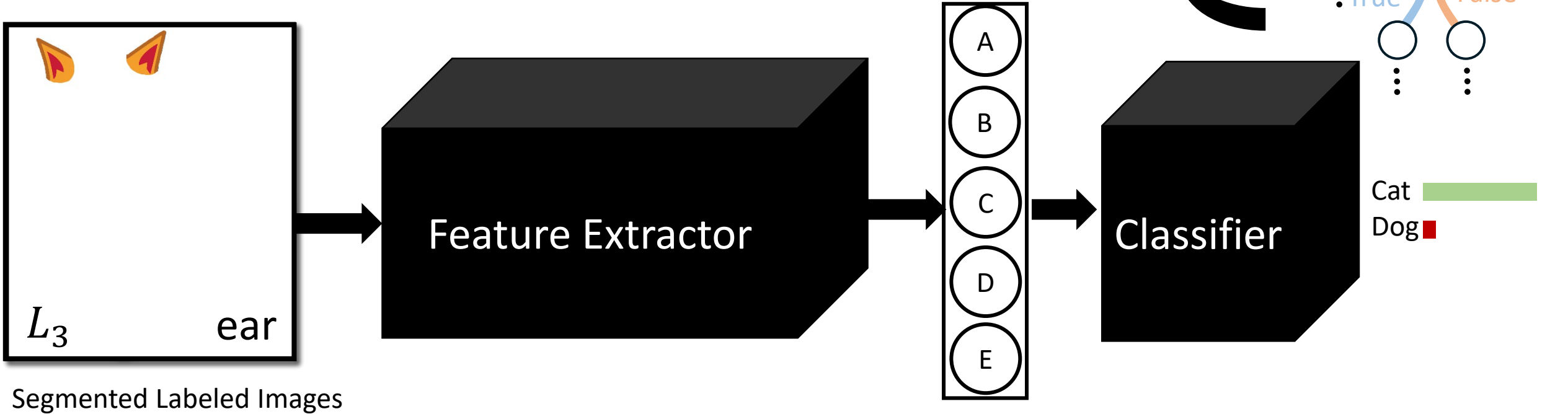
3 Match Features to Concepts¹



1. Network Dissection: Quantifying Interpretability of Deep Visual Representations (CVPR 2017), David Bau, Bolei Zhou, et.al

Our Strategy:

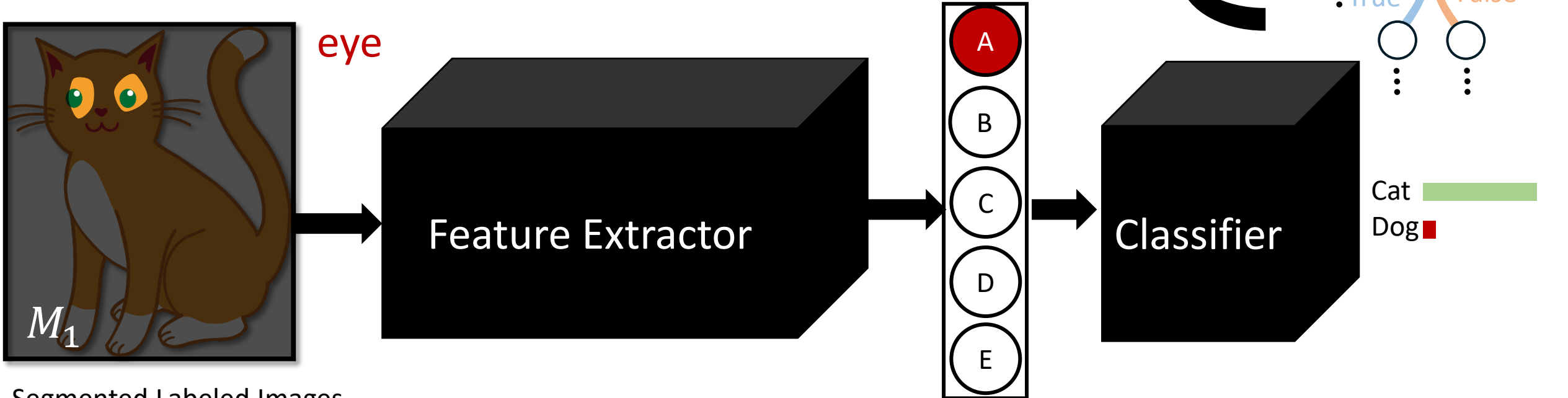
3 Match Features to Concepts¹



1. Network Dissection: Quantifying Interpretability of Deep Visual Representations (CVPR 2017), David Bau, Bolei Zhou, et.al

Our Strategy:

3 Match Features to Concepts¹



$$IoU(M_i, L_j) = \frac{M_i \cap L_j}{M_i \cup L_j}$$

Area of Overlap

Area of Union



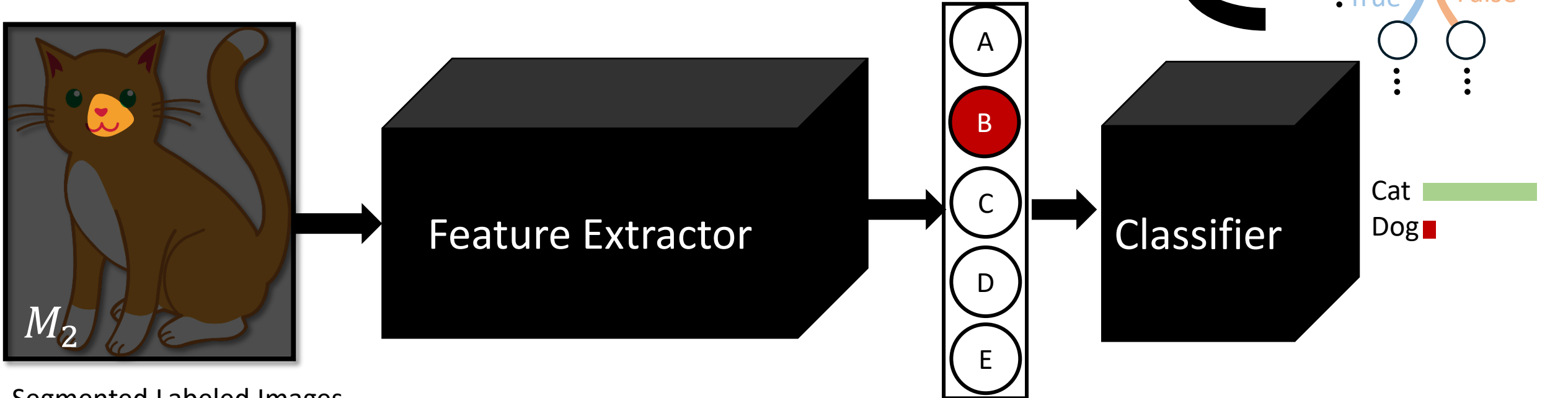
$IoU(M_1, L_1) > IoU(M_1, L_2) > IoU(M_1, L_3)$

eye > muzzle > ear !

1. Network Dissection: Quantifying Interpretability of Deep Visual Representations (CVPR 2017), David Bau, Bolei Zhou, et.al

Our Strategy:

3 Match Features to Concepts¹



Segmented Labeled Images

$$IoU(M_i, L_j) = \frac{M_i \cap L_j}{M_i \cup L_j}$$

Area of Overlap
Area of Union

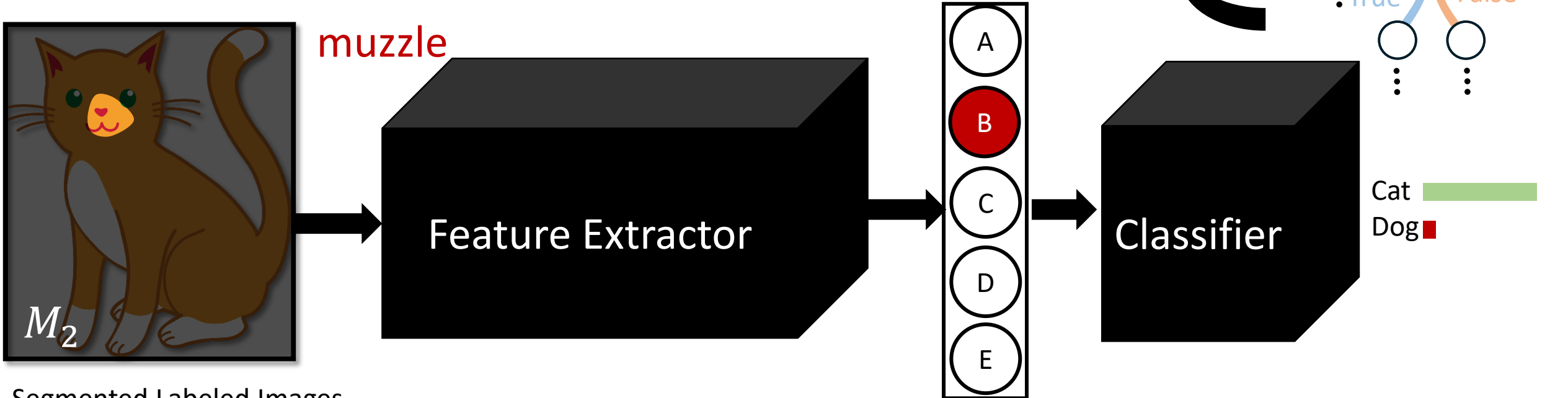


$IoU(M_2, L_2) > IoU(M_2, L_1) > IoU(M_2, L_3)$
muzzle > eye > ear !

1. Network Dissection: Quantifying Interpretability of Deep Visual Representations (CVPR 2017), David Bau, Bolei Zhou, et.al

Our Strategy:

3 Match Features to Concepts¹



$$IoU(M_i, L_j) = \frac{M_i \cap L_j}{M_i \cup L_j}$$

Area of Overlap
Area of Union

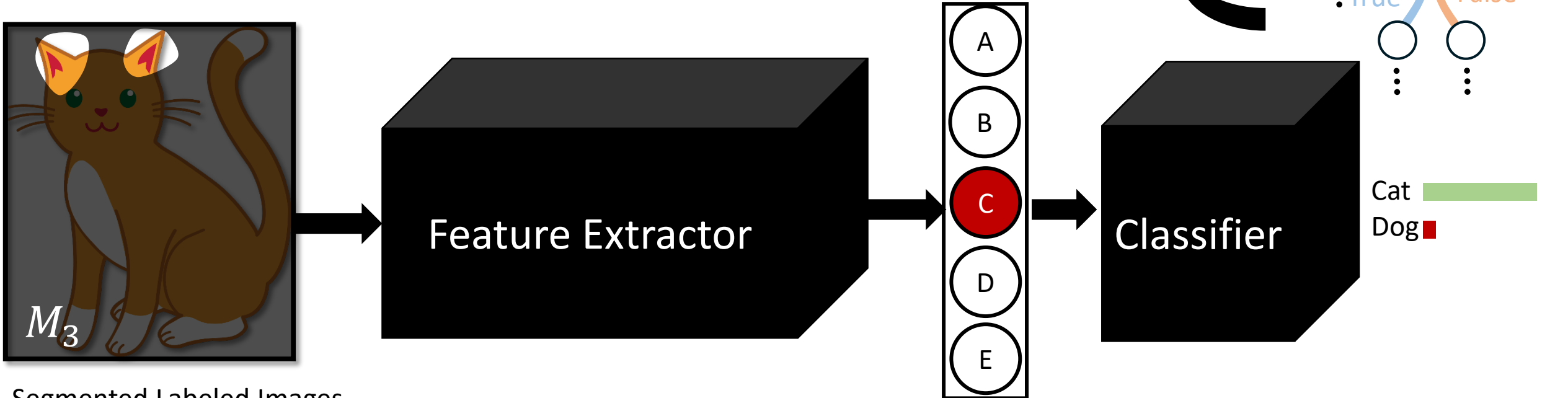


$IoU(M_2, L_2) > IoU(M_2, L_1) > IoU(M_2, L_3)$
muzzle > eye > ear !

