

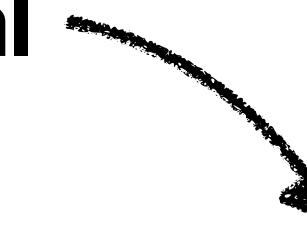


Visualization-Based Neural Network Introspection

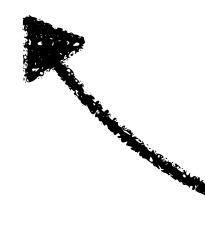
Alex Bäuerle

Defense, 21.12.2022

Potential

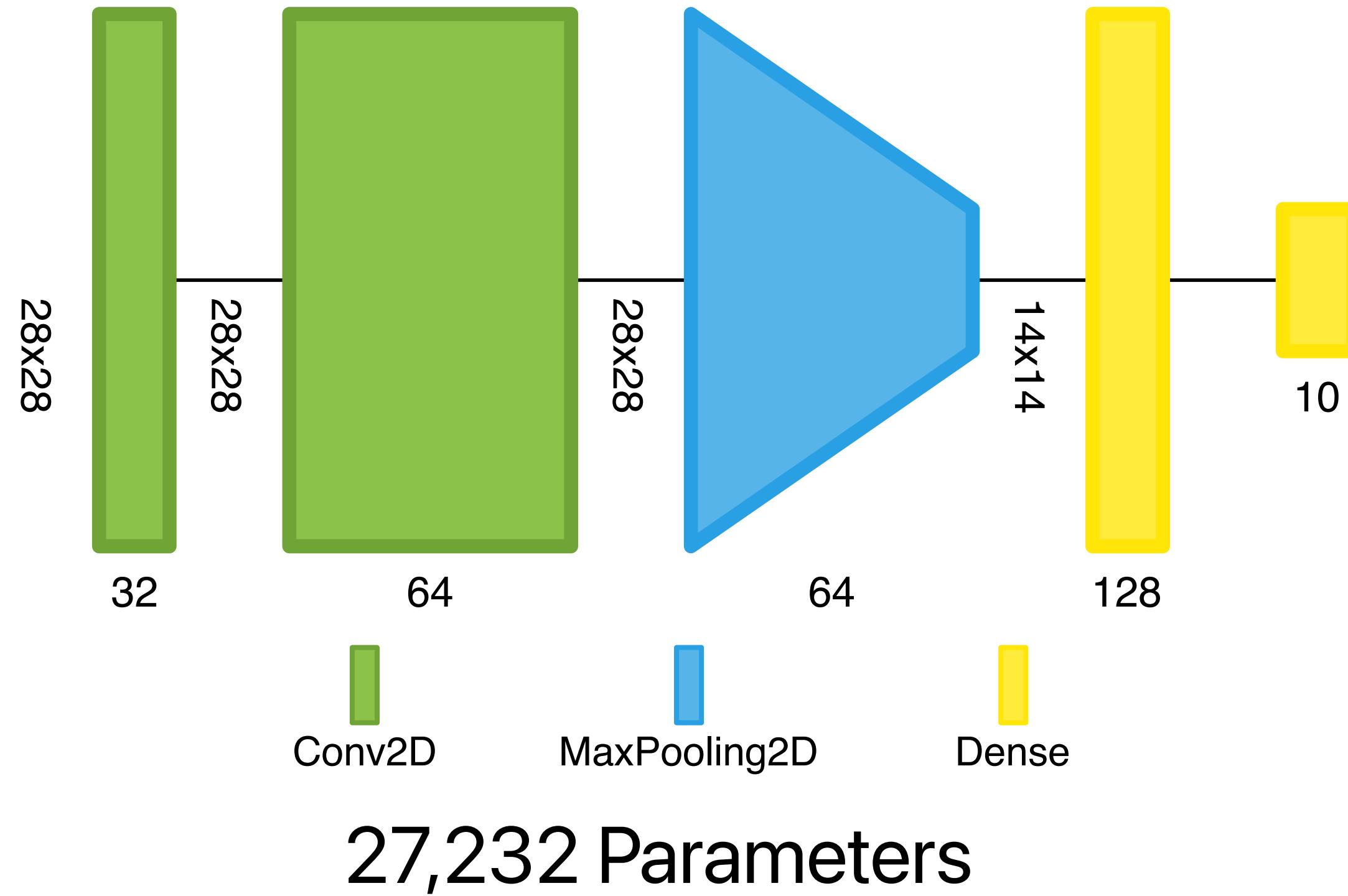


“Recent and rapid advances in AI research and technologies are widely expected to bring about pervasive and far-reaching social transformation on a global scale. The celebratory rhetoric that accompanied the emergence of these technologies several years ago [...] has recently become considerably more muted in the face of rising public anxiety about the possible adverse effects [...]”



Risk

Karen Yeung 2020, International Legal Materials, Recommendation of the council on artificial intelligence (OECD)



ChatGPT and How AI Disrupts Industries

Ajay Agrawal, Joshua Gans, and Avi Goldfarb, Harvard Business Review, December 2022

How Will AI Impact the Transportation Industry?

Maayan, Dataversity, April 2021

How Machine Learning Will Transform Biomedicine

Goecks, et al. "How machine learning will transform biomedicine.", Cell, 2020

AI Is a Game-Changer in the Fight Against Hunger and Poverty. Here's Why

Bennington-Castro, NBCNews, June 2021

AI Art Is Challenging the Boundaries of Curation

Raphaël Millière, WIRED, July 2022

AI bot ChatGPT writes smart essays – should professors worry?

Chris Stokel-Walker, Nature News Explainer, December 2022

How AI-Driven Technology Is Increasing Food Security, And Improving The Lives Of Farmers Worldwide

Forbes, August 2021

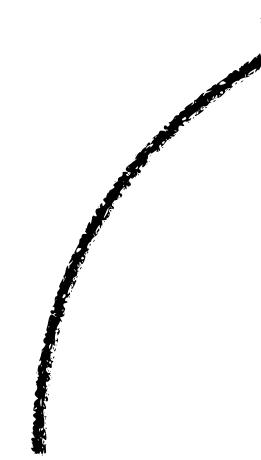
Here's how AI can help fight climate change

World Economic Forum, August 2021

“The full cycle of a machine learning project is not just modeling. It is finding the right data, deploying it, monitoring it, feeding data back, showing safety [...]. [That goes] **beyond doing well on the test set**, which fortunately or unfortunately is what we in machine learning are great at.”

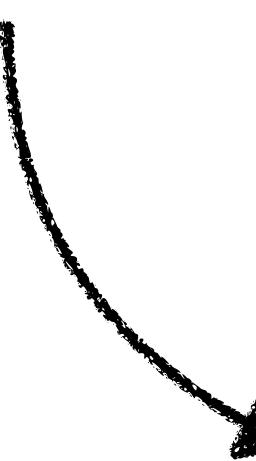
Andrew Ng, IEEE Spectrum, May 2021

Introspection



Develop methods for ML introspection and make them accessible through visualization.

Visualization



What has a model learned?

How was my model trained?

What does my model do?

Where does my model fail?

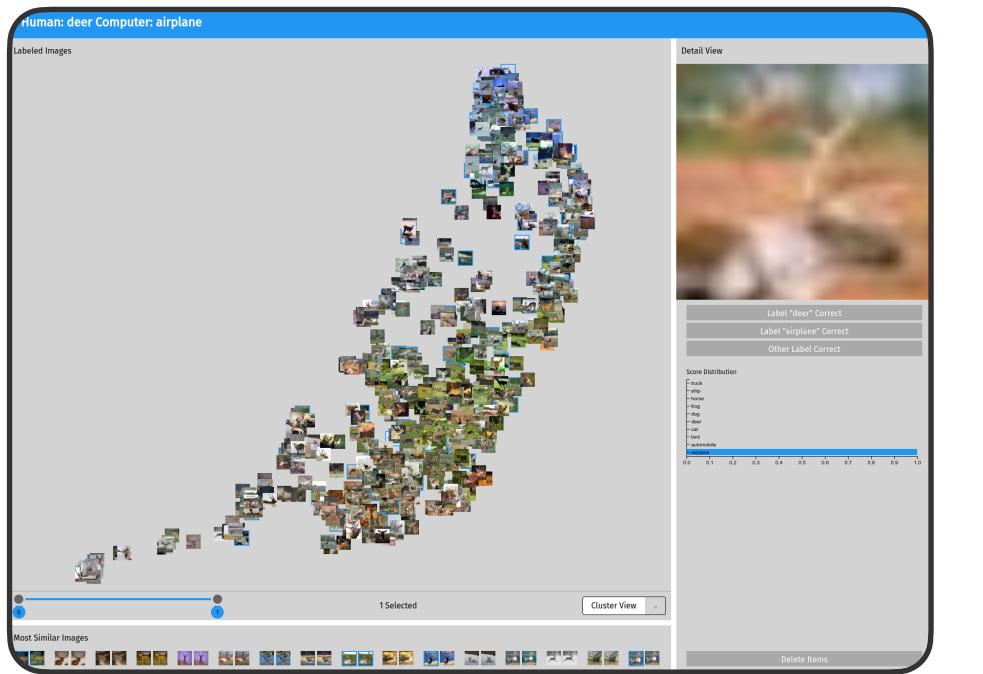
How do we communicate results?

How do we make introspection usable?

How do we explore data?

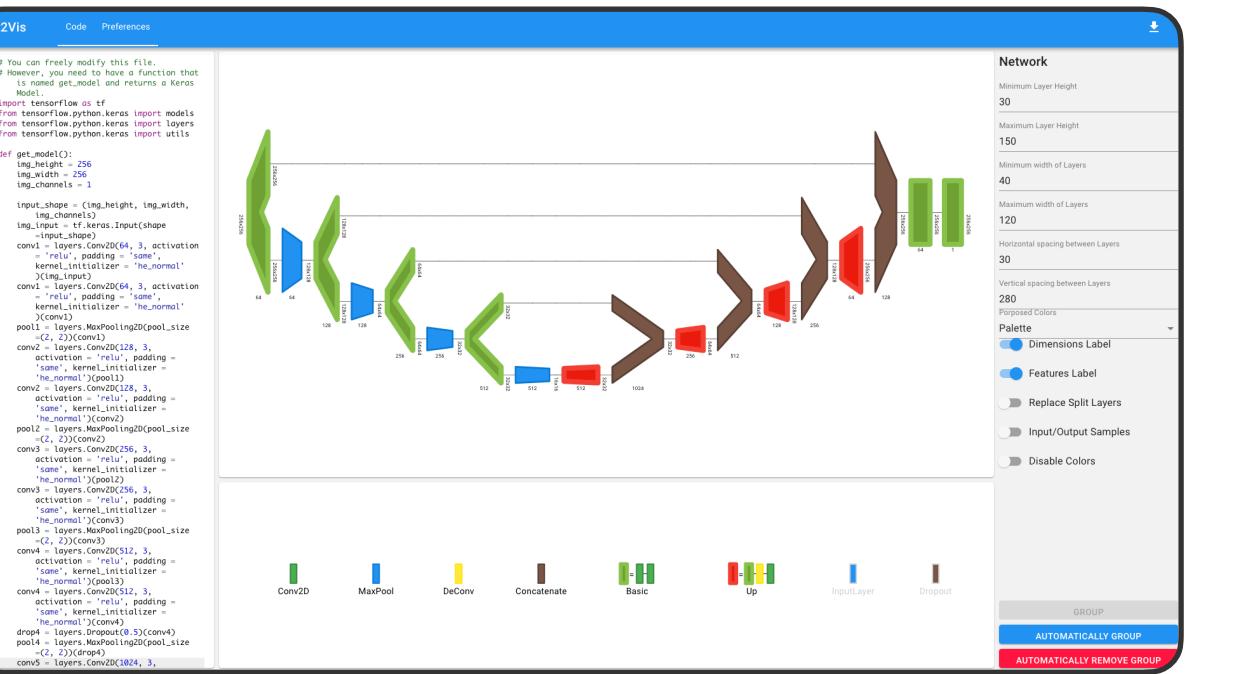
Can we see how a model works?

Quality Assurance



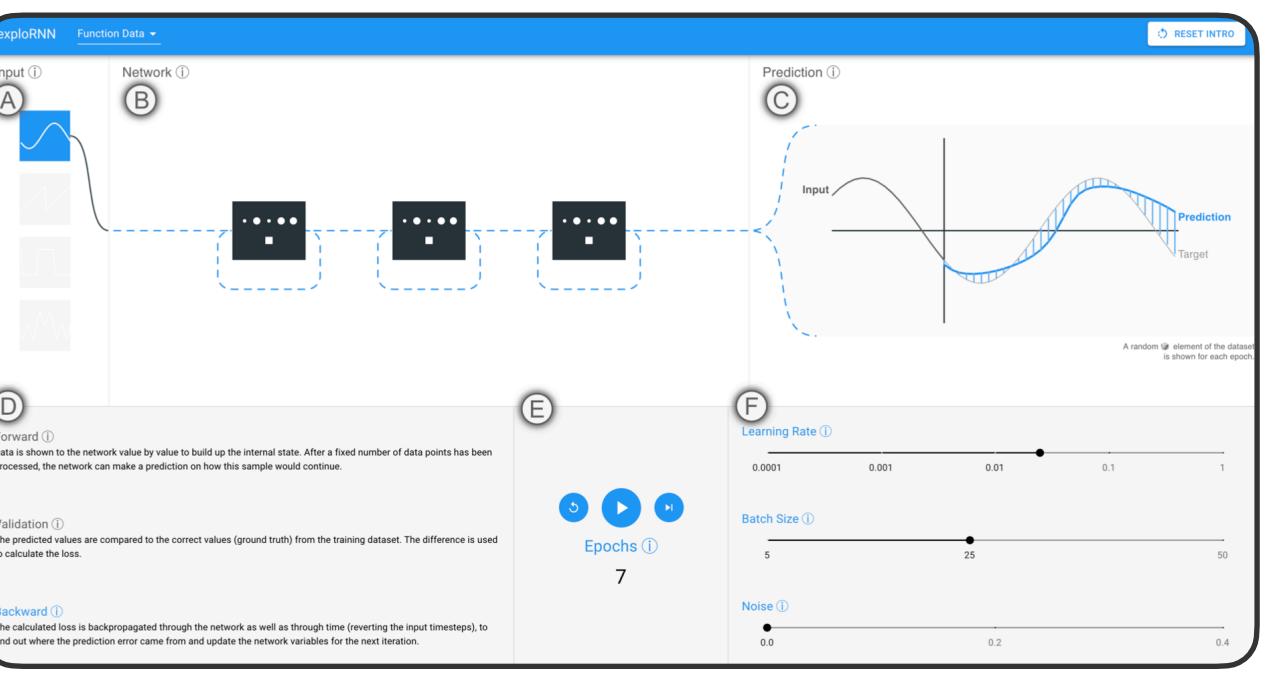
Bäuerle et al. 2020, CGF
Classifier-Guided Visual Correction of Noisy Labels for Image Classification Tasks