1. Network Dissection: Quantifying Interpretability of Deep Visual Representations (CVPR 2017), David Bau, Bolei Zhou, et.al

Our Strategy:

3 Match Features to Concepts¹



Segmented Labeled Images

$$IoU(M_i, L_j) = \frac{M_i \cap L_j}{M_i \cup L_j}$$

Area of Overlap
Area of Union

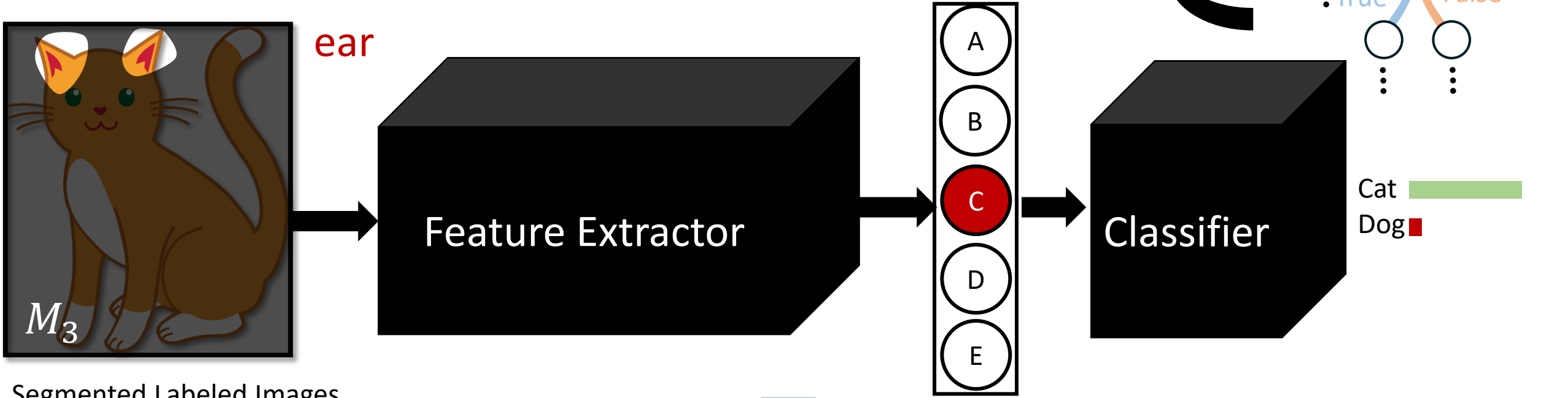


$IoU(M_3, L_3) > IoU(M_3, L_1) > IoU(M_3, L_2)$
ear > eye > muzzle !

1. Network Dissection: Quantifying Interpretability of Deep Visual Representations (CVPR 2017), David Bau, Bolei Zhou, et.al

Our Strategy:

3 Match Features to Concepts¹



Segmented Labeled Images

$$IoU(M_i, L_j) = \frac{M_i \cap L_j}{M_i \cup L_j}$$

Area of Overlap

Area of Union



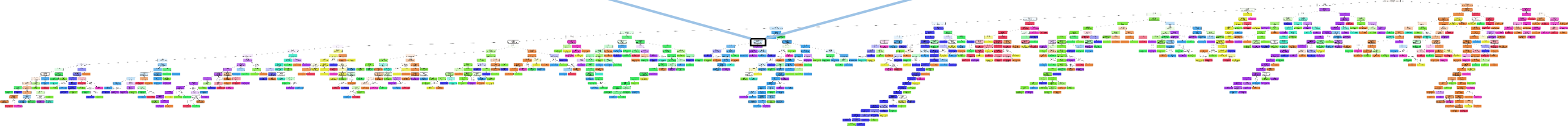
$IoU(M_3, L_3) > IoU(M_3, L_1) > IoU(M_3, L_2)$

ear > eye > muzzle !

1. Network Dissection: Quantifying Interpretability of Deep Visual Representations (CVPR 2017), David Bau, Bolei Zhou, et.al

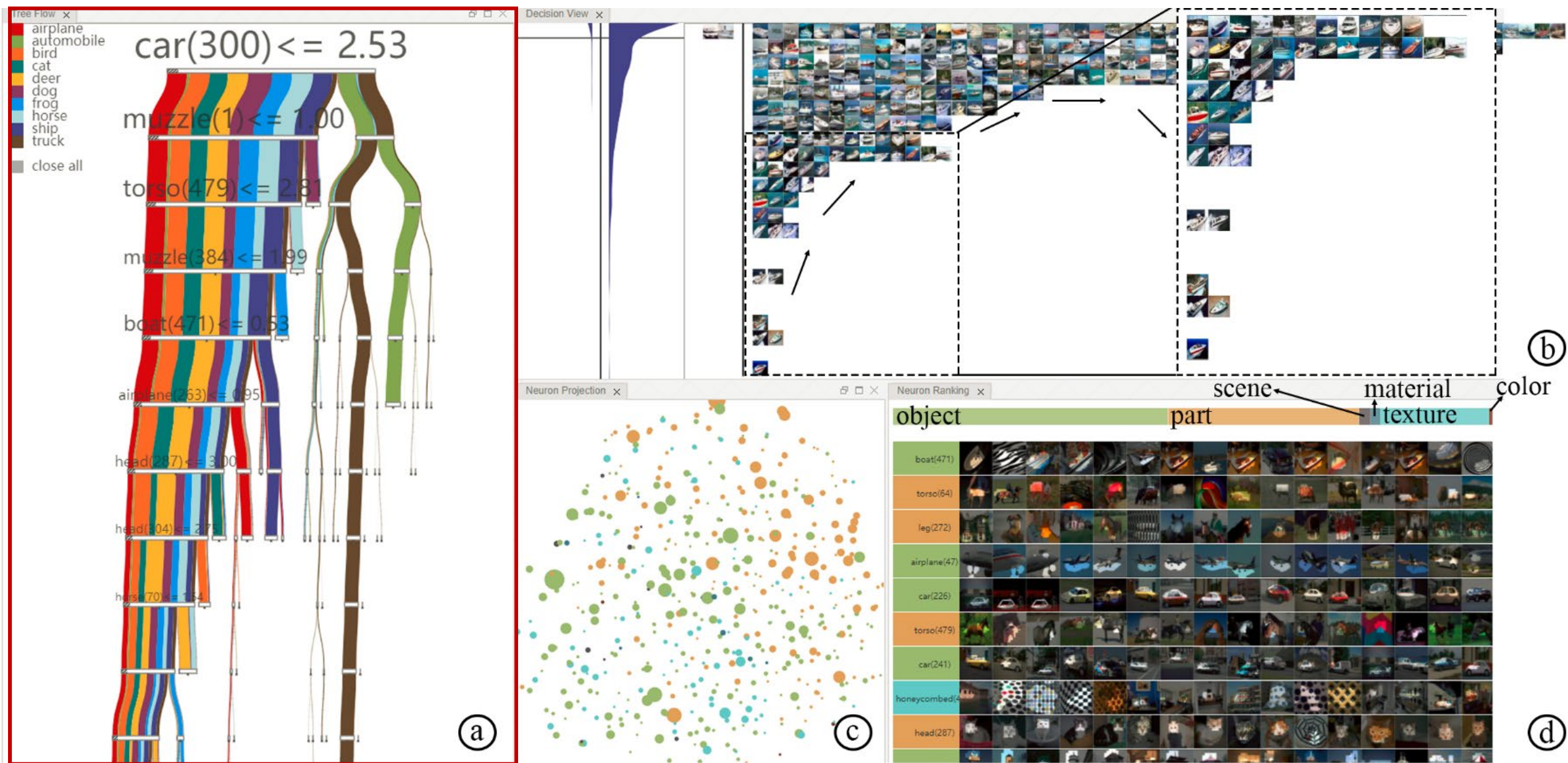
Challenges:

Neuron (28) \leq 46.5589
Samples = 71
Value = [3, 0, 1, 3, 1, 35, 4, 24, 0, 0]

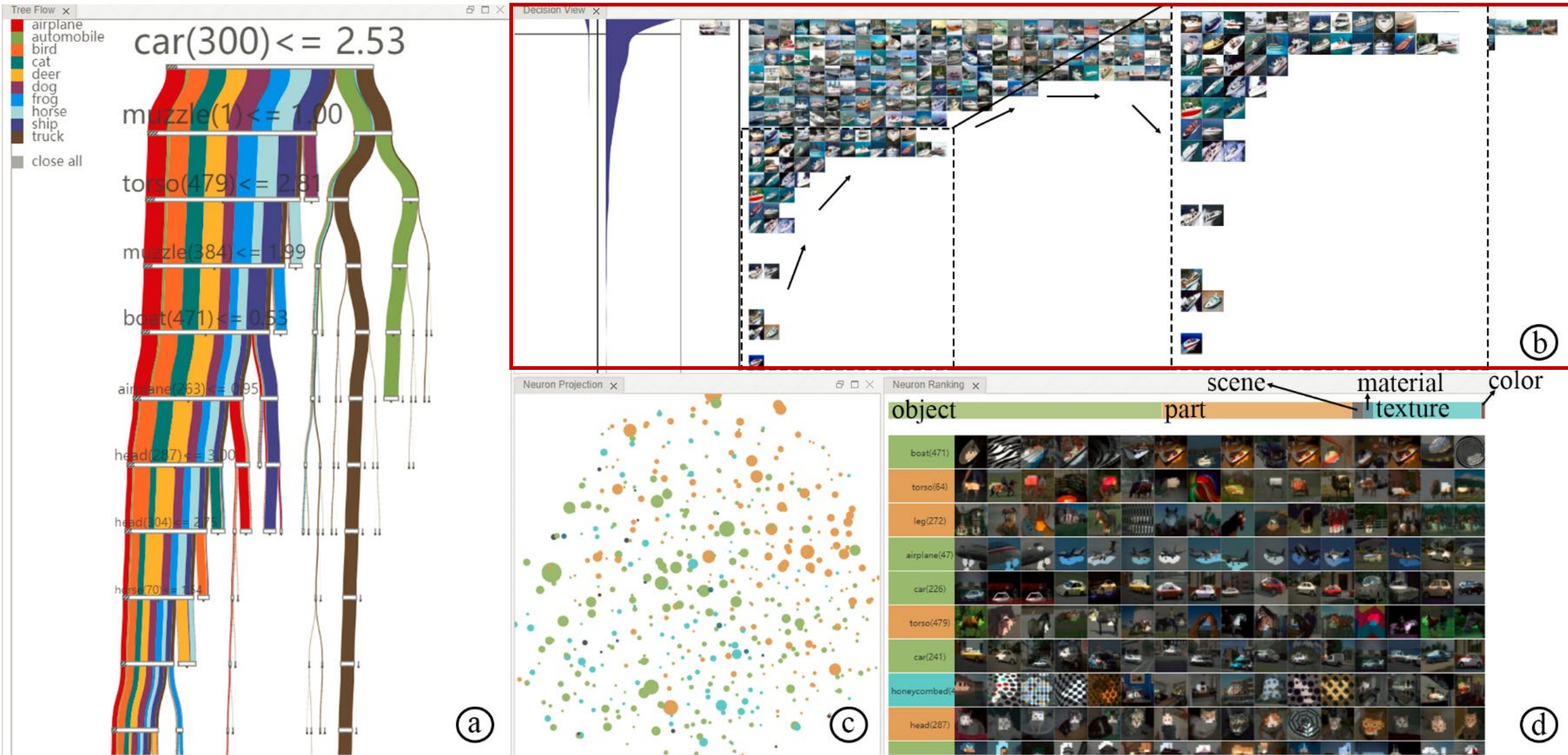


- 1 Scalability of large decision trees
- 2 Hard to grasp information of each rule
- 3 Hard to reason misclassification

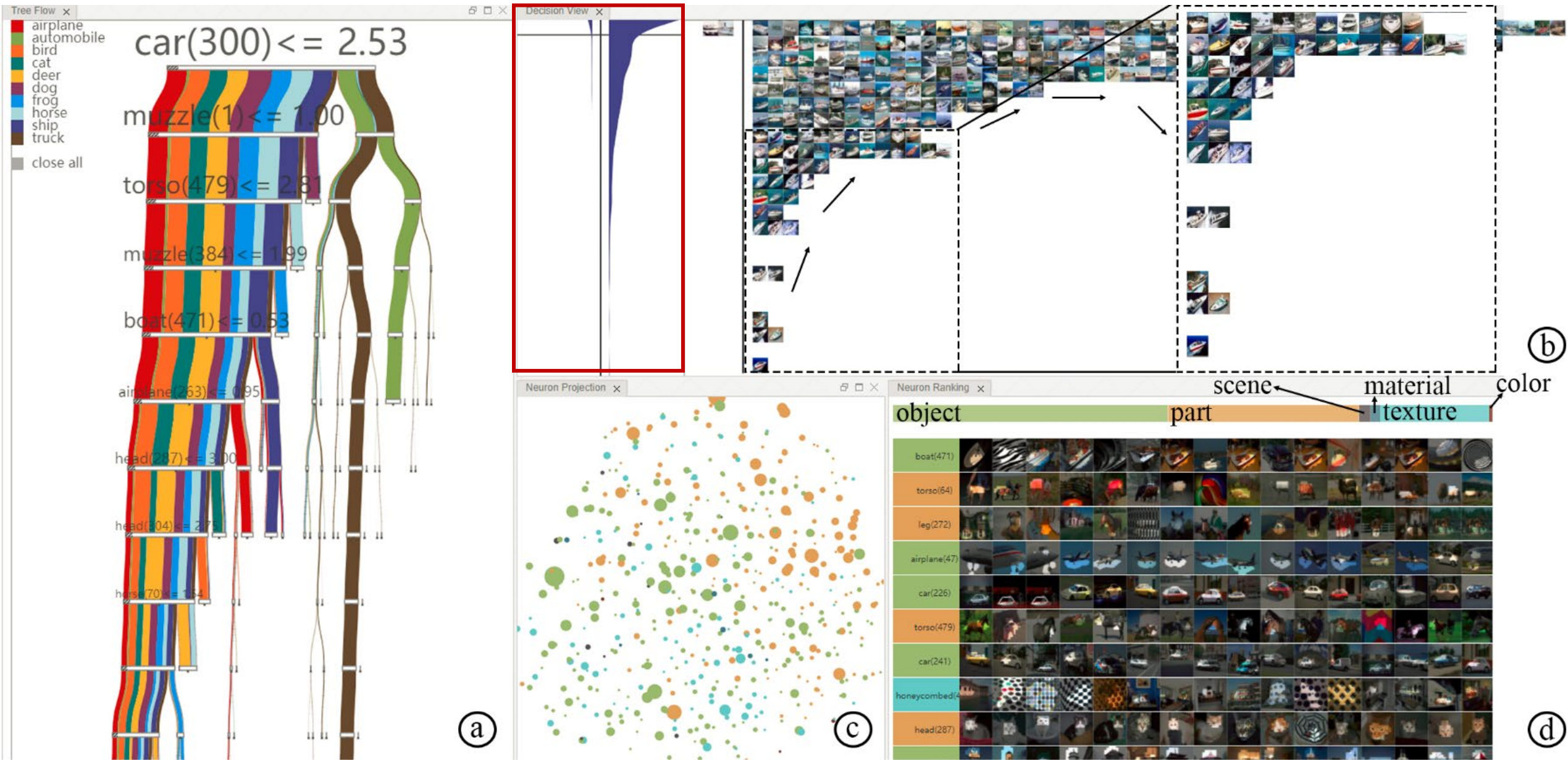
TreeFlow: dataflow along the decision tree



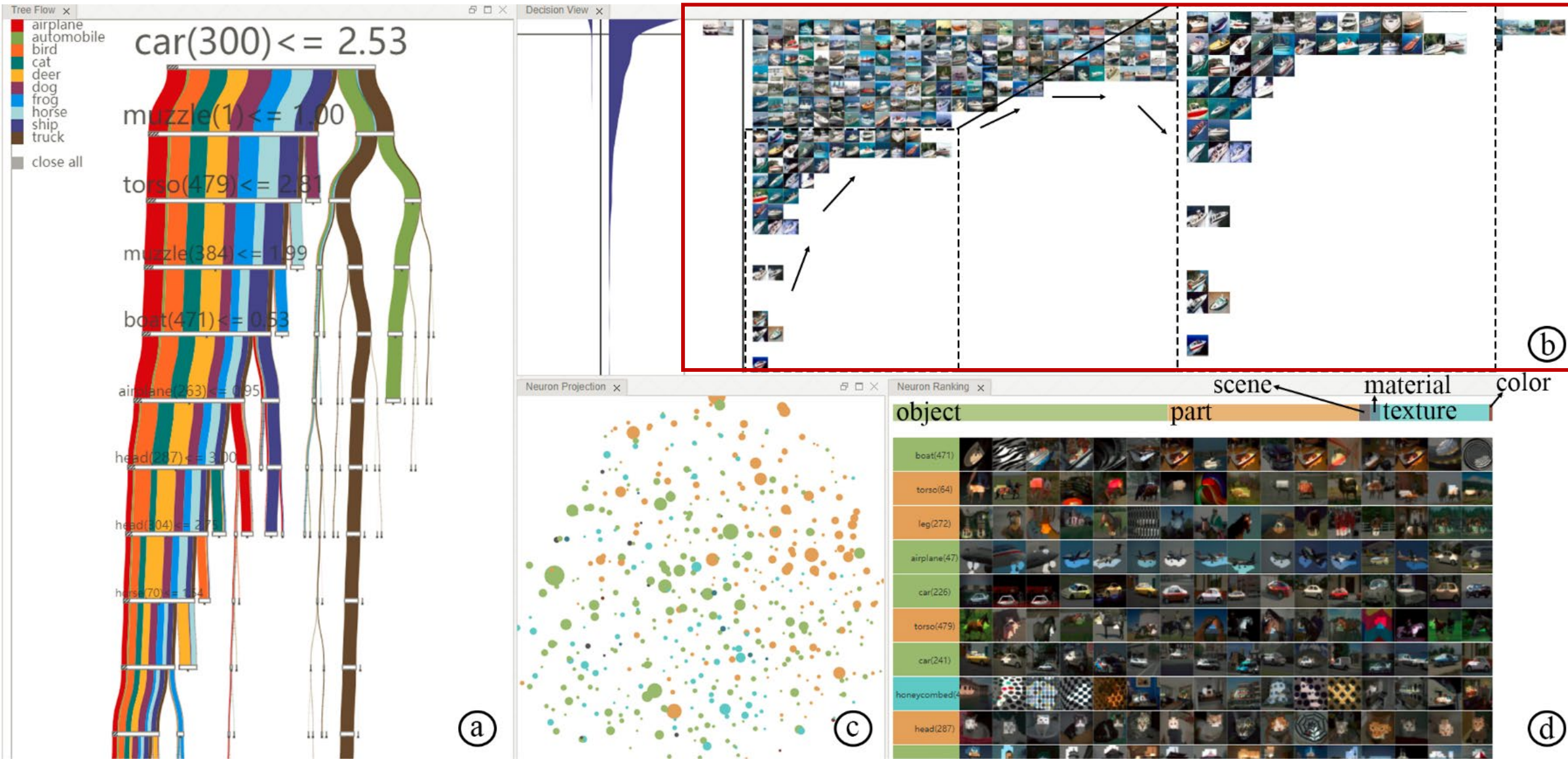
Decision View: data distribution of decision node



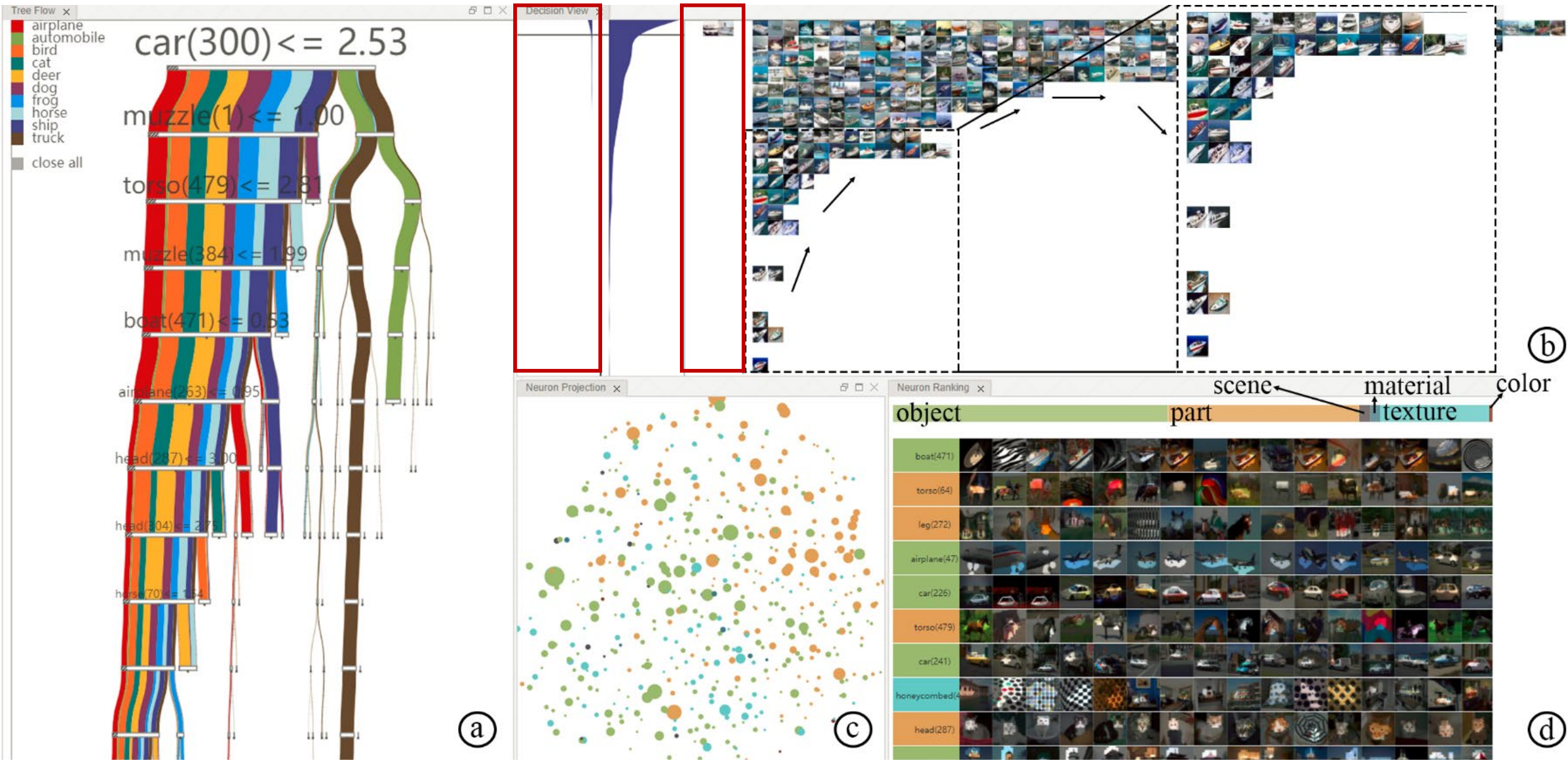
Decision View: data distribution of decision node