Communication

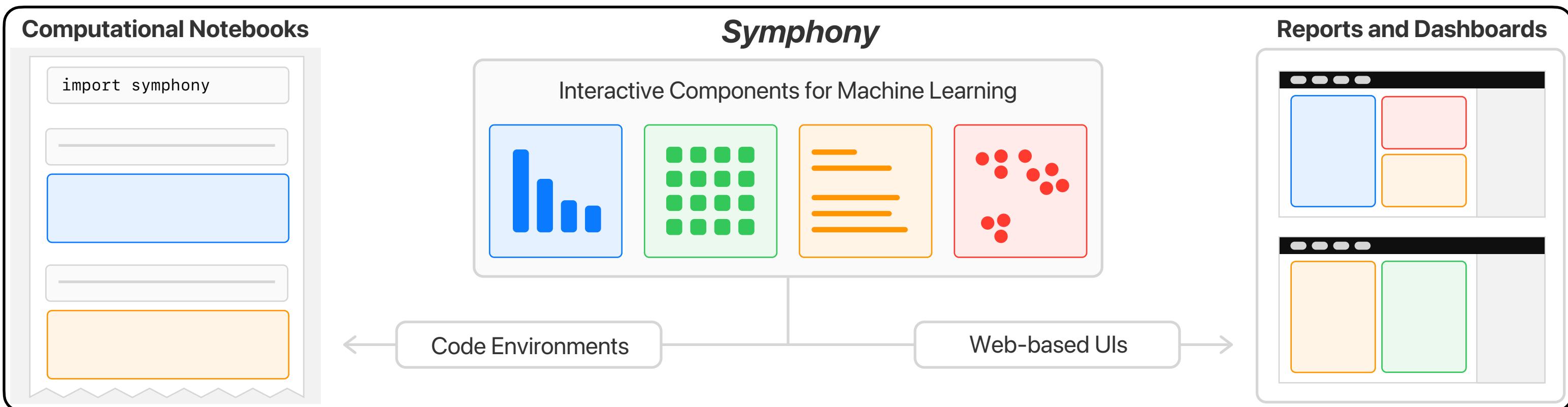


Bäuerle et al. 2021, TVCG
Net2Vis - A Visual Grammar for Automatically Generating Publication-Tailored CNN Architecture Visualizations

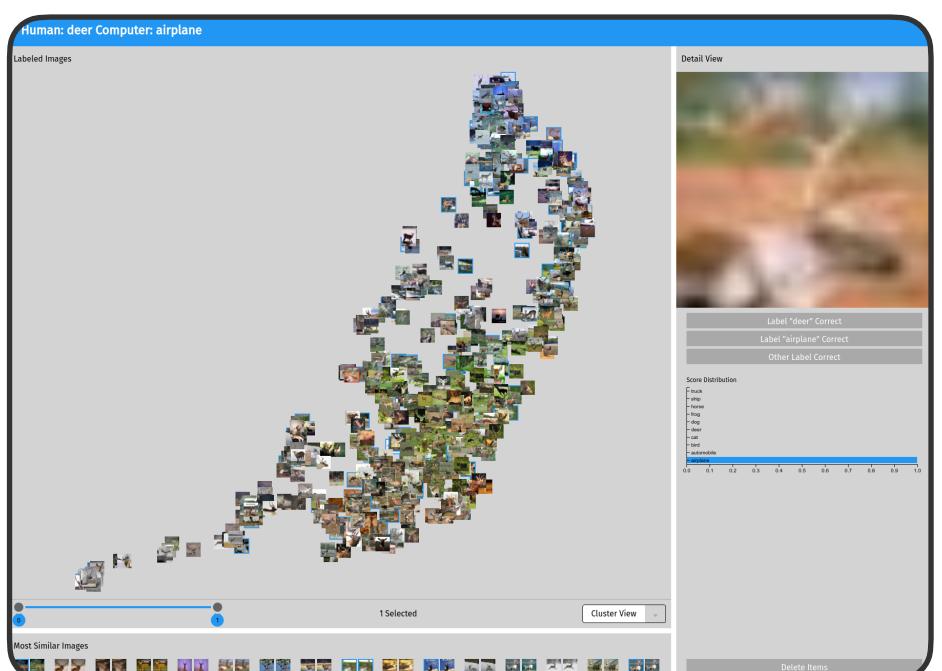
Education



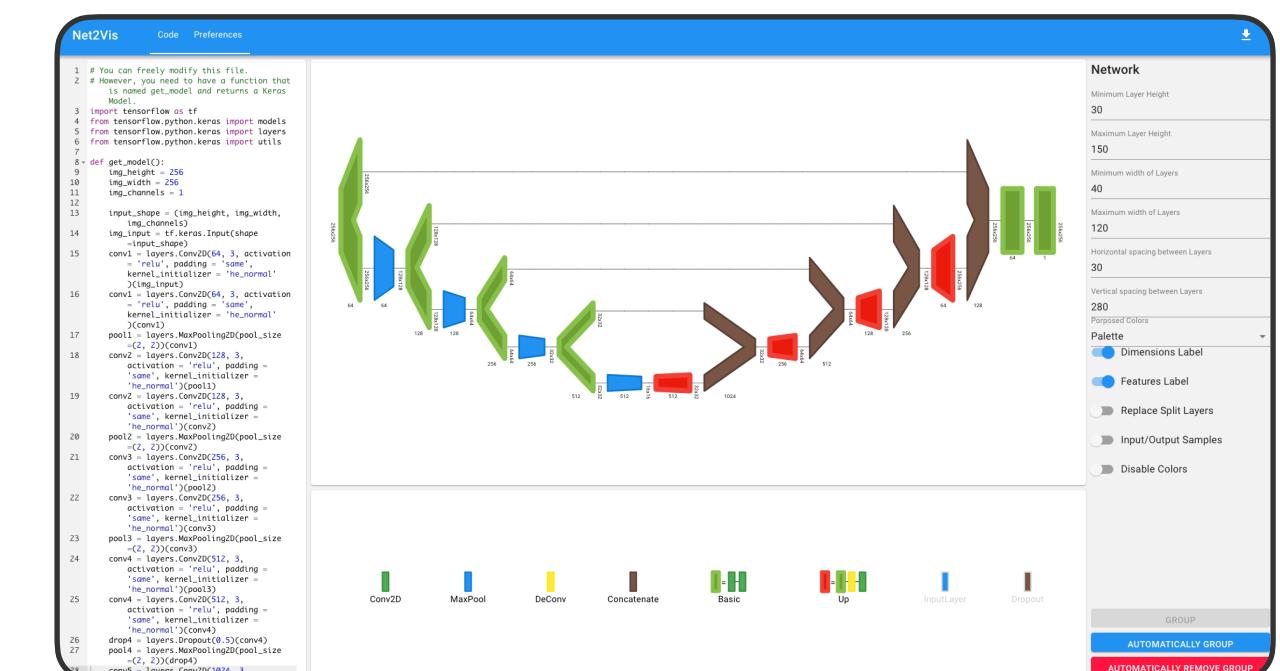
Bäuerle et al. 2022, The Visual Computer
exploRNN: Understanding Recurrent Neural Networks through Visual Exploration



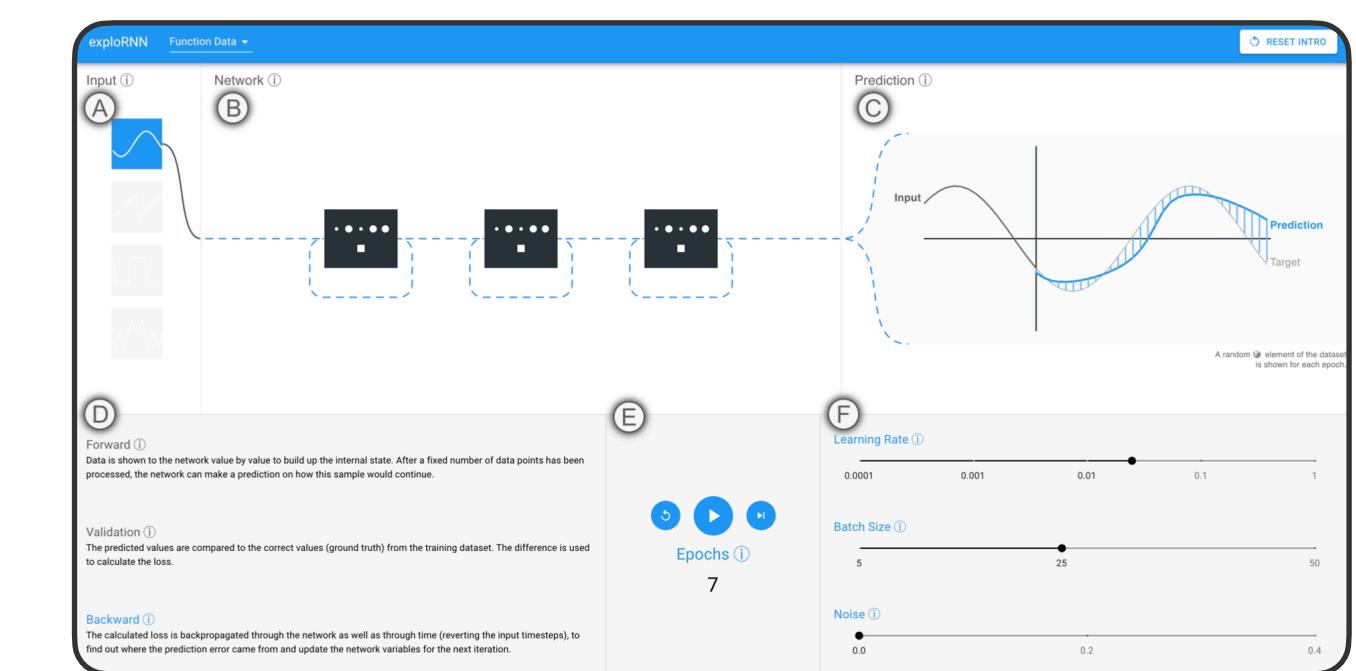
Bäuerle and Cabrera et al. 2022, CHI
Symphony: Composing Interactive Interfaces for Machine Learning



Bäuerle et al. 2020, CGF
Classifier-Guided Visual Correction of Noisy Labels for Image Classification Tasks

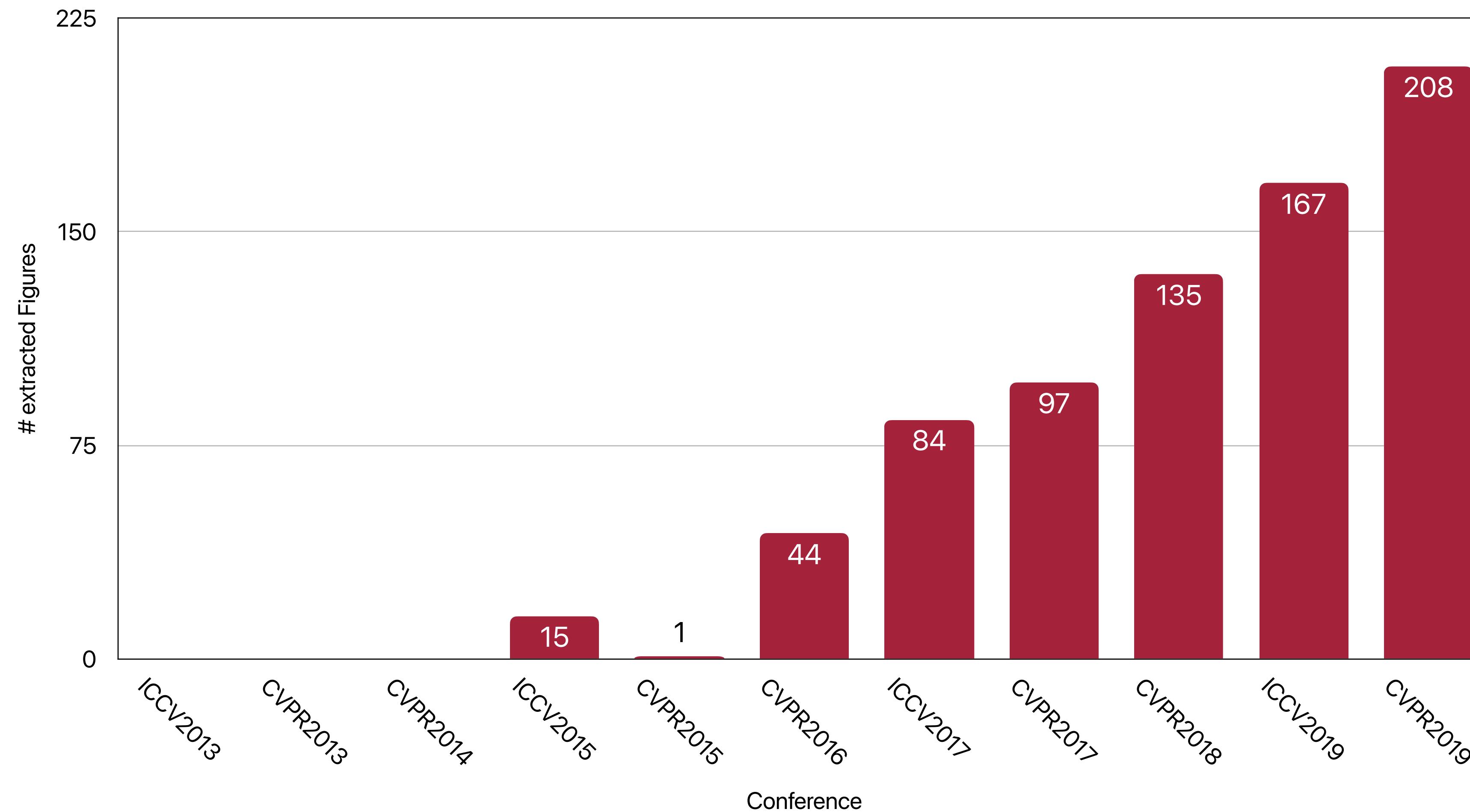


Bäuerle et al. 2021, TVCG
Net2Vis - A Visual Grammar for Automatically Generating Publication-Tailored CNN Architecture Visualizations

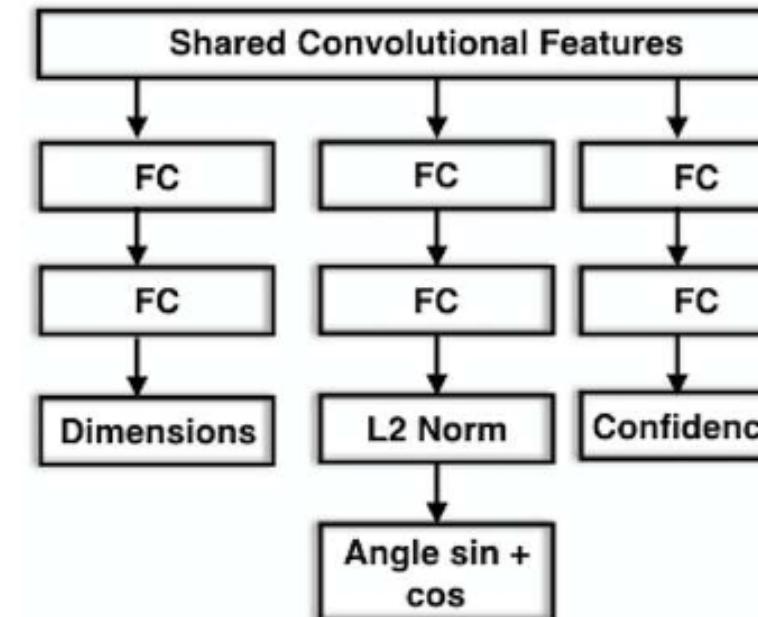


Bäuerle et al. 2022, The Visual Computer
exploRNN: Understanding Recurrent Neural Networks through Visual Exploration