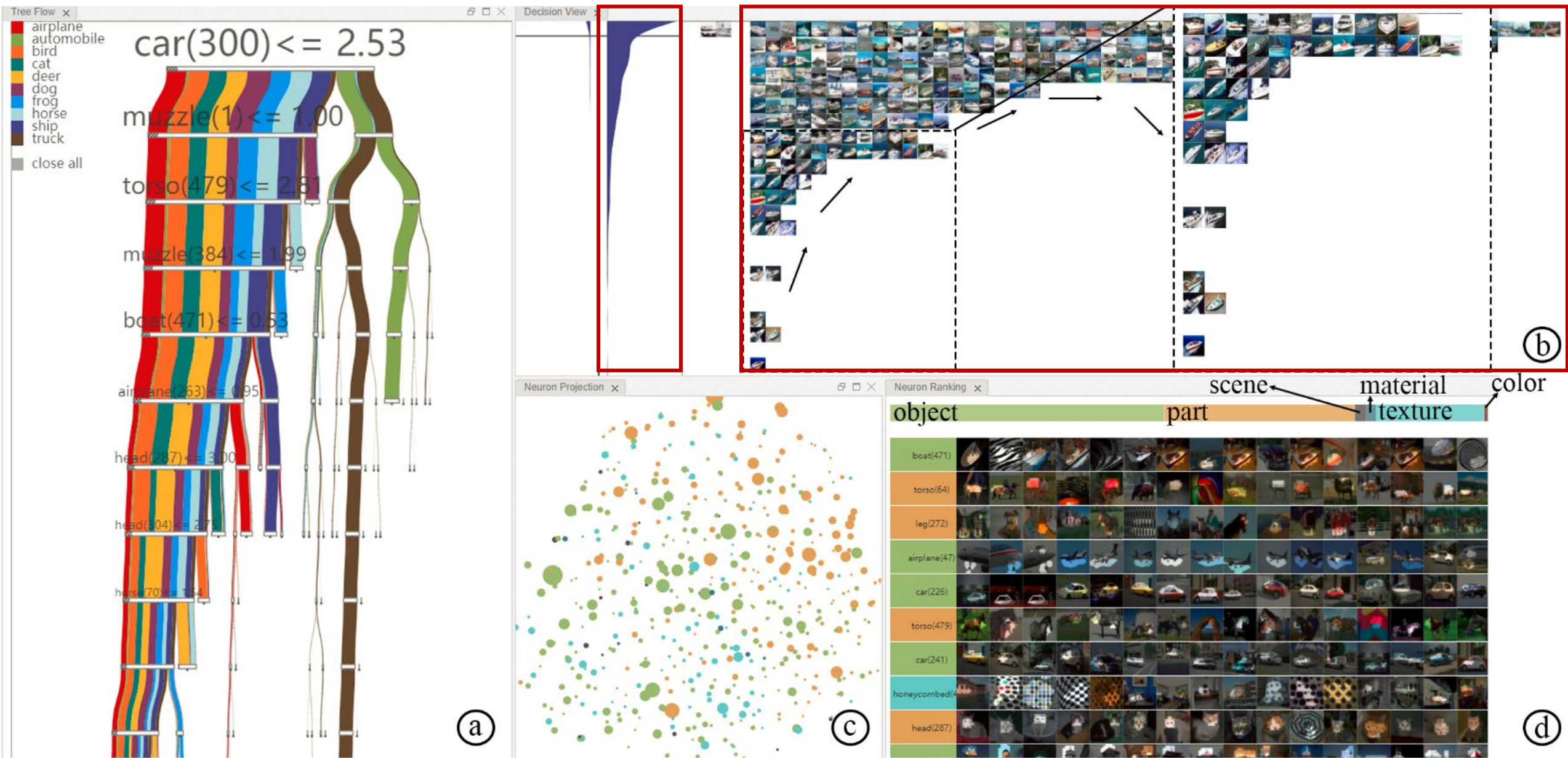
Decision View: data distribution of decision node



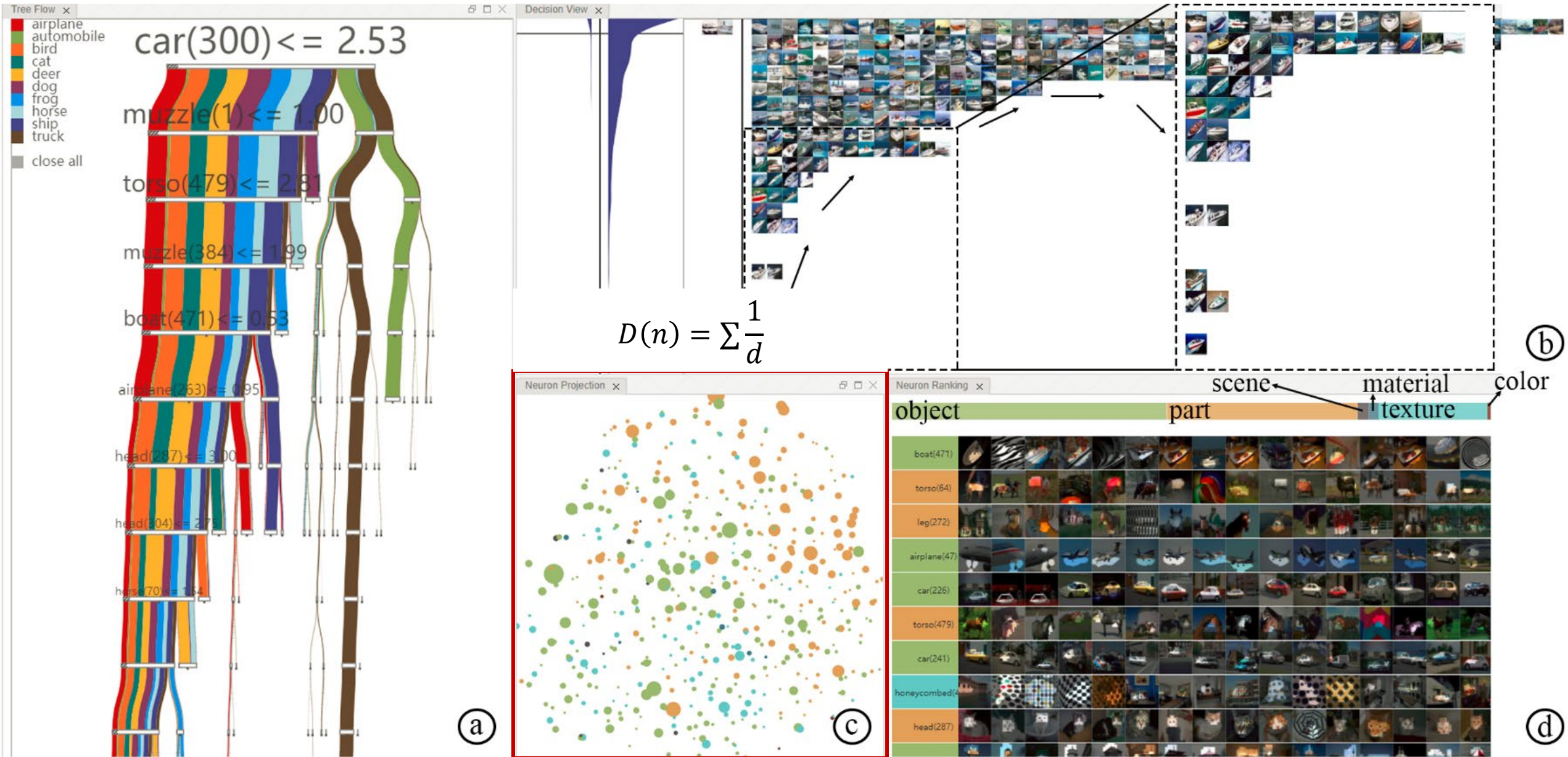
Decision View: data distribution of decision node



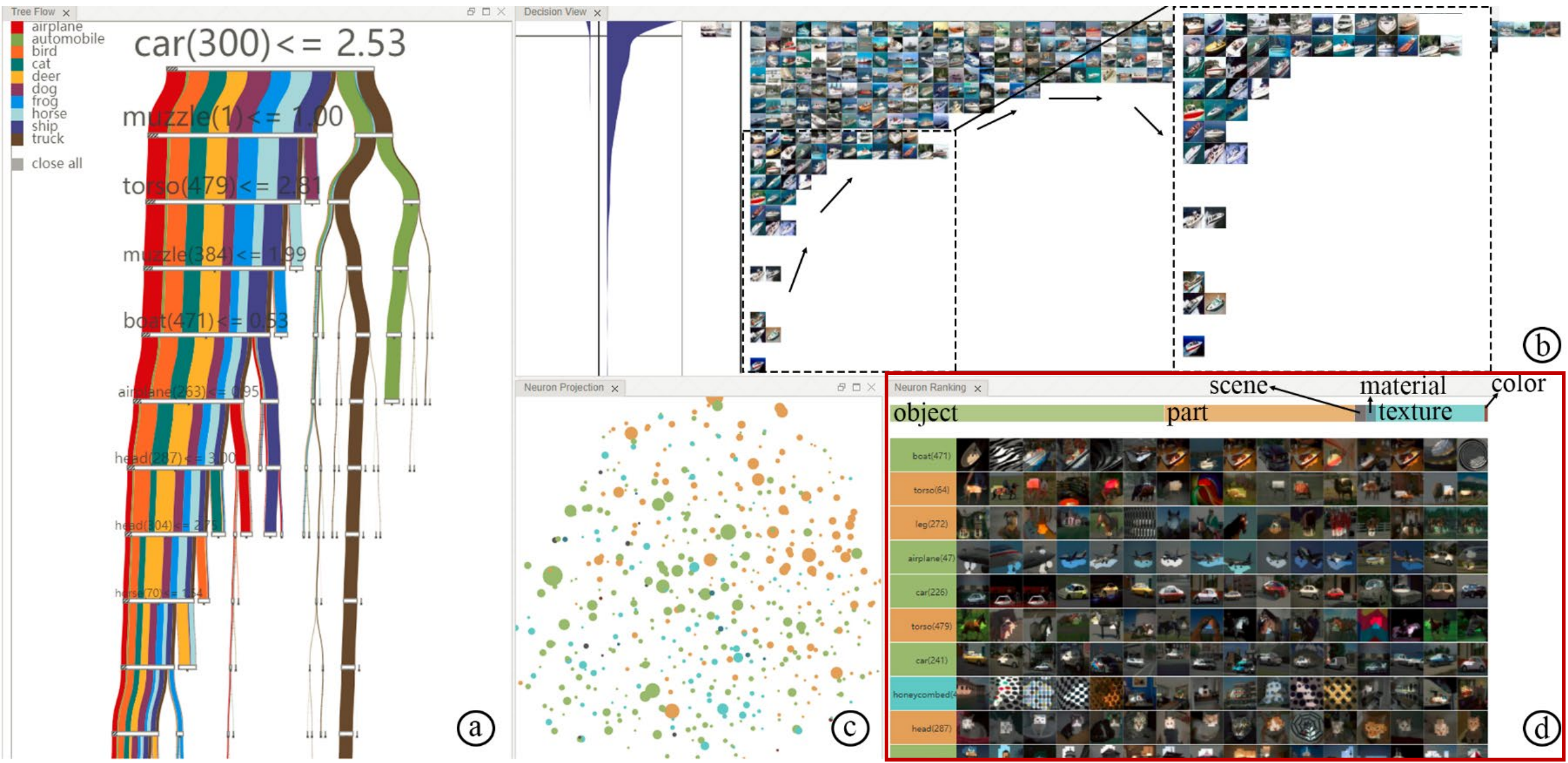
Decision View: data distribution of decision node