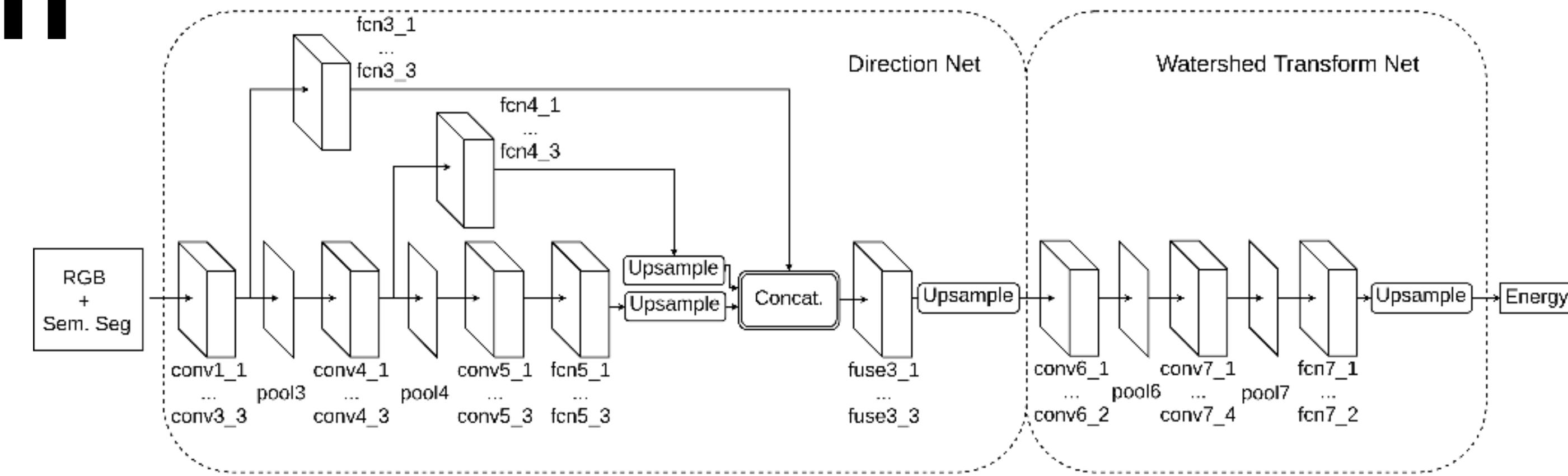
Architecture Figures in Publications



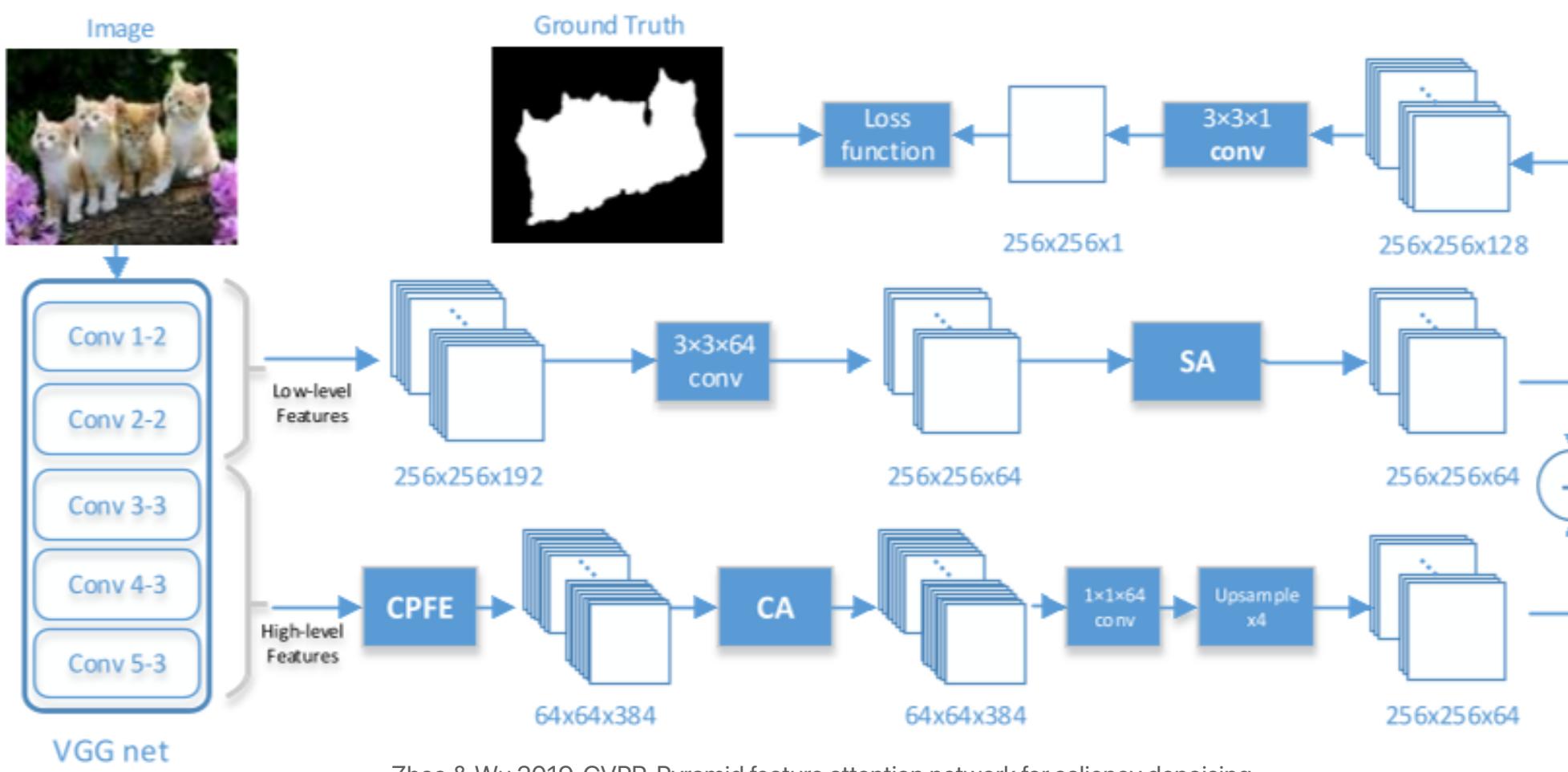
Current Situation



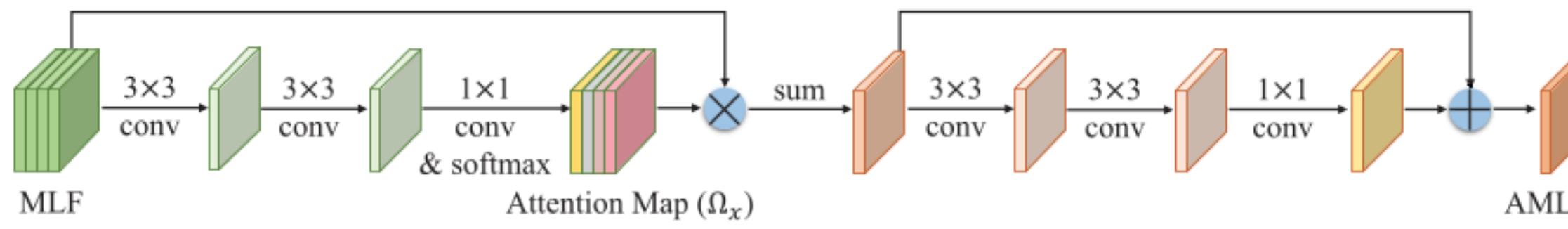
Mousavian et. al. 2017, CVPR, 3d bounding box estimation using deep learning and geometry.



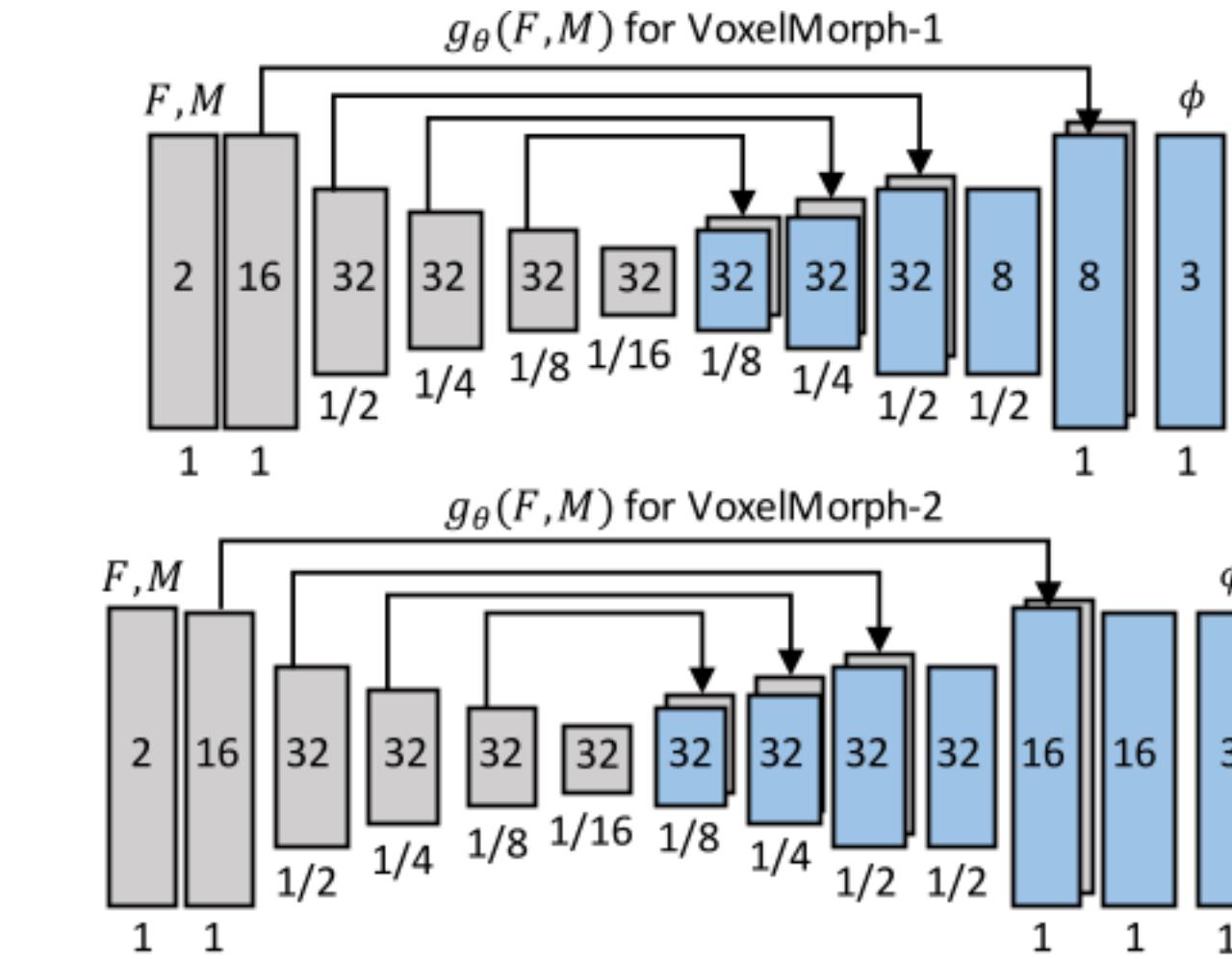
Bai & Urtasun 2017, CVPR, Deep watershed transform for instance segmentation.



Zhao & Wu 2019, CVPR, Pyramid feature attention network for saliency denoising.



Deng et. al. 2019, CVPR, Deep multi-model fusion for single image dehazing.



Balakrishnan et. al. 2018, CVPR, An unsupervised learning model for deformable medical image registration.

Net2Vis

IEEE TRANSACTIONS ON VISUALIZATION AND COMPUTER GRAPHICS

1

Net2Vis – A Visual Grammar for Automatically Generating Publication-Ready CNN Architecture Visualizations

Alex Bäuerle, Christian van Onzenoodt, and Timo Ropinski

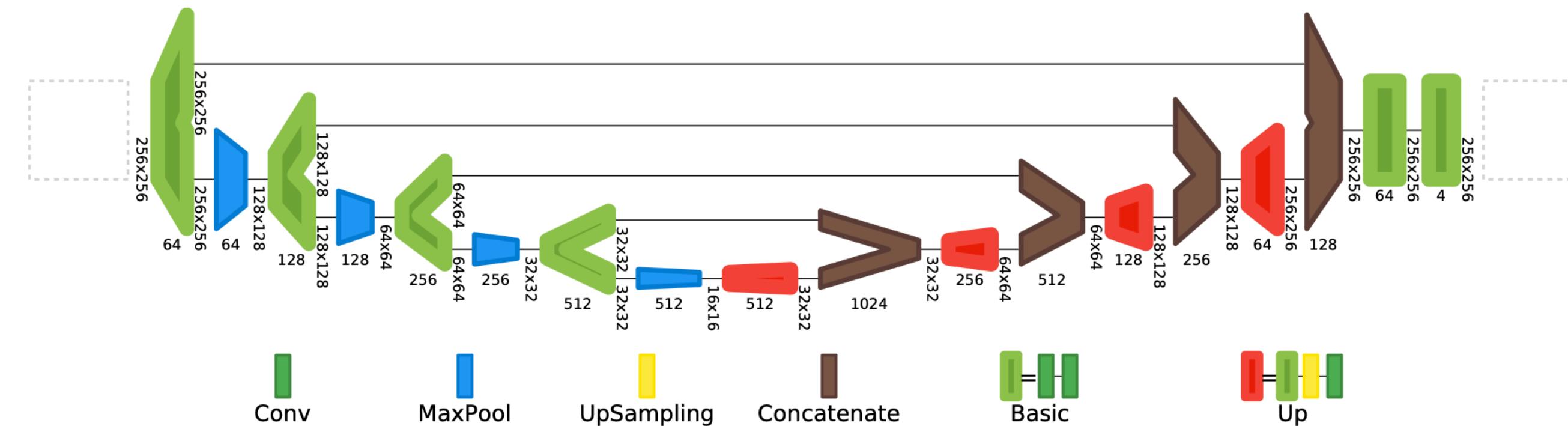


Fig. 1: Visualization of a U-Net variant automatically generated using our approach, based on the Keras code describing the architecture. Data flows from left to right. Glyphs represent layers or aggregates, while lines represent connections. Glyph widths communicate feature size, while heights communicate the spatial resolution. Both values are also given through labels, while dashed boxes on the left and right serve as placeholders to provide input and output samples. The legend communicates layer types and the composition of aggregates.

Abstract—To convey neural network architectures in publications, appropriate visualizations are of great importance. While most current deep learning papers contain such visualizations, these are usually handcrafted just before publication, which results in a lack

Visualization Analysis



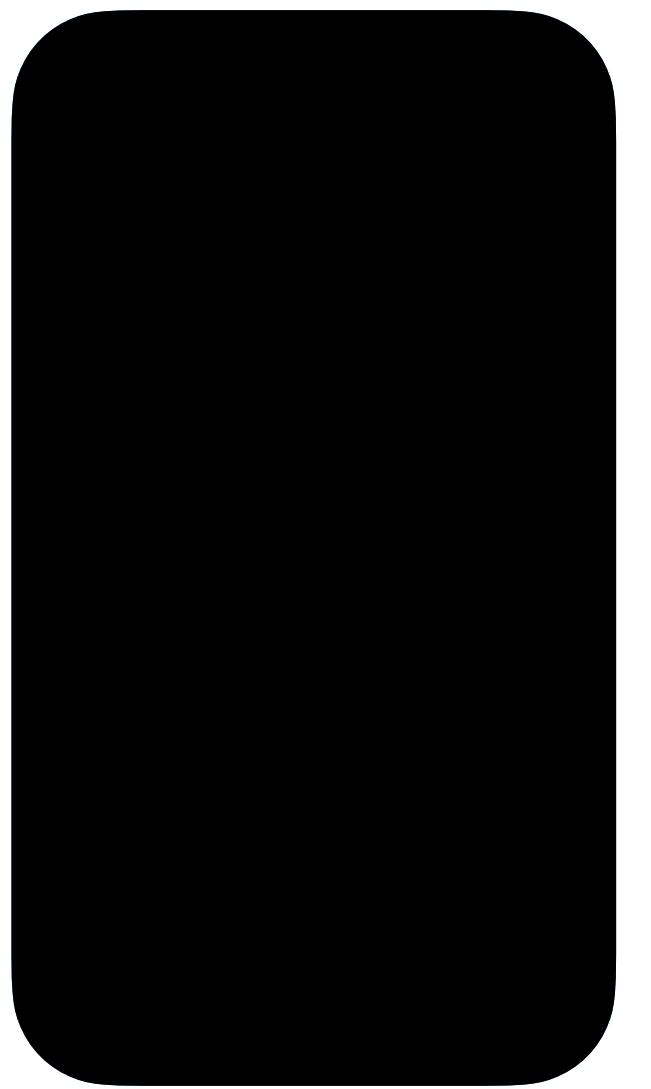
Layer Properties



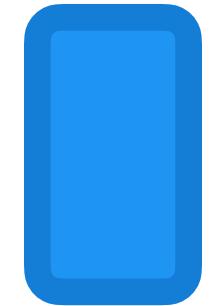
Conv2D



Conv2D

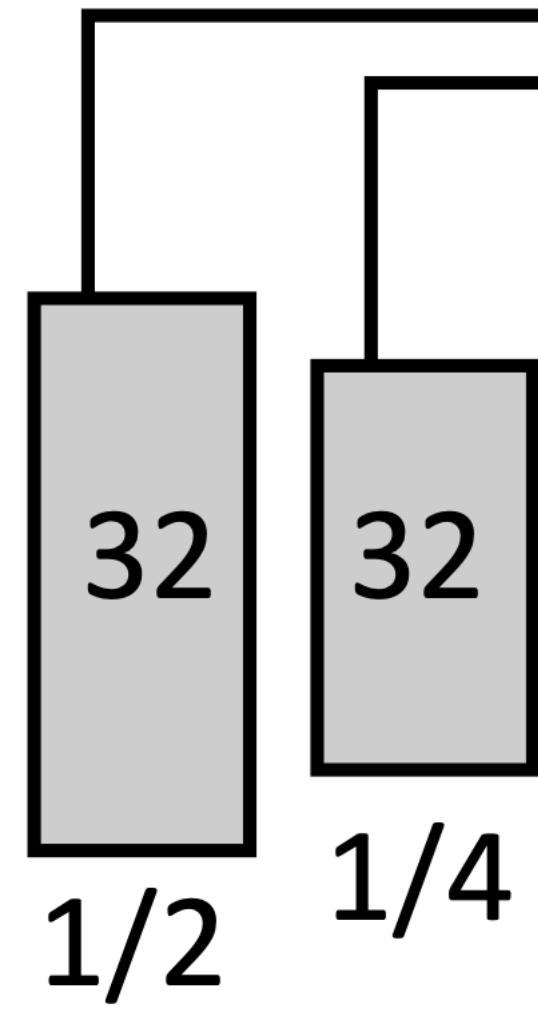


Conv2D

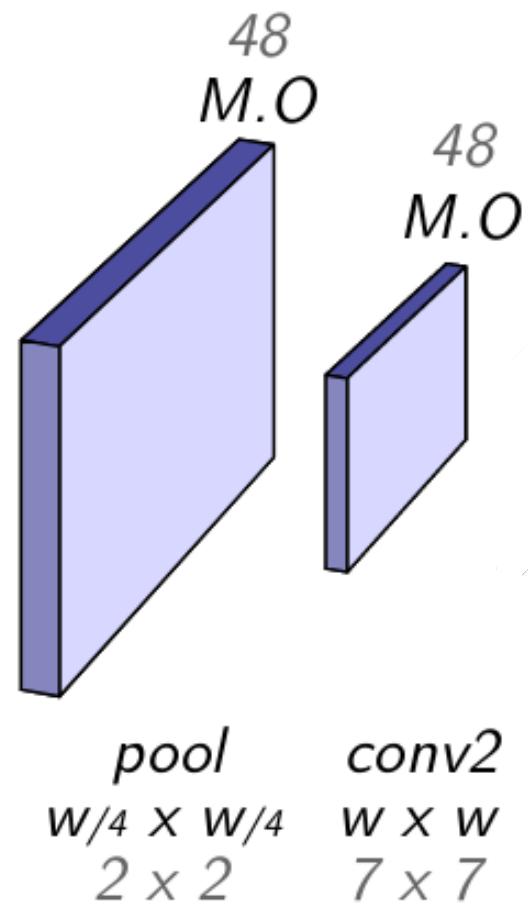


Conv2D

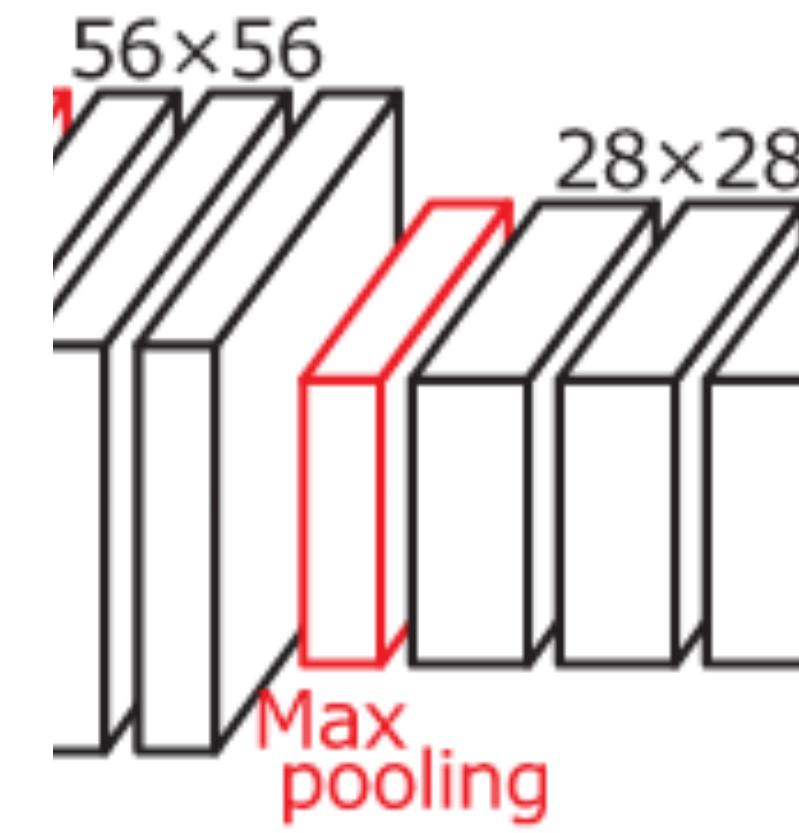
Layer Properties



Balakrishnan et. al., 2018, CVPR, An unsupervised learning model for deformable medical image registration.

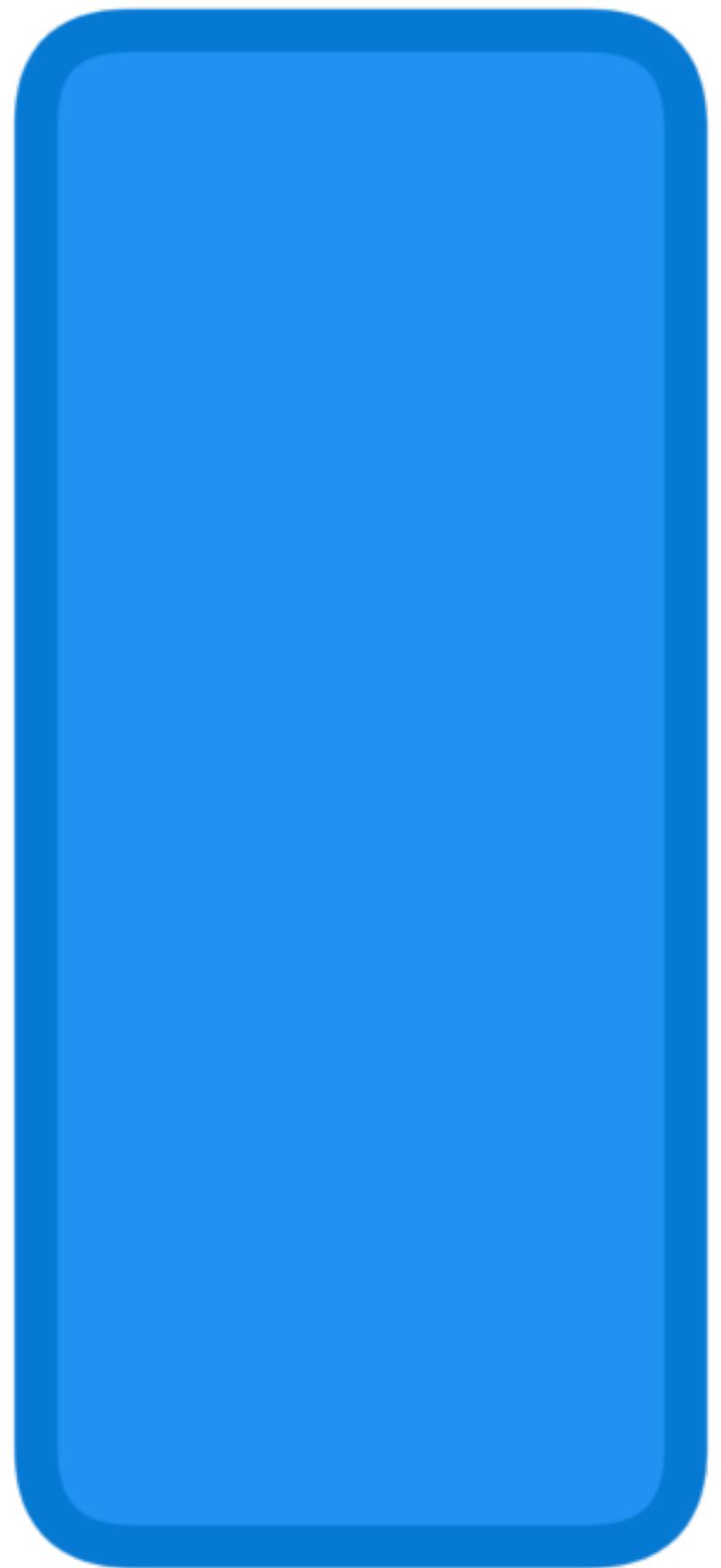


Teney & Martial 2016, ACCV, Leaning to extract motion from videos in convolutional neural networks.



Noh et. al. 2015, ICCV, Learning deconvolution network for semantic segmentation.