Neuron View: t-SNE projection of neurons based on semantics similarity



Neuron Ranking: a collection of semantic labels and saliency maps



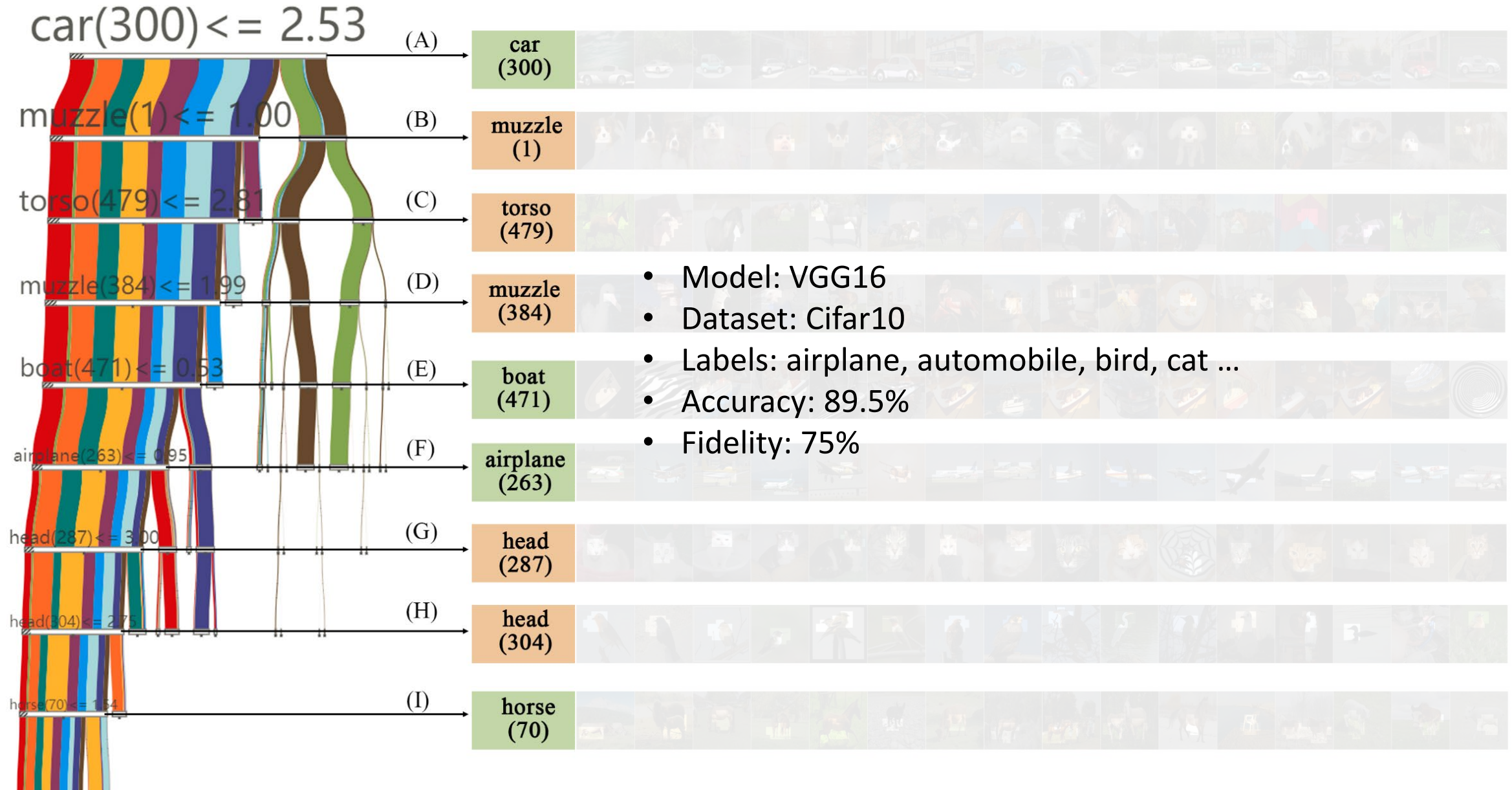
Basic Interactions

Use Case 1: Interpreting Surrogate Decision Tree

data class

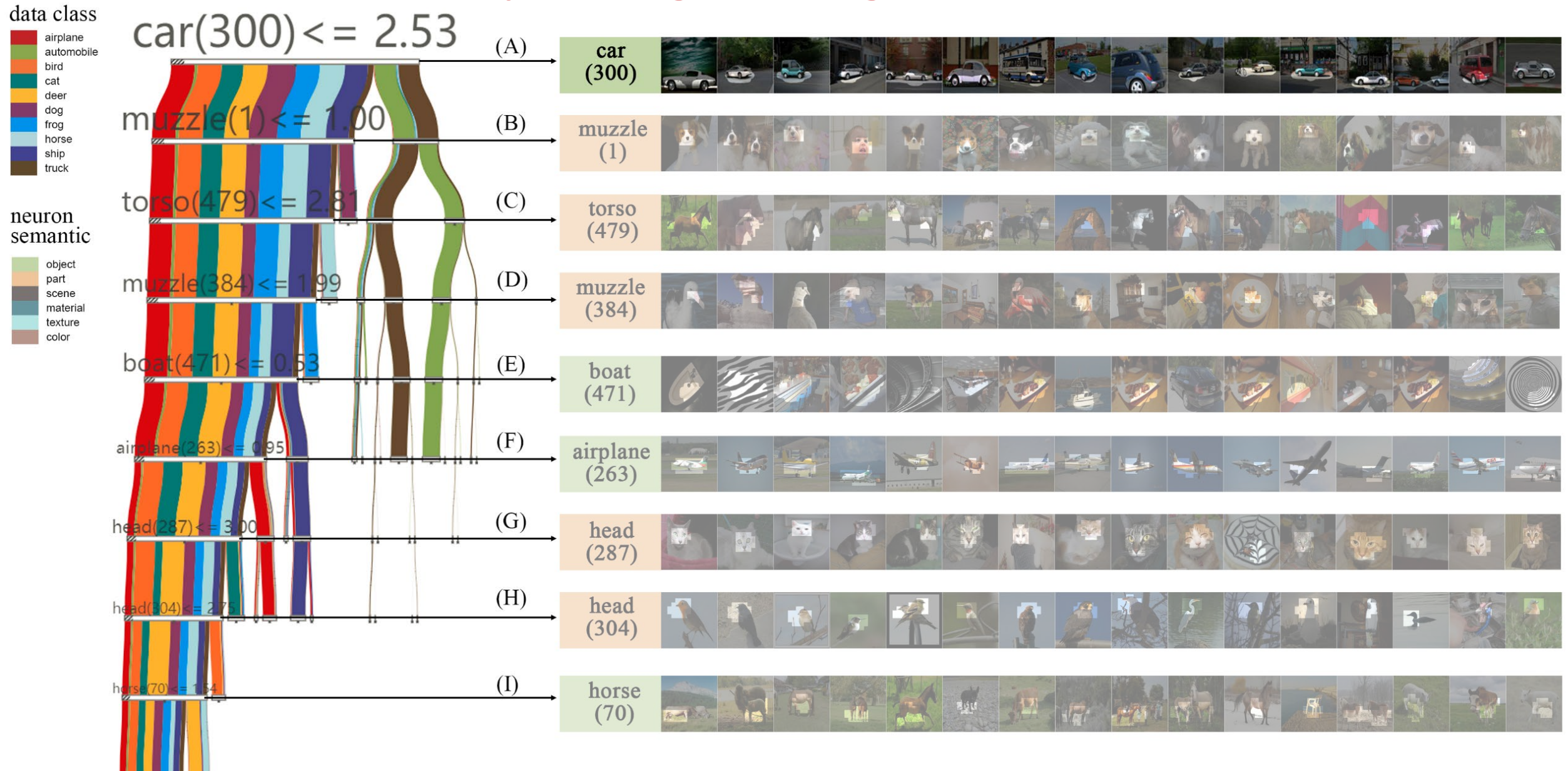


neuron semantic



- Model: VGG16
- Dataset: Cifar10
- Labels: airplane, automobile, bird, cat ...
- Accuracy: 89.5%
- Fidelity: 75%

Use Case 1: Interpreting Surrogate Decision Tree



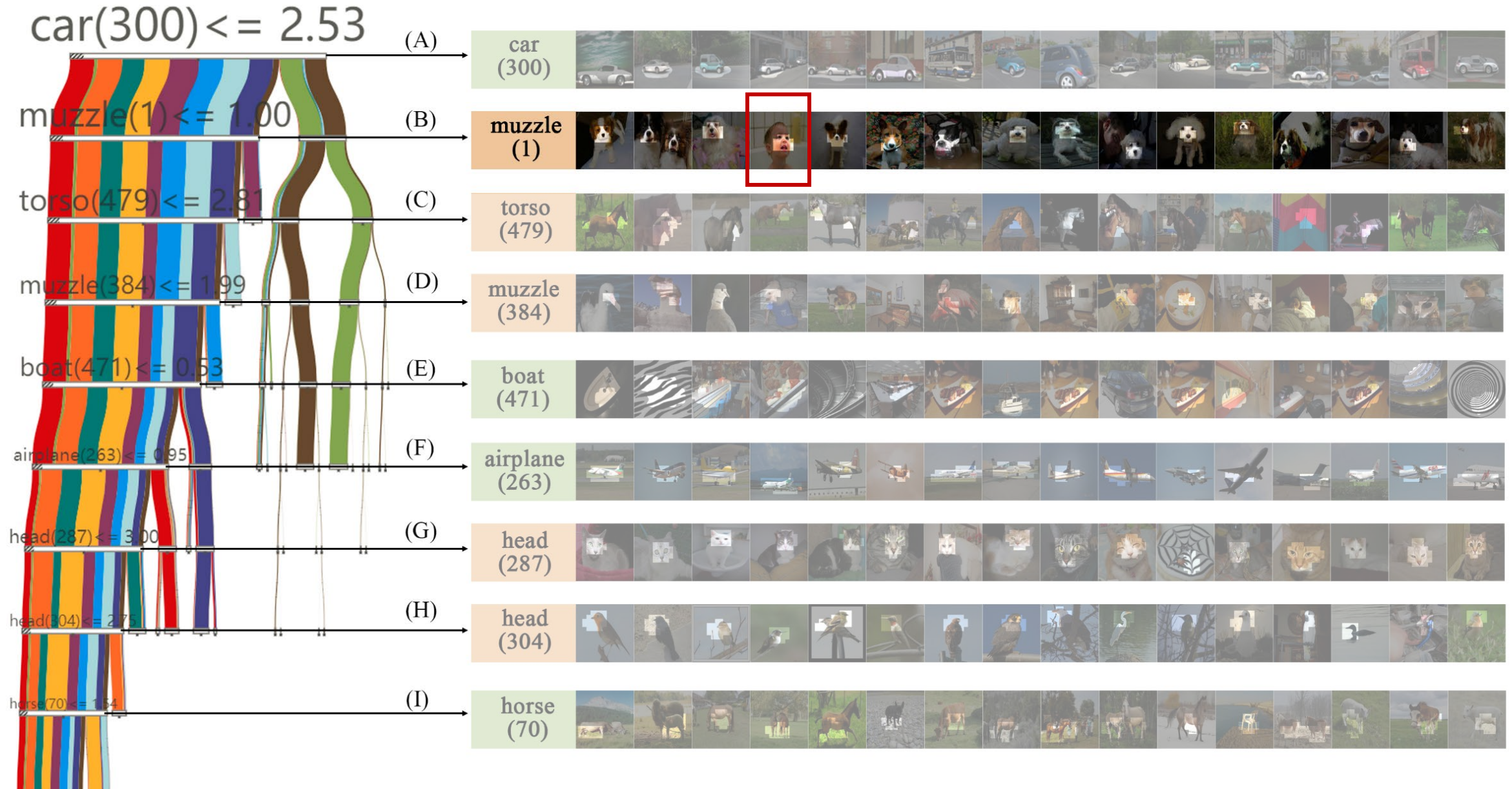
Use Case 1: Interpreting Surrogate Decision Tree

data class

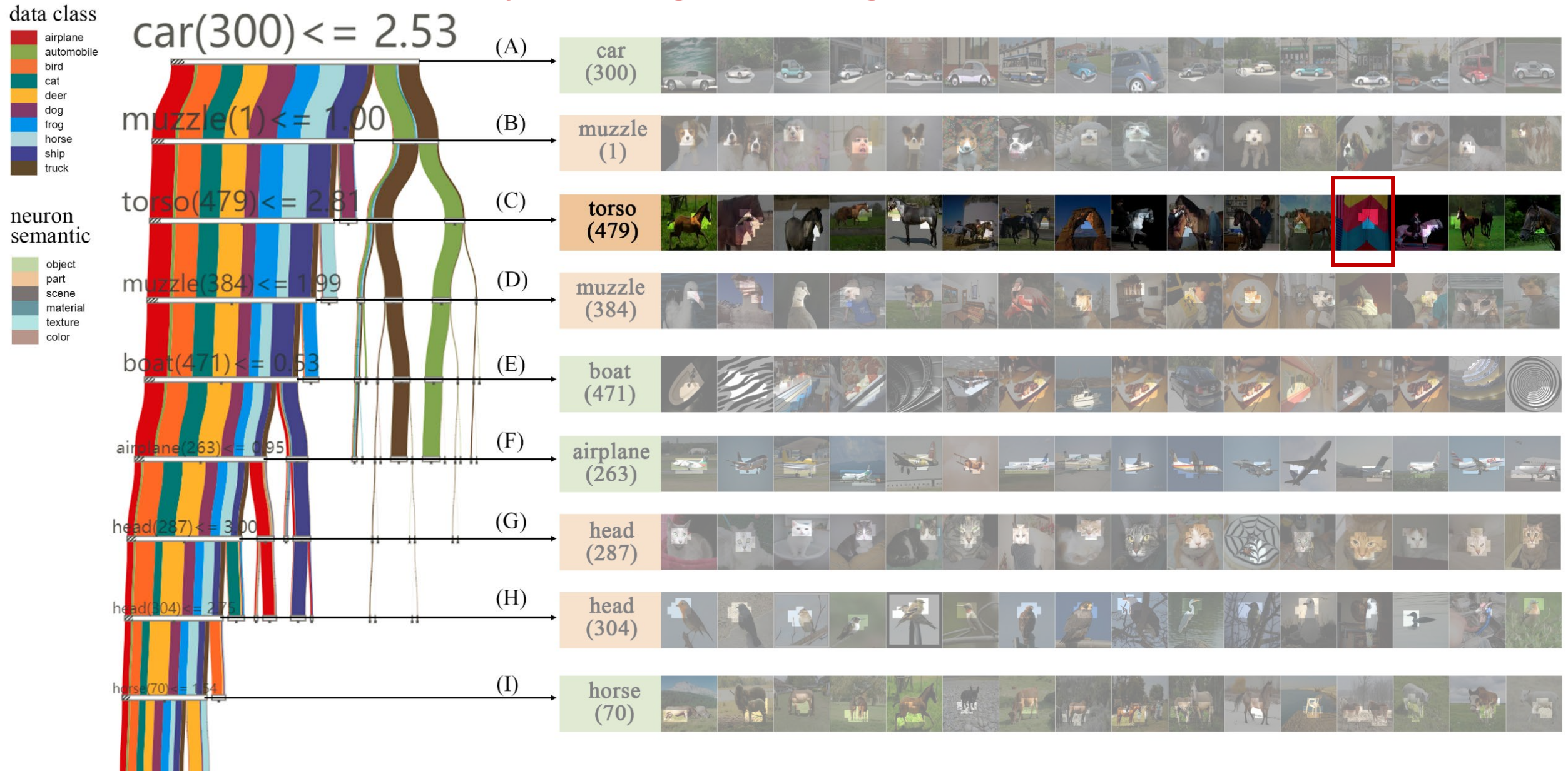
- airplane
- automobile
- bird
- cat
- deer
- dog
- frog
- horse
- ship
- truck

neuron semantic

- object
- part
- scene
- material
- texture
- color



Use Case 1: Interpreting Surrogate Decision Tree



Use Case 1: Interpreting Surrogate Decision Tree

data class

- airplane
- automobile
- bird
- cat
- deer
- dog
- frog
- horse
- ship
- truck

neuron semantic

- object
- part
- scene
- material
- texture
- color



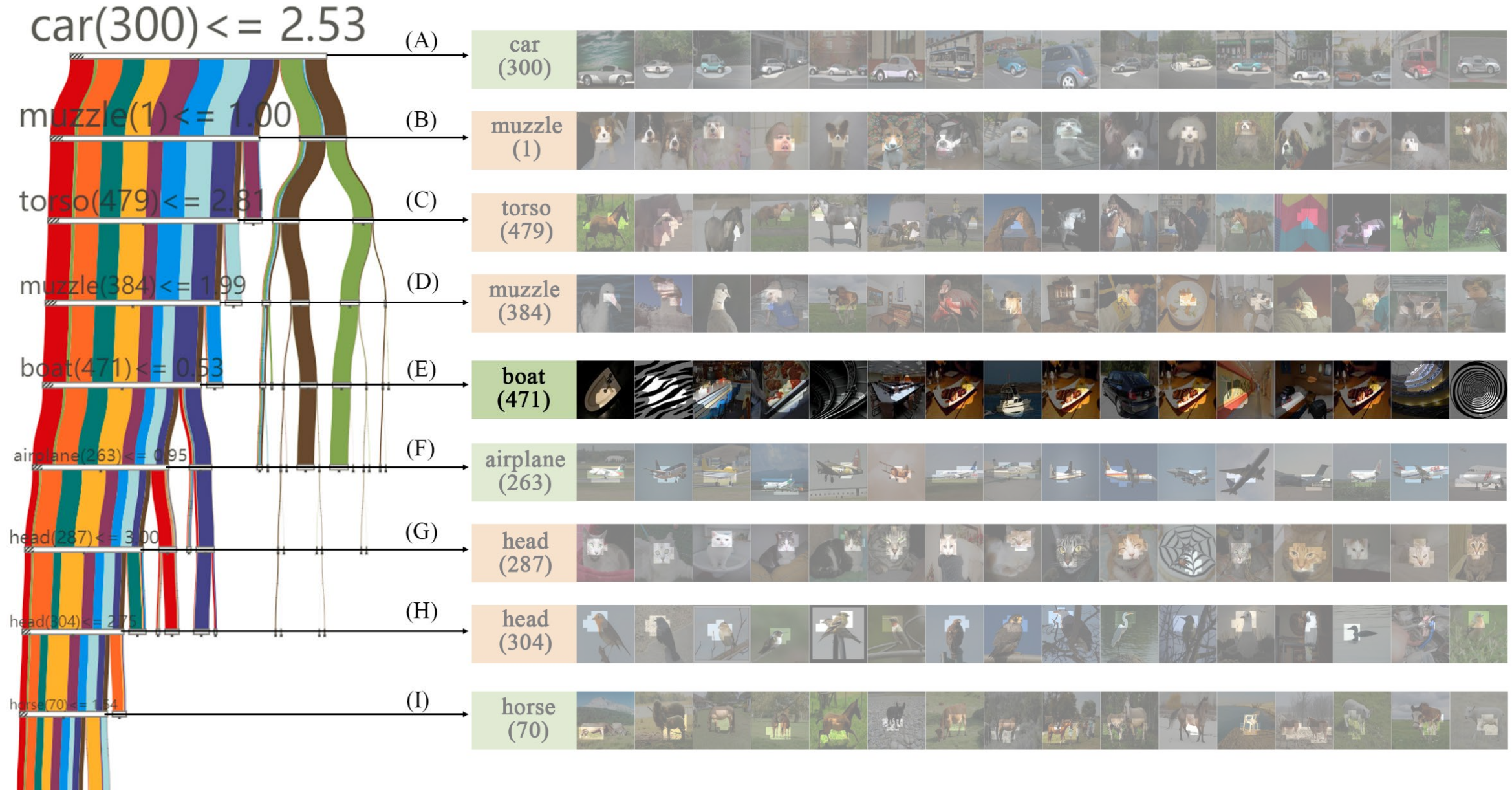
Use Case 1: Interpreting Surrogate Decision Tree

data class

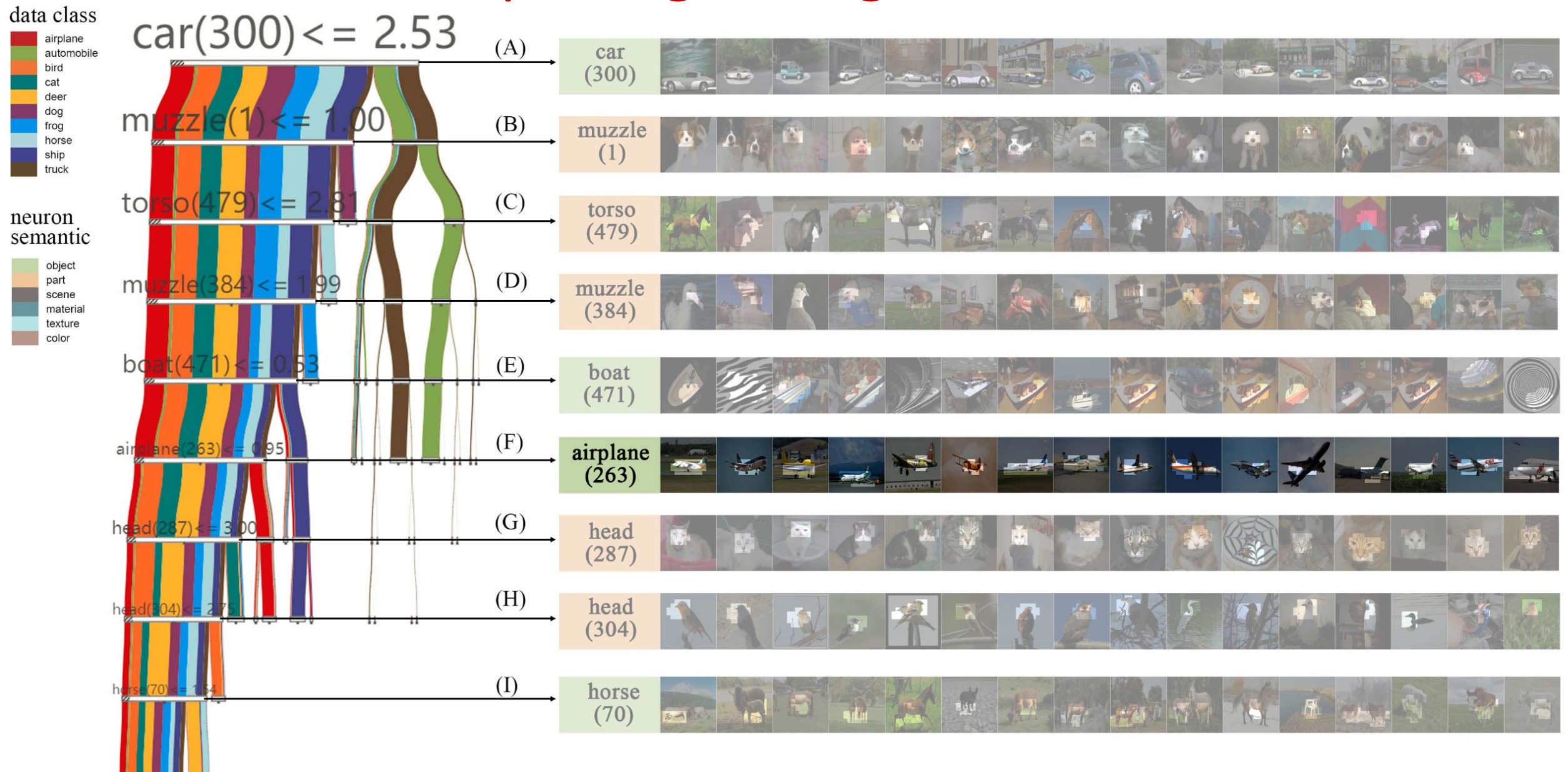
- airplane
- automobile
- bird
- cat
- deer
- dog
- frog
- horse
- ship
- truck