Layer Properties



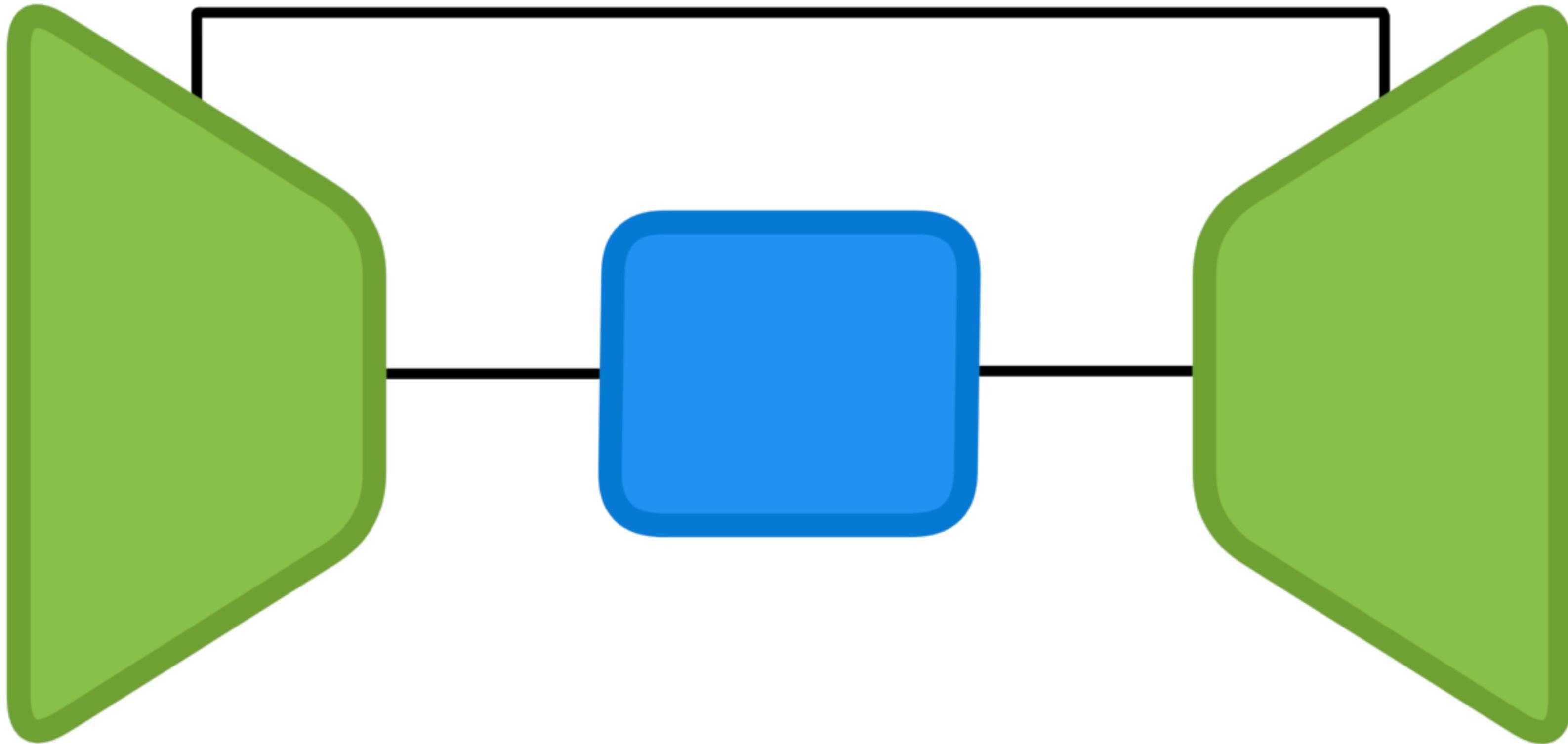
Layer Properties



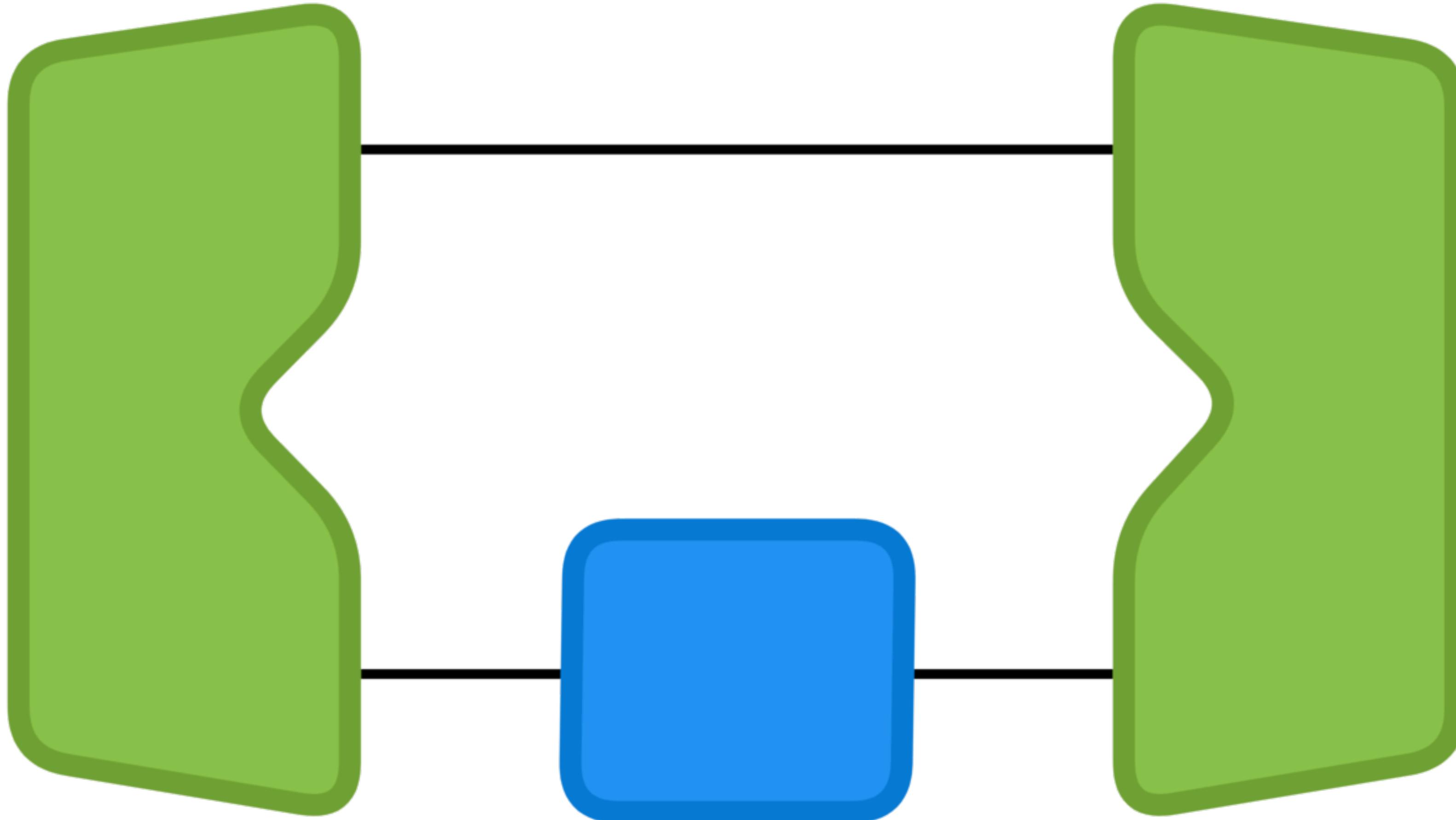
Model Properties



Model Properties



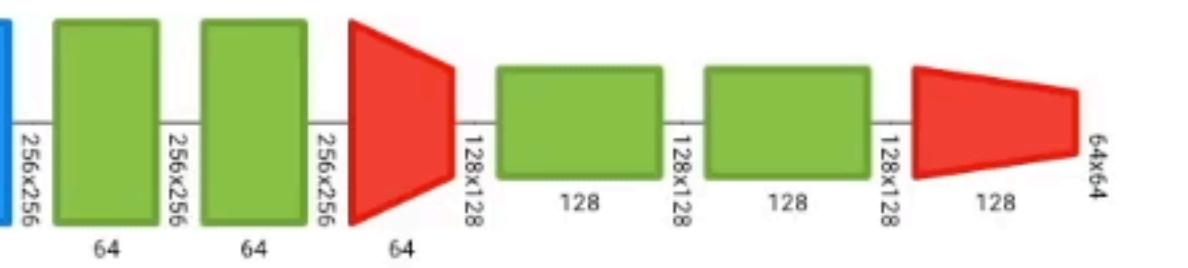
Model Properties



```

1 # You can freely modify this file.
2 # However, you need to have a function that
3 # is named get_model and returns a Keras
4 # Model.
5
6 import keras as k
7 from keras import models
8 from keras import layers
9 from keras import utils
10
11 def get_model():
12     img_height = 256
13     img_width = 256
14     img_channels = 1
15
16     input_shape = (img_height, img_width,
17                     img_channels)
18     img_input = k.Input(shape=input_shape)
19     conv1 = layers.Conv2D(64, 3, activation =
20                         'relu', padding = 'same',
21                         kernel_initializer = 'he_normal'
22                         )(img_input)
23     conv1 = layers.Conv2D(64, 3, activation =
24                         'relu', padding = 'same',
25                         kernel_initializer = 'he_normal'
26                         )(conv1)
27     pool1 = layers.MaxPooling2D(pool_size=(2,
28                                     2))(conv1)
29     conv2 = layers.Conv2D(128, 3, activation
30                         = 'relu', padding = 'same',
31                         kernel_initializer = 'he_normal'
32                         )(pool1)
33     conv2 = layers.Conv2D(128, 3, activation
34                         = 'relu', padding = 'same',
35                         kernel_initializer = 'he_normal'
36                         )(conv2)
37     pool2 = layers.MaxPooling2D(pool_size=(2,
38                                     2))(conv2)
39
40     model = models.Model(img_input, pool2)
41
42     return model
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75

```



InputLayer



Conv2D



MaxPooling2D

Network

Minimum Layer Height

30

Maximum Layer Height

100

Minimum width of Layers

20

Maximum width of Layers

80

Horizontal spacing between Layers

20

Vertical spacing between Layers

200

Proposed Colors

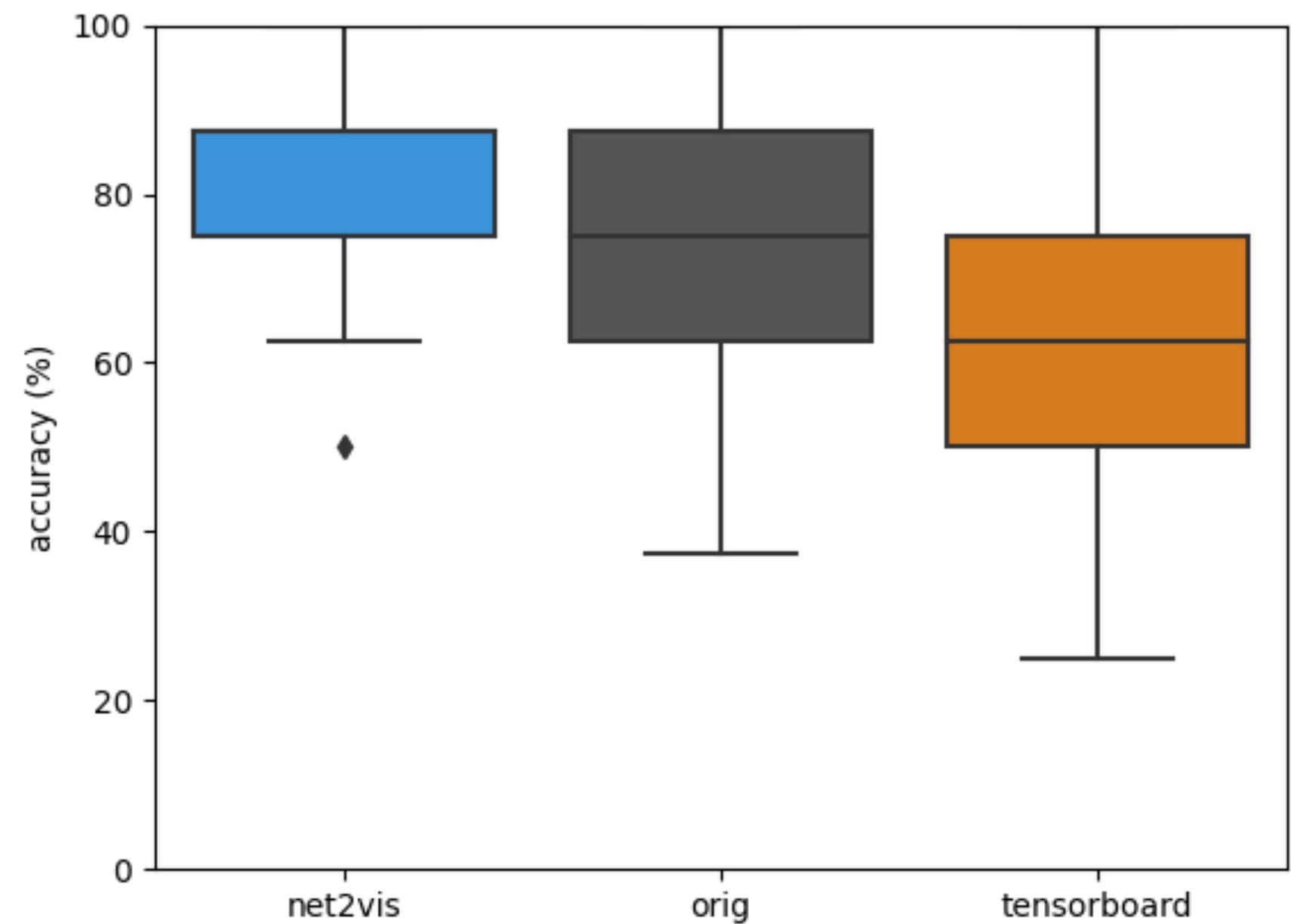
Palette Dimensions Label Features Label Name Label Replace Split Layers Input/Output Samples Disable Colors

Quantitative User Study

10 students after DL course

Architecture visualizations from
Net2Vis, TensorBoard, Paper

1. Architecture questions - best accuracy
2. Comparison questions - preferred
3. Own architecture - excellent usability



Adoption

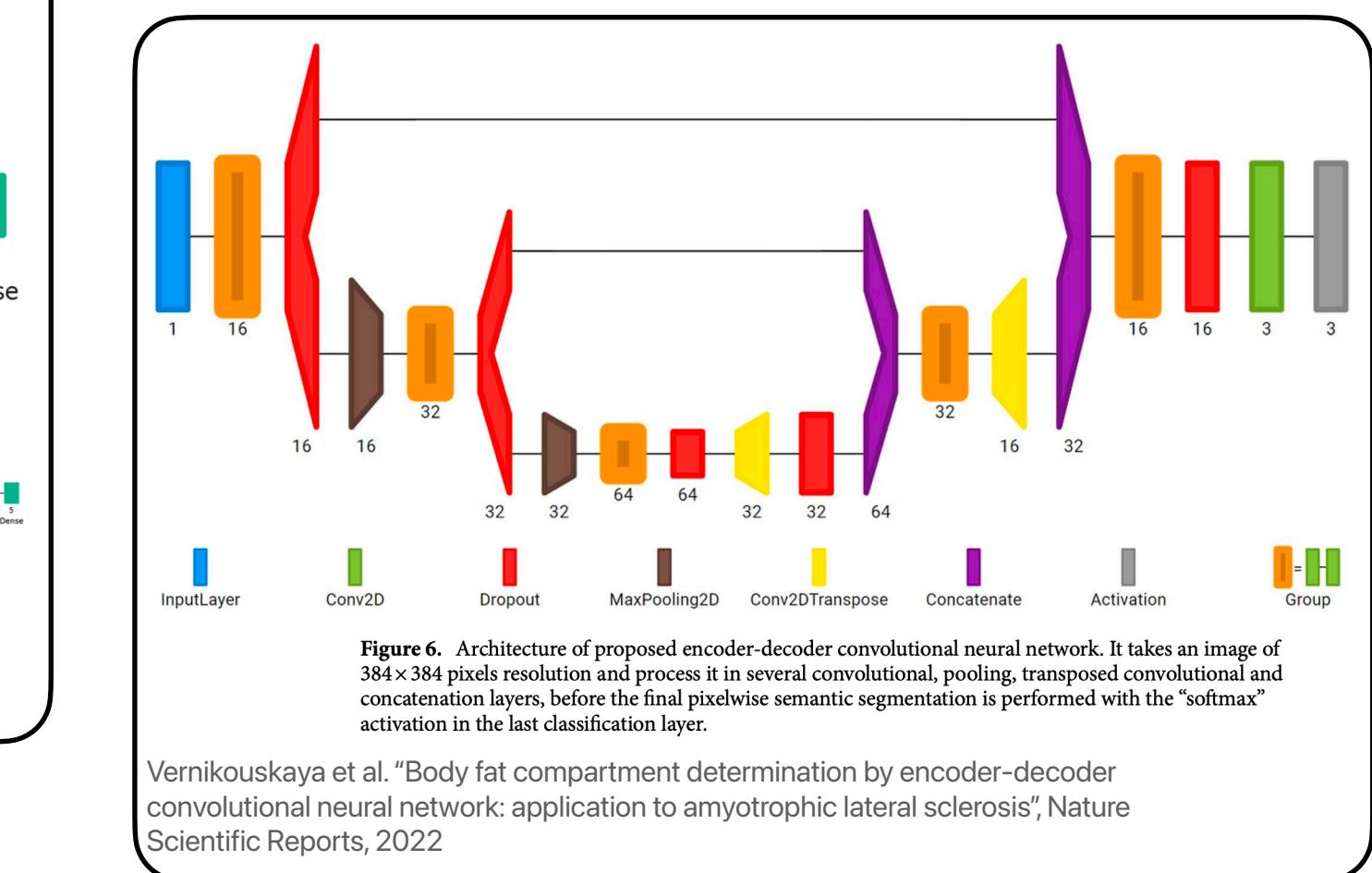
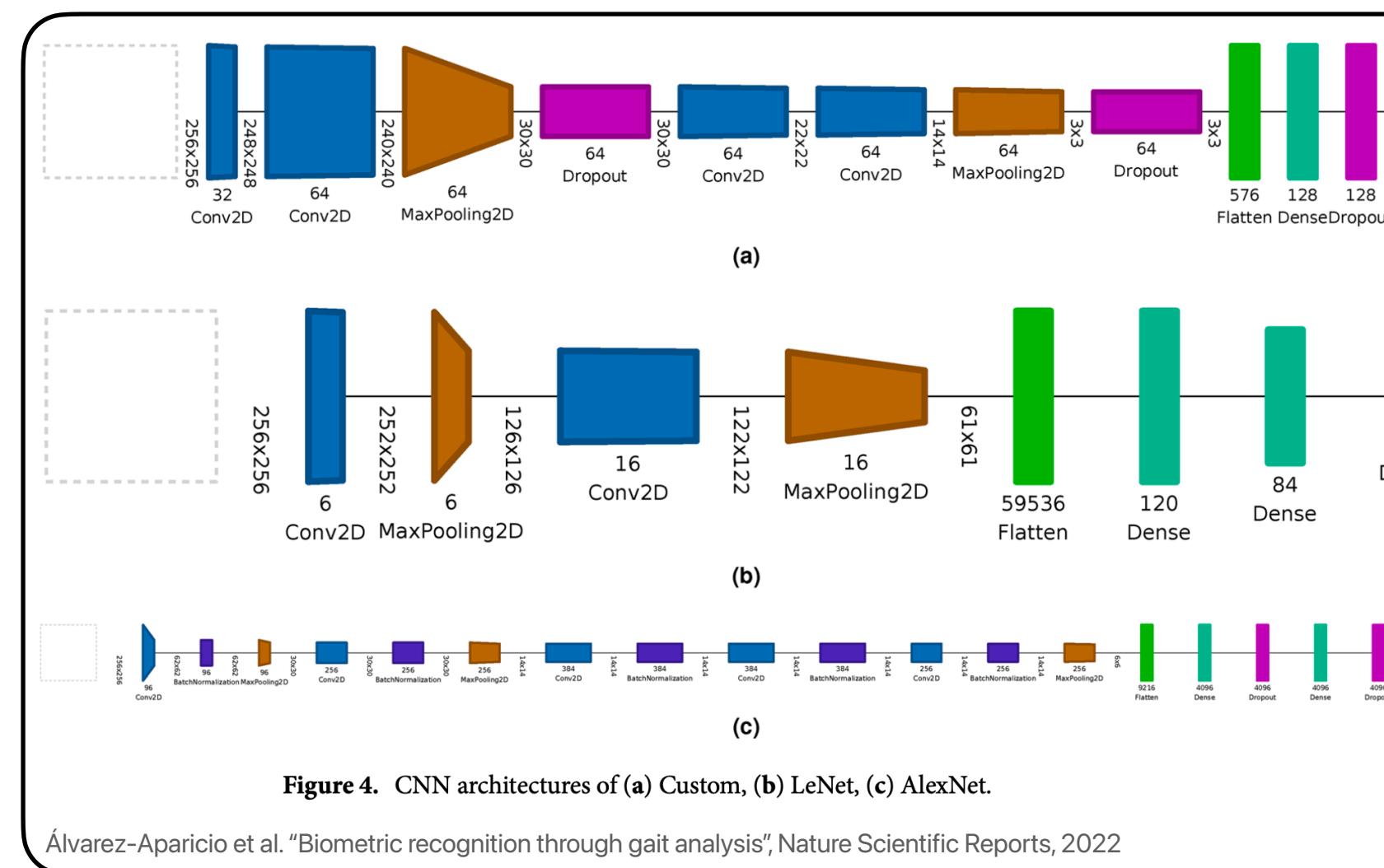
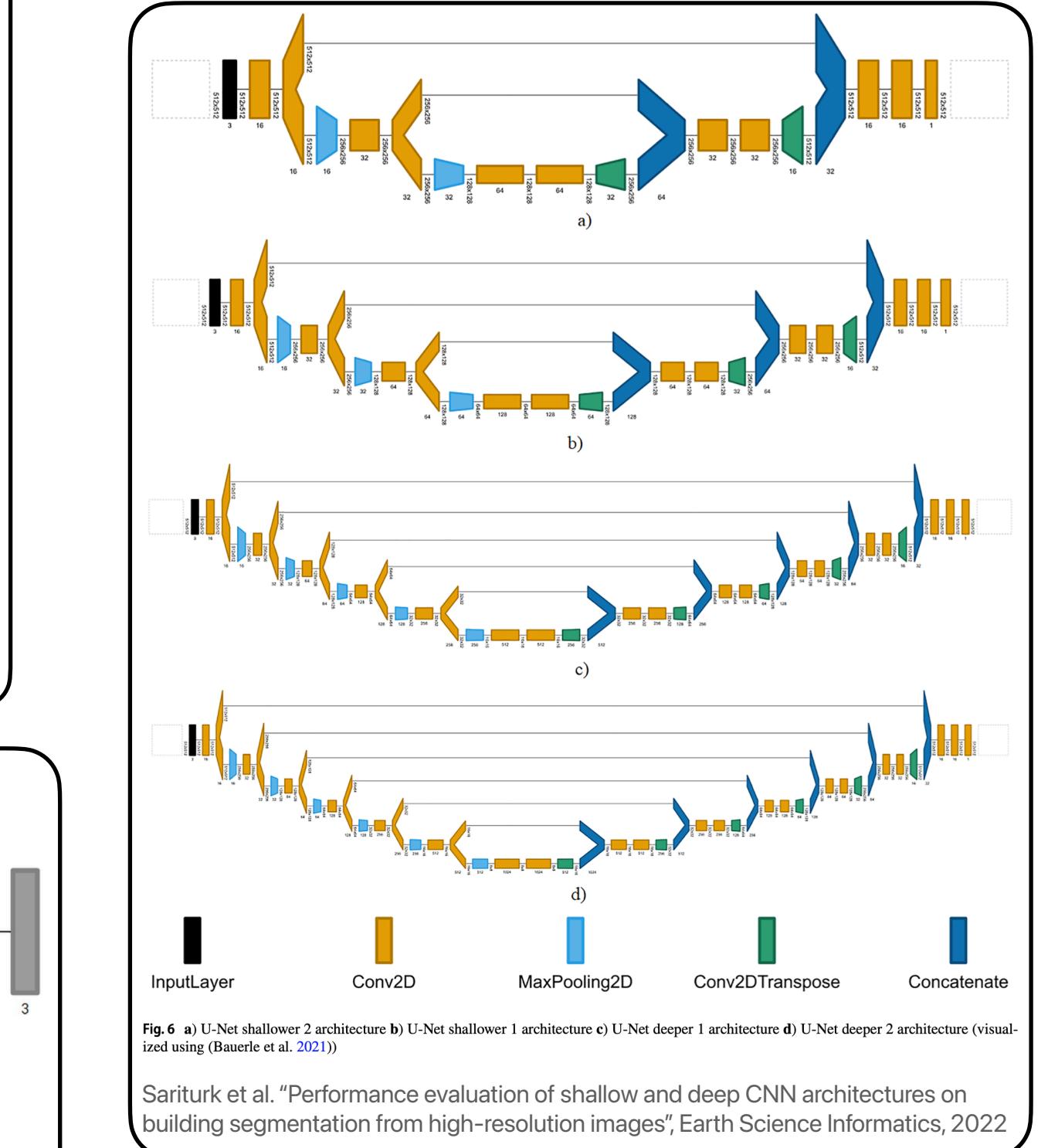
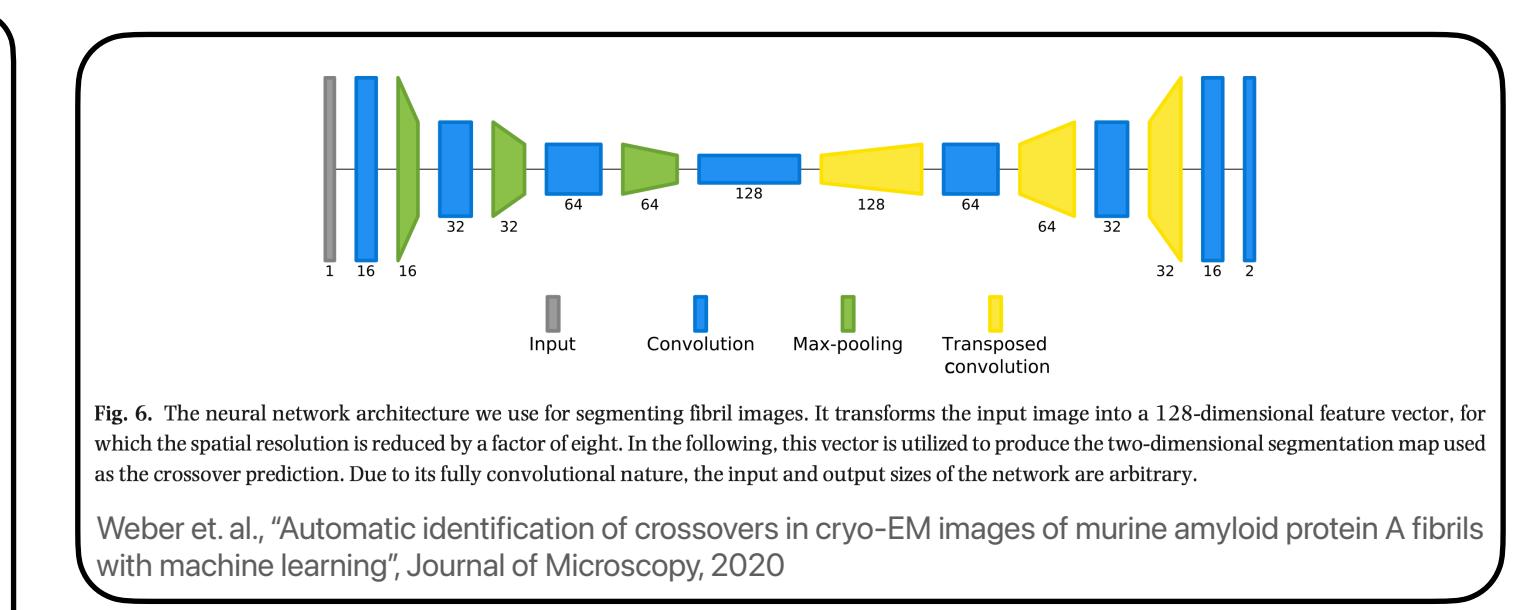
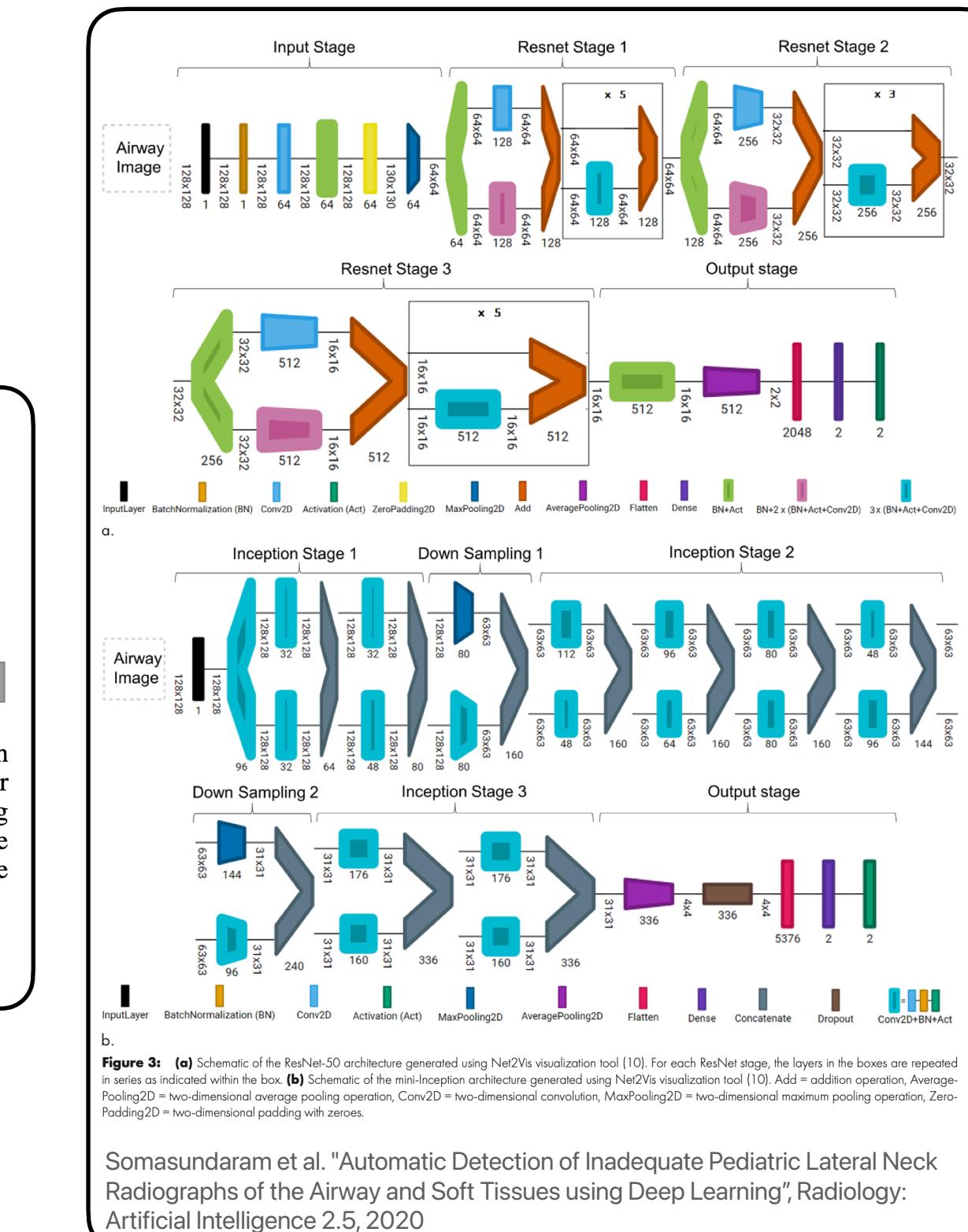
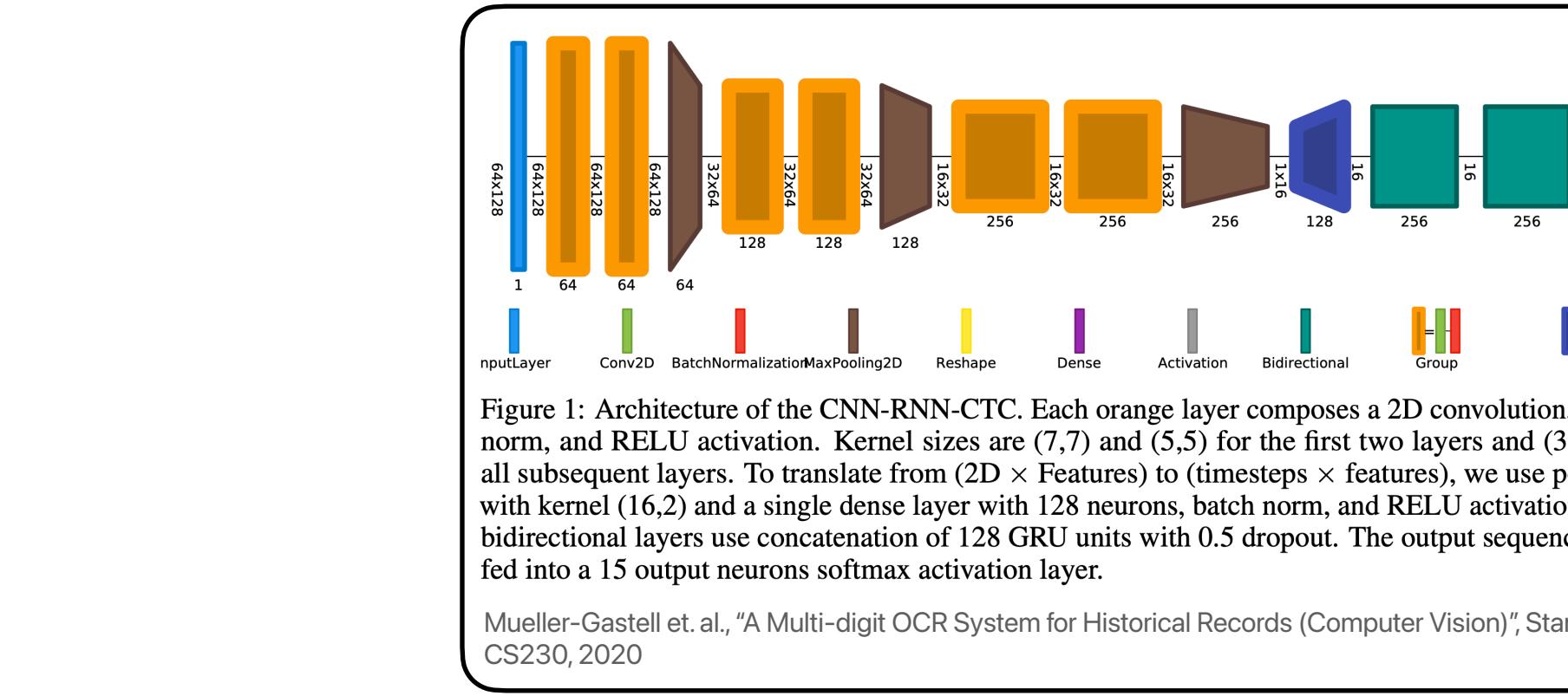


Figure 5: Architecture of proposed encoder-decoder convolutional neural network. It takes an image of 384x384 pixels resolution and process it in several convolutional, pooling, transposed convolutional and concatenation layers, before the final pixelwise semantic segmentation is performed with the "softmax" activation in the last classification layer.

Vernikouskaya et al. "Body fat compartment determination by encoder-decoder convolutional neural network: application to amyotrophic lateral sclerosis", Nature Scientific Reports, 2022

Improving research communication through automatic visualization- generation.