neuron semantic

- object
- part
- scene
- material
- texture
- color



Use Case 1: Interpreting Surrogate Decision Tree



Use Case 1: Interpreting Surrogate Decision Tree

data class

- airplane
- automobile
- bird
- cat
- deer
- dog
- frog
- horse
- ship
- truck

neuron semantic

- object
- part
- scene
- material
- texture
- color



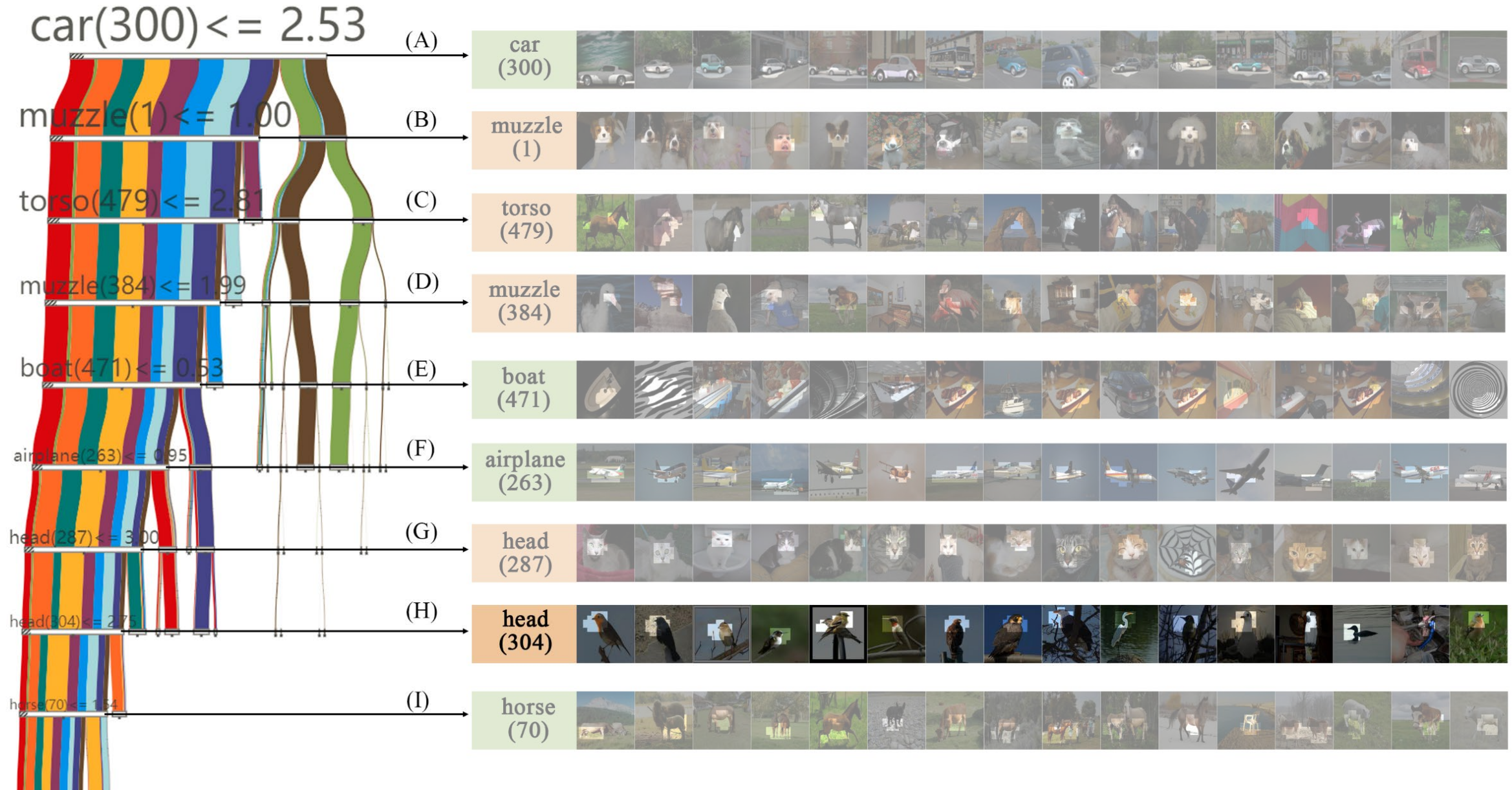
Use Case 1: Interpreting Surrogate Decision Tree

data class

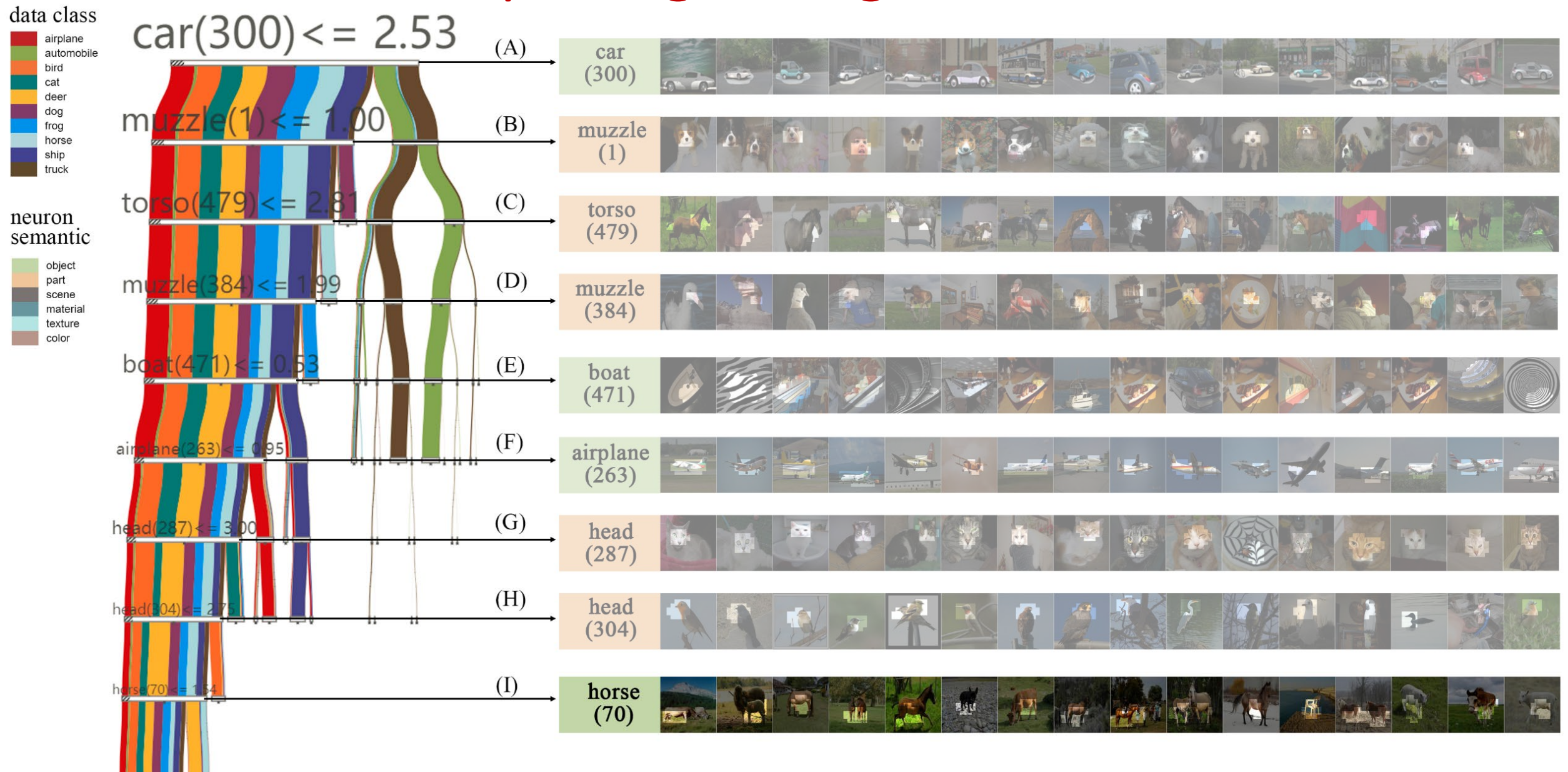
- airplane
- automobile
- bird
- cat
- deer
- dog
- frog
- horse
- ship
- truck