Quality Assurance

Human: deer Computer: airplane



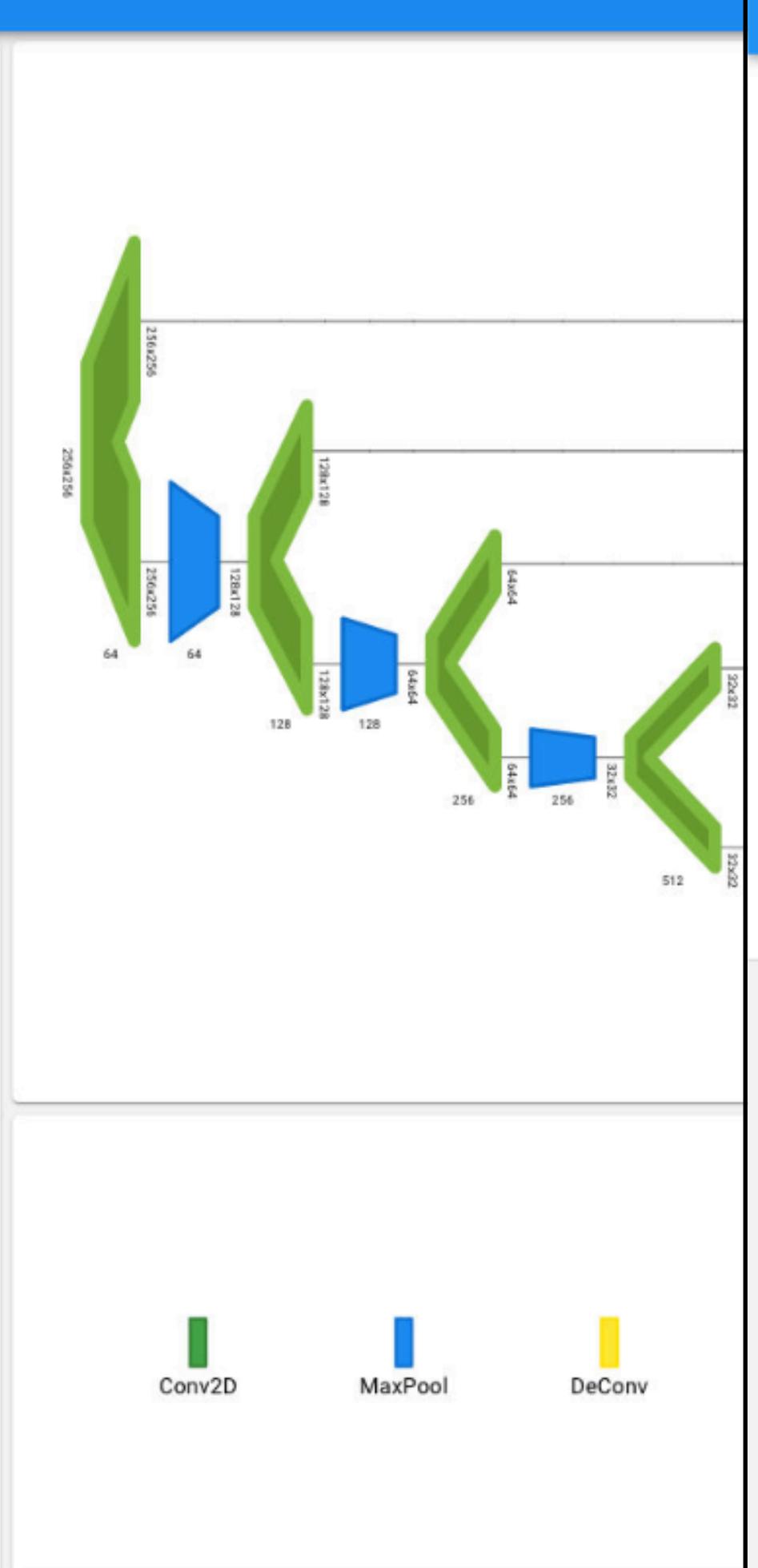
Communication

Net2Vis

Code

Preferences

```
1 # You can freely modify this file.
2 # However, you need to have a function that
3 # is named get_model and returns a Keras
4 # Model.
5
6 import tensorflow as tf
7 from tensorflow.python.keras import models
8 from tensorflow.python.keras import layers
9 from tensorflow.python.keras import utils
10
11 def get_model():
12     img_height = 256
13     img_width = 256
14     img_channels = 1
15
16     input_shape = (img_height, img_width,
17                   img_channels)
18     img_input = tf.keras.Input(shape =
19                               =input_shape)
20     conv1 = layers.Conv2D(64, 3, activation =
21                         = 'relu', padding = 'same',
22                         kernel_initializer = 'he_normal'
23                         )(img_input)
24     conv1 = layers.Conv2D(64, 3, activation =
25                         = 'relu', padding = 'same',
26                         kernel_initializer = 'he_normal'
27                         )(conv1)
28     pool1 = layers.MaxPooling2D(pool_size =
29                         =(2, 2))(conv1)
30     conv2 = layers.Conv2D(128, 3,
31                         activation = 'relu', padding =
32                         = 'same', kernel_initializer =
33                         = 'he_normal')(pool1)
34     conv2 = layers.Conv2D(128, 3,
35                         activation = 'relu', padding =
36                         = 'same', kernel_initializer =
37                         = 'he_normal')(conv2)
38     pool2 = layers.MaxPooling2D(pool_size =
39                         =(2, 2))(conv2)
40     conv3 = layers.Conv2D(256, 3,
41                         activation = 'relu', padding =
42                         = 'same', kernel_initializer =
43                         = 'he_normal')(pool2)
44     conv3 = layers.Conv2D(256, 3,
45                         activation = 'relu', padding =
46                         = 'same', kernel_initializer =
47                         = 'he_normal')(conv3)
48     pool3 = layers.MaxPooling2D(pool_size =
49                         =(2, 2))(conv3)
50     conv4 = layers.Conv2D(512, 3,
51                         activation = 'relu', padding =
52                         = 'same', kernel_initializer =
53                         = 'he_normal')(pool3)
54     conv4 = layers.Conv2D(512, 3,
55                         activation = 'relu', padding =
56                         = 'same', kernel_initializer =
57                         = 'he_normal')(conv4)
58     drop4 = layers.Dropout(0.5)(conv4)
59     pool4 = layers.MaxPooling2D(pool_size =
60                         =(2, 2))(drop4)
61     conv5 = layers.Conv2D(1024, 3,
```



Education

exploRNN

Function Data ▾

Input ⓘ



Network ⓘ



Forward ⓘ

Data is shown to the network value by value to build up the internal state. After a fixed number of data points has been processed, the network can make a prediction on how this sample would continue.

Validation ⓘ

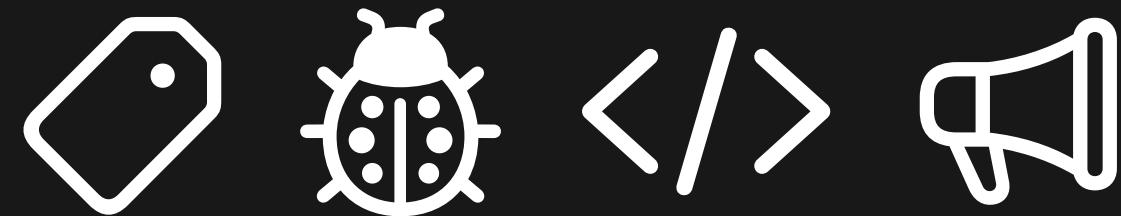
The predicted values are compared to the correct values (ground truth) from the training dataset. The difference is used to calculate the loss.

Backward ⓘ

The calculated loss is backpropagated through the network as well as through time (reverting the input timesteps), to find out where the prediction error came from and update the network variables for the next iteration.

How Are ML Interfaces Used Today?

Interviews with 9 Apple ML practitioners



What data reporting tools do you use?

What tools are you lacking?

Do you use visualizations?

Where do you store your data?

What would your ideal workflow be?

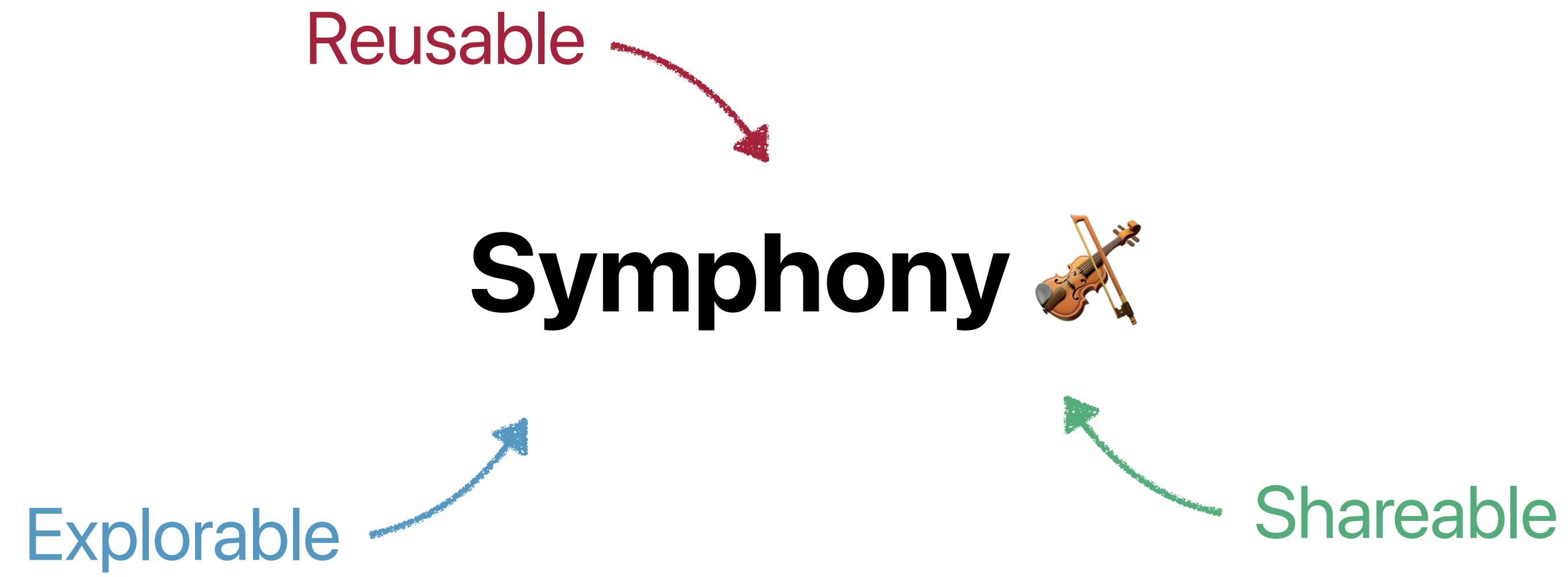
How are ML interfaces used today?

Ad-hoc tools
and analyses

Limitations of
existing ML interfaces

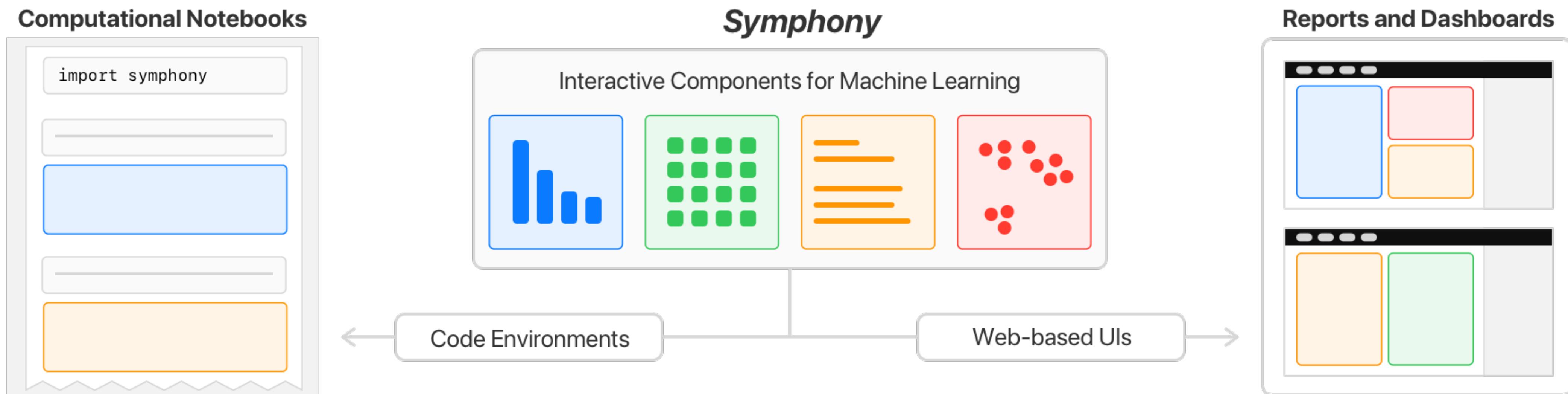
Lack of communication
between stakeholders

ML interface framework with **reusable**,
explorable, and **shareable** visualizations.



Symphony 🎻

Modular, Interactive, and Shareable Components



File Edit View Insert Cell Kernel Widgets Help

Trusted Python 3 (ipykernel) ○

Run Cell Markdown Cell Kernel Help

Cifar 10 Symphony

```
In [1]: import pandas as pd

from symphony import Symphony
from symphony_list import SymphonyList
from symphony_summary import SymphonySummary
from symphony_familiarity import SymphonyFamiliarity
from symphony_scatterplot import SymphonyScatterplot
from symphony_duplicates import SymphonyDuplicates
from symphony_hierarchical_confusion_matrix import SymphonyHierarchicalConfusionMatrix
from symphony_fairvis import SymphonyFairVis
from symphony_markdown import SymphonyMarkdown
```

```
In [2]: DATA_PATH = 'cifar/'
df = pd.read_parquet('table.parquet')
```

```
In [3]: df.head()
```

```
Out[3]:
```

	id	label	dataset	output	duplicates_embedding	projection_embedding_x	projection_embedding_y	familiarity_embedding	splitFamiliarity_embedding_byAttr_label	splitFamiliarity_embedding
0	train/ship/image194.png	ship	train	ship	-1	0.409252	9.471022	24.430313	{'airplane': 12.975905777838438, 'automobile': ...}	{'test': 21.55013580}
1	train/deer/image610.png	deer	train	deer	-1	15.026260	4.096645	5.905596	{'airplane': -29.83306951371854, 'automobile': ...}	{'test': 6.171497820}
2	train/truck/image2259.png	truck	train	truck	-1	-0.561973	5.591904	19.569766	{'airplane': 10.05607537711247, 'automobile': ...}	{'test': 19.87549481}
3	train/cat/image2898.png	cat	train	cat	-1	6.223715	0.279402	16.975622	{'airplane': -31.109959897103085, 'automobile': ...}	{'test': 10.82361368}
4	test/dog/image447.png	dog	test	horse	-1	8.475064	3.314383	19.254085	{'airplane': 12.839210831518892, 'automobile': ...}	{'test': 25.41012778}

Create and Show Visualizations

We can now explore the individual report widgets!

```
In [ ]: report = Symphony(df, files_path='cifar/')
```

```
In [ ]: report.widget(SymphonyMarkdown, page="Overview", width="M", content=open("README.md").read())
```

```
In [ ]: report.widget(SymphonySummary, page="Overview", width="M")
```

```
In [ ]: report.widget(SymphonyList, page="Overview")
```

```
In [ ]: report.widget(SymphonyDuplicates, page="Data Analysis", width="M", height="L")
```

Description ? Edit Download**Overview**

The CIFAR-10 dataset consists of 60000 32x32 colour images in 10 classes, with 6000 images per class. There are 50000 training images and 10000 test images.

Dataset Collection

The original tiny images dataset was collected by researchers at MIT and NYU. The researchers used keywords from the WordNet database to search for images on search platforms like Google, Flickr, and Altavista. They then removed perfect duplicates and images with a significant amount of white pixels that tended to be synthetic images. The images were then downsampled to 32x32 pixels.

To create the CIFAR-10 dataset, researchers at the University of Toronto took a subset of the tiny images dataset and labeled them across 10 object classes. The dataset has a total of 60,000 images with 6,000 images in each of the 10 classes.