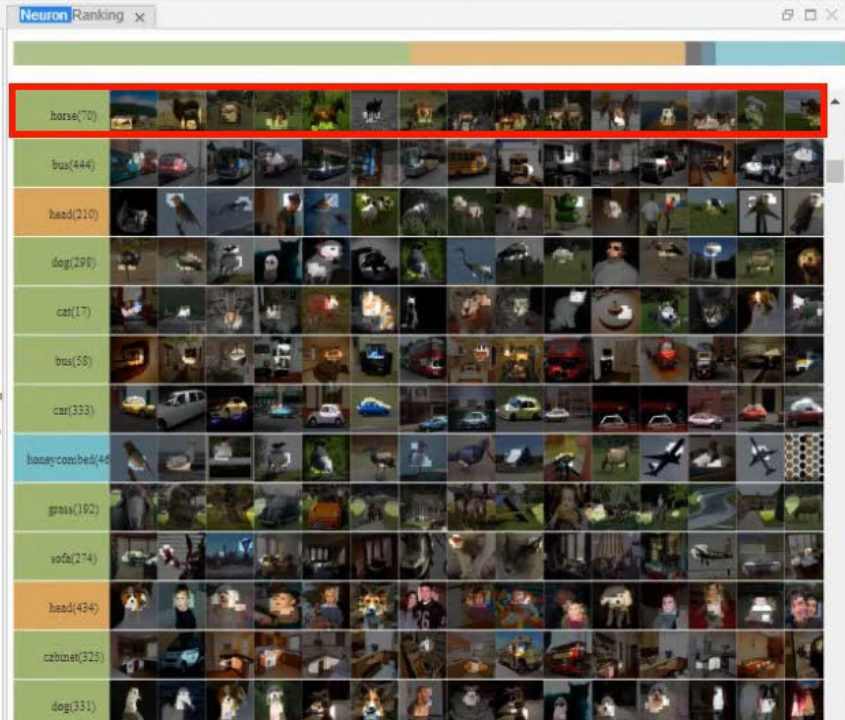
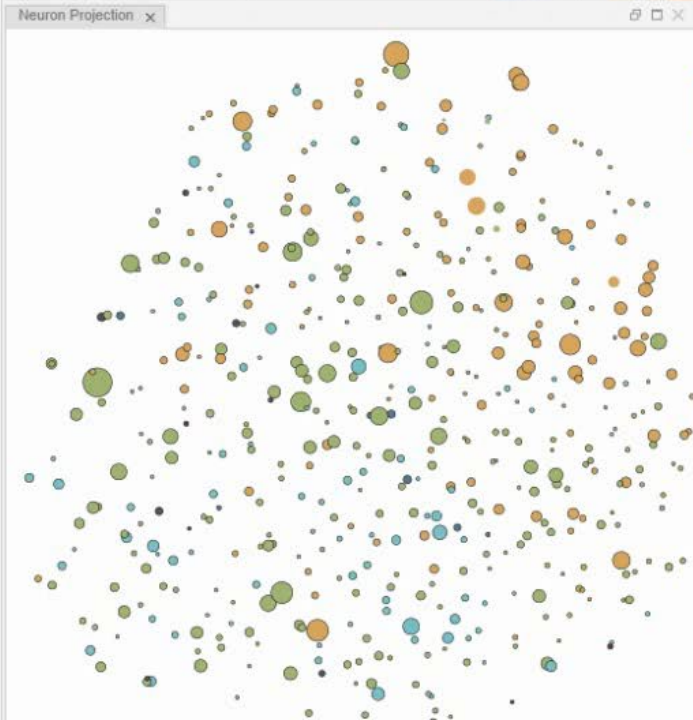
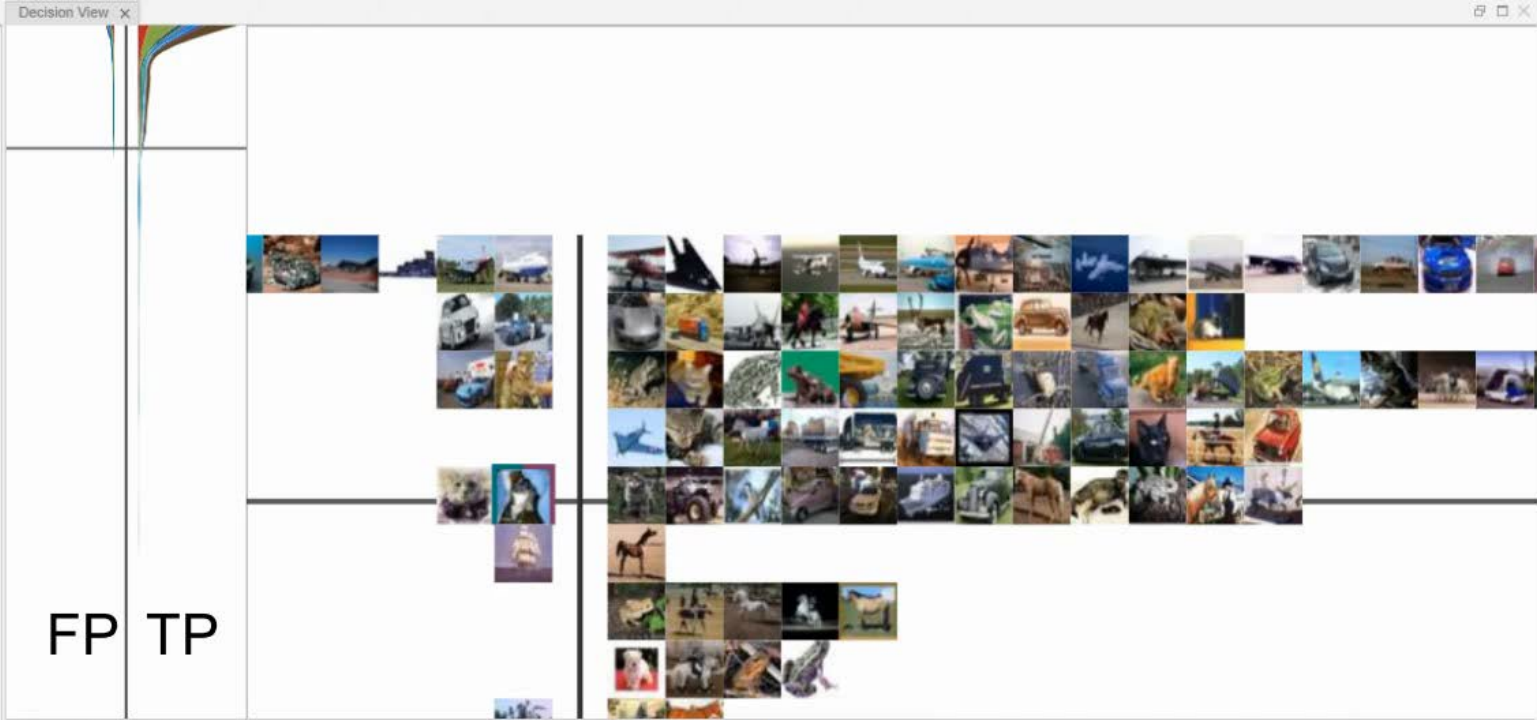
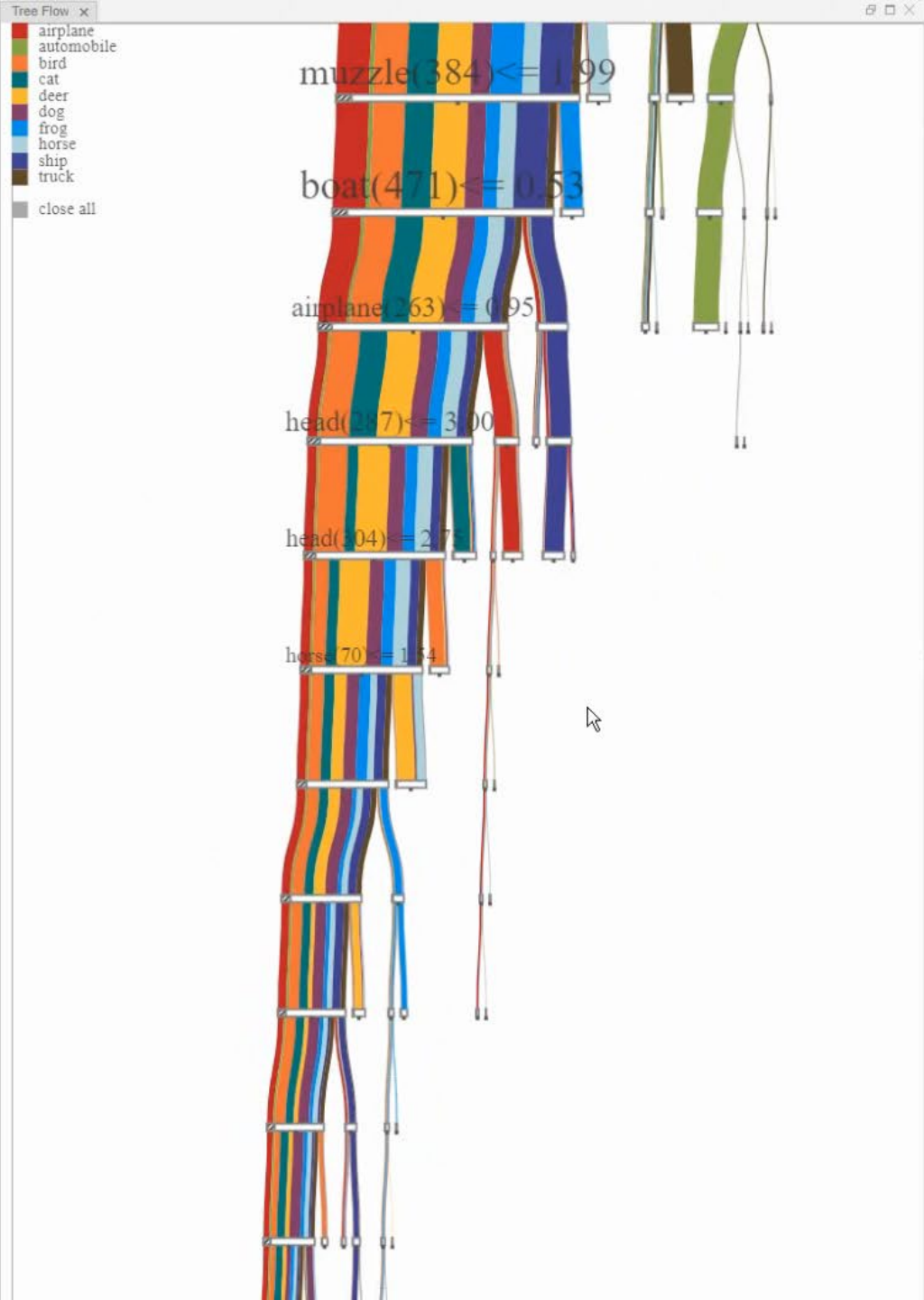
neuron semantic

- object
- part
- scene
- material
- texture
- color



Use Case 1: Interpreting Surrogate Decision Tree





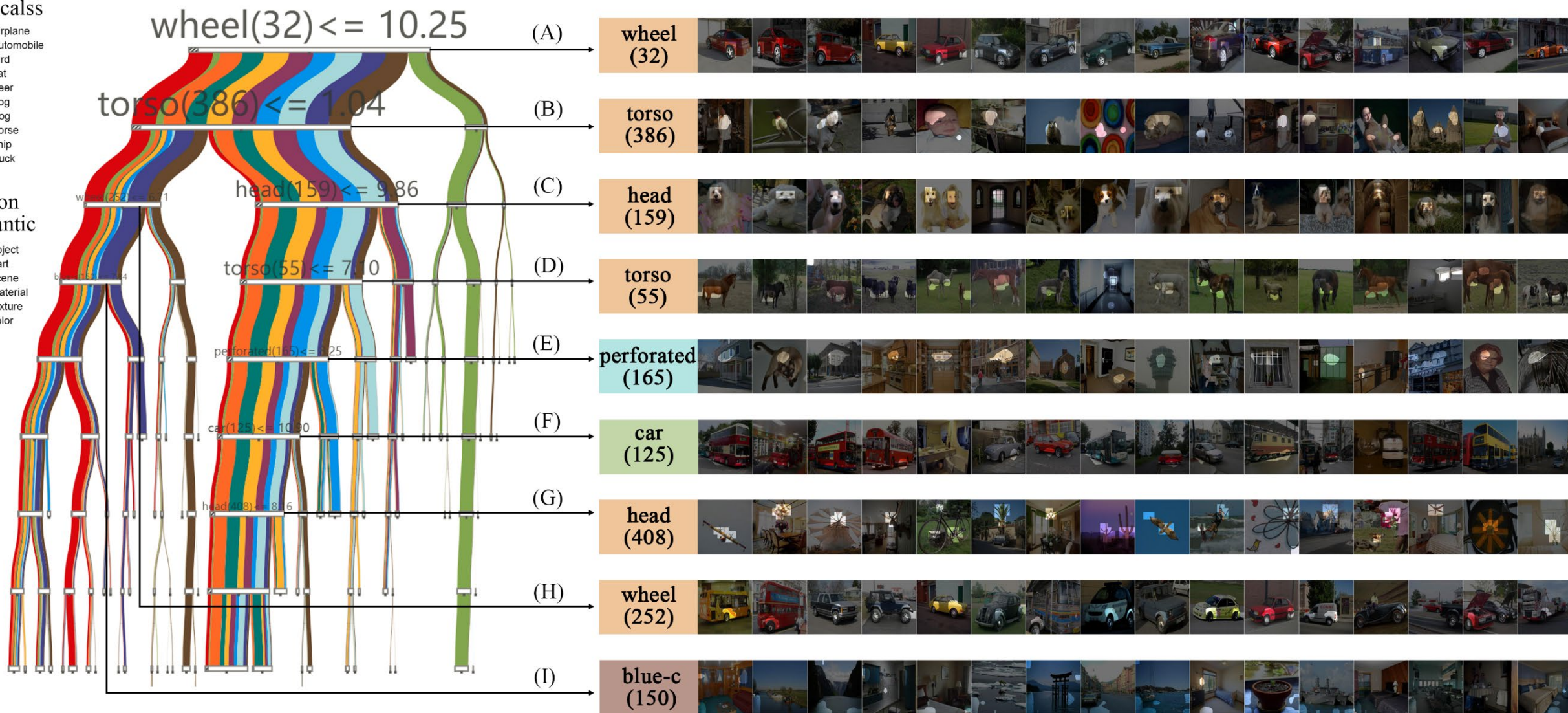
Use Case 2: Comparing Surrogate Decision Trees

data calss

- airplane
- automobile
- bird
- cat
- deer
- dog
- frog
- horse
- ship
- truck

neuron semantic

- object
- part
- scene
- material
- texture
- color



Conclusion & Discussions

1 Conclusion

- New strategy to convert CNNs to surrogate decision trees
- CNN2DT, a visual analytics system
- Use cases and user study

2 Potential Users

- Non-DL experts & ML practitioners

3 Limitations & Future Work

- Scalability for color encoding, number of classes
- Dependent on semantic database
- Multi-semantics of each neuron | Non-interpretable neurons
- Further evaluations on other architectures