Dataset Labeling

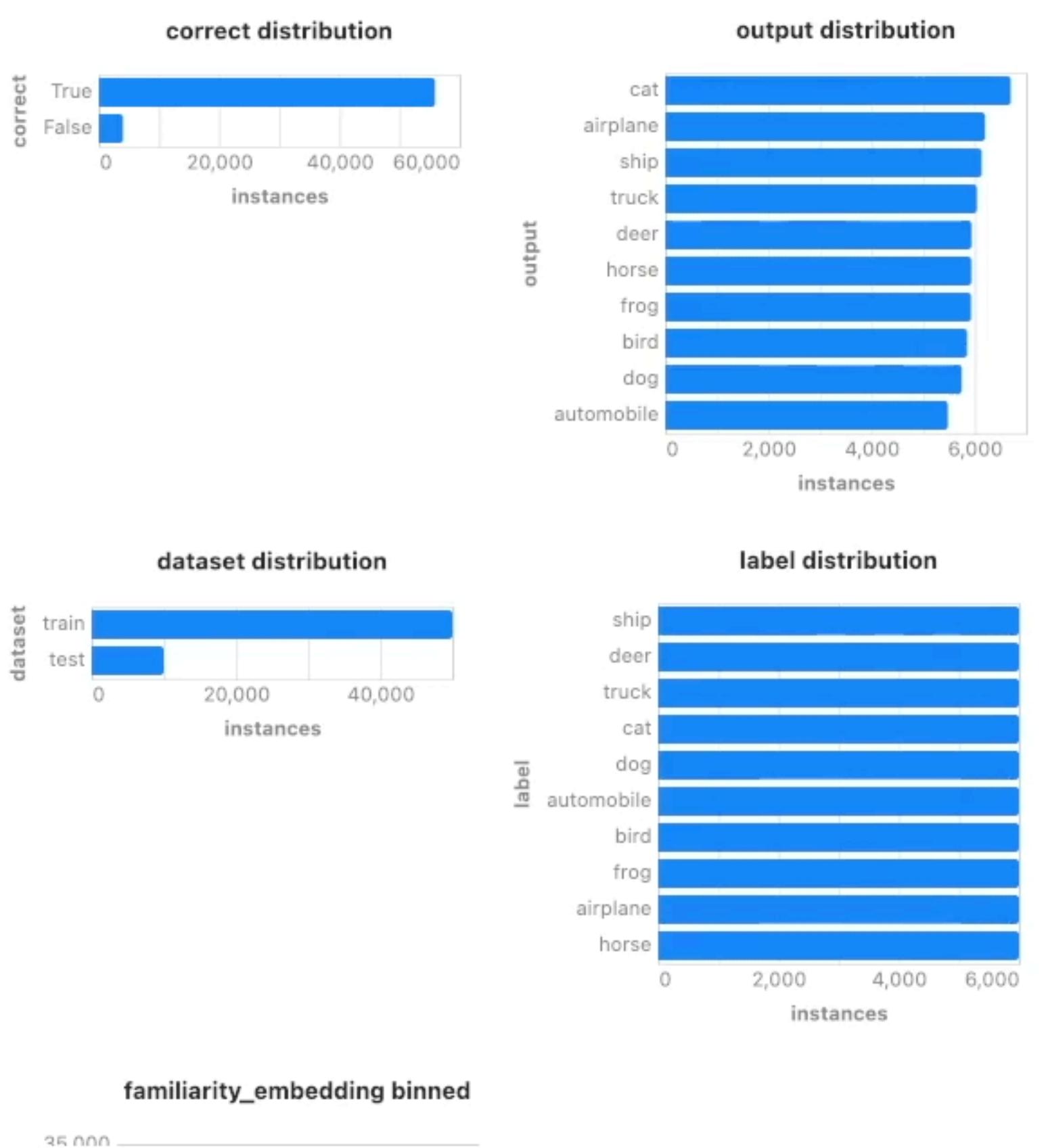
The dataset was labeled by paid students. The 10 labels are airplanes, cars, birds, cats, deer, dogs, frogs, horses, ships, and trucks.

Accessing the Dataset

The dataset can be downloaded from [the website](#) of its original authors.

Summary Statistics ? Edit

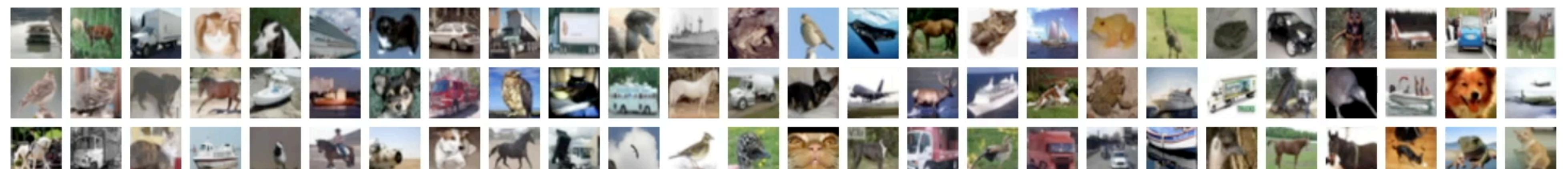
of Instances: 60,000 duplicates_embedding: 1,434

 Download**Settings**

Show unfiltered charts

Samples per page

150

Filter Apply**Group**X Group**List View** ? Download**Instance List**

Duplicates

Download

Candidate Group 0



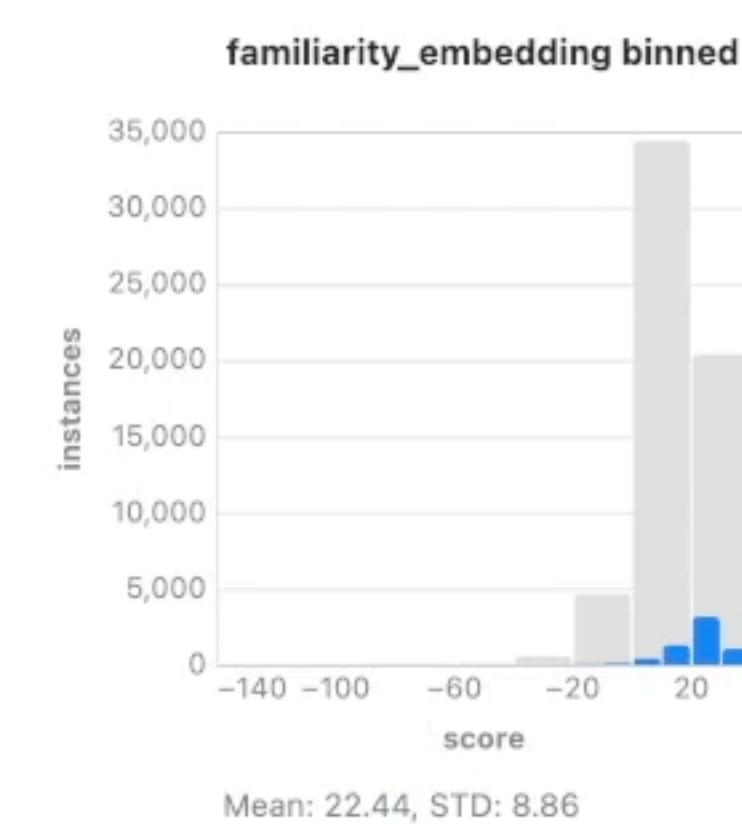
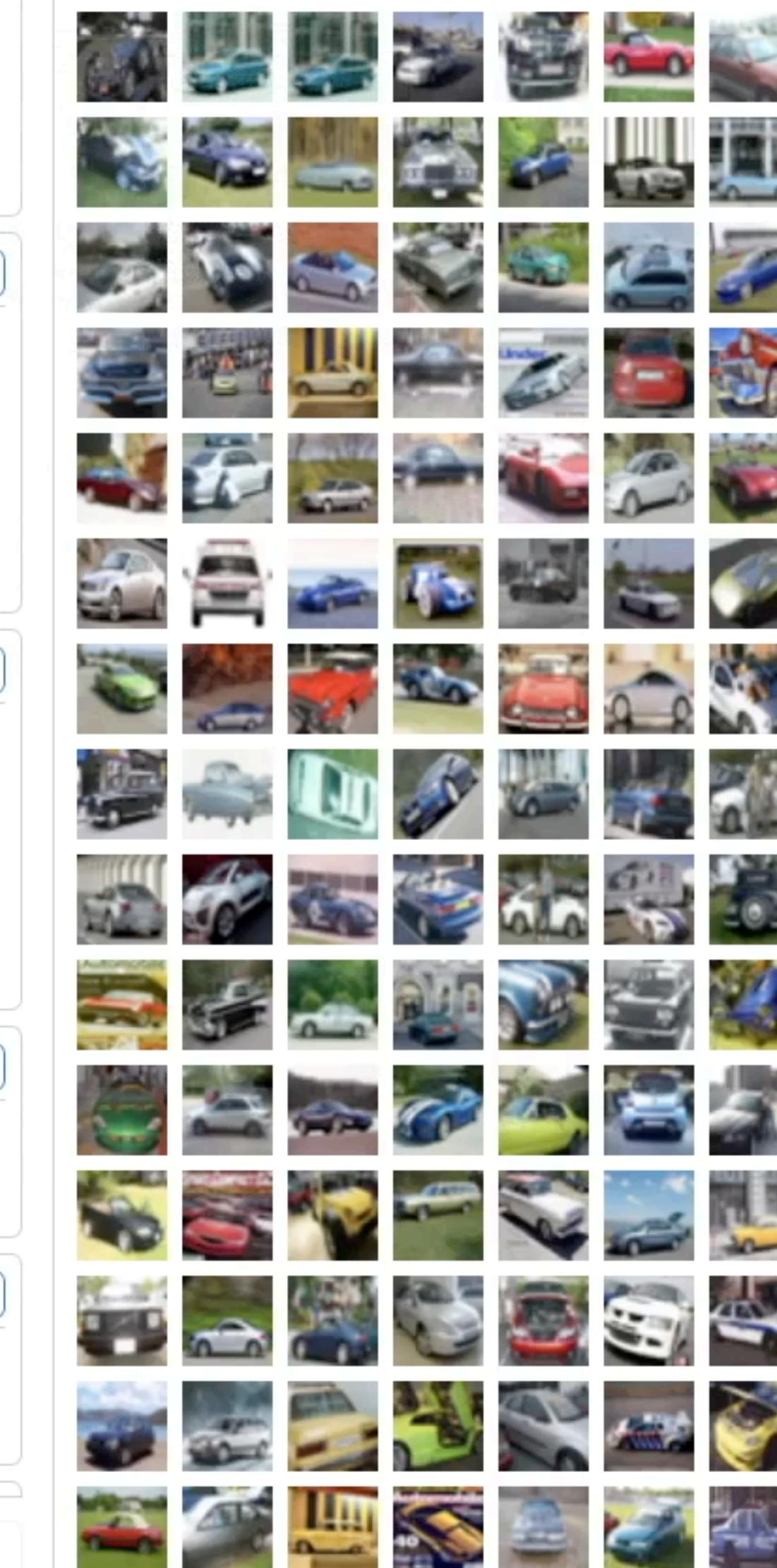
Show 15 more

Familiarity By Attribute

Download

Unfamiliar to Familiar

Least Familiar Instances



Candidate Group 1



Show 8 more

Candidate Group 2



Show 2 more

Candidate Group 3



Candidate Group 4



Settings

Show unfiltered charts

Samples per page

150



Filter

`d.label == 'automobile'`

Clear

Apply

Group



Hierarchical Confusion Matrix ⓘ[Download](#)

Dimensions

root ↕

Shelf Enable and disable different dimensions of the data. The order of dimension defines the nesting level.

All dimensions are already in use.

Where Condition the confusion matrix on the value of a given label.

Hover over cells to show more information.

Counts**Observed**

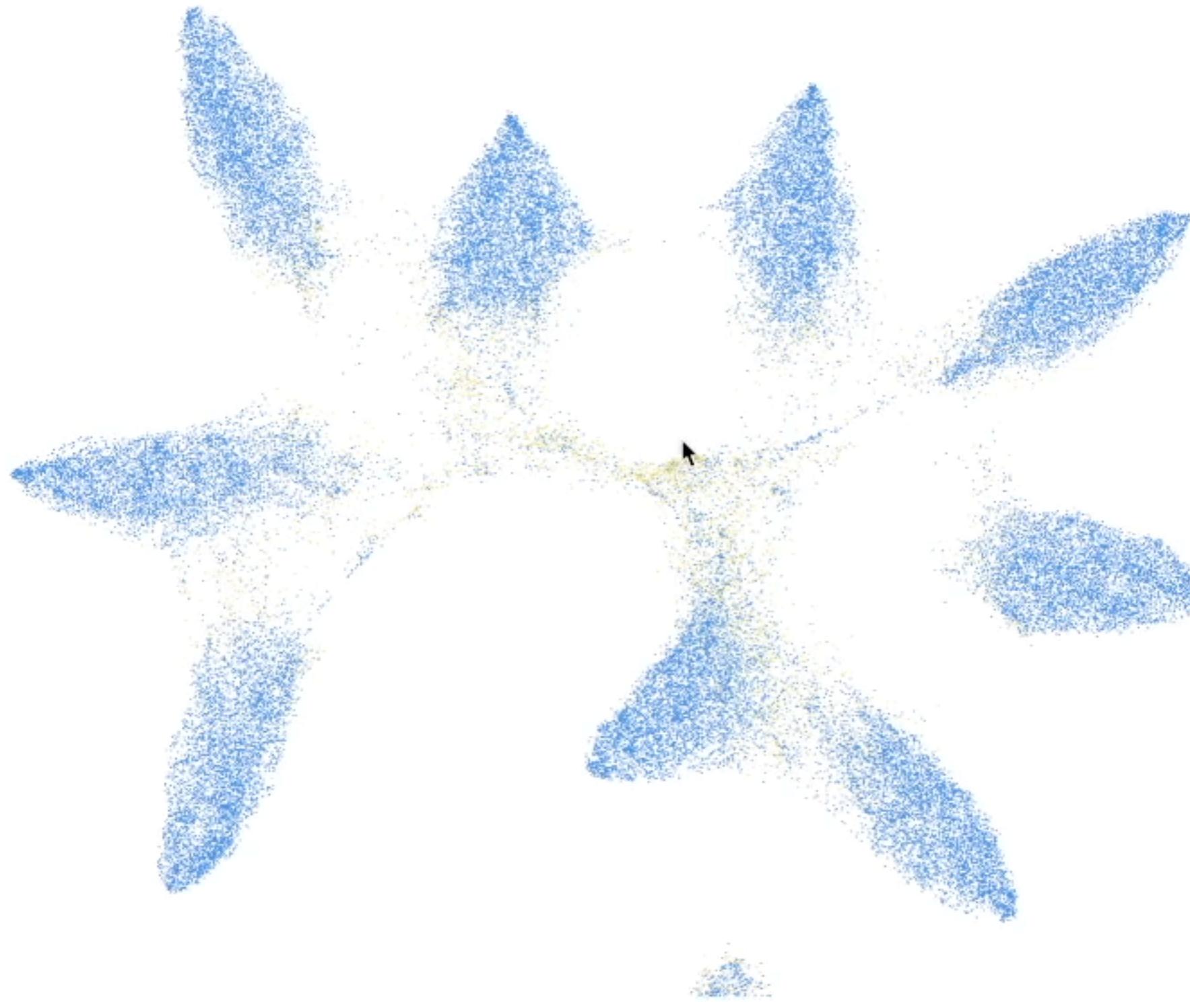
Actual

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

	Precision		Recall	Accuracy
	Actual	Predicted		
✓ root	0.93	0.93	0.87	
ship	0.94	0.97	0.99	
deer	0.95	0.94	0.99	
truck	0.93	0.94	0.99	
cat	0.83	0.92	0.97	
dog	0.93	0.89	0.98	
automobile	0.98	0.90	0.99	
bird	0.94	0.92	0.98	
frog	0.96	0.95	0.99	
airplane	0.92	0.95	0.99	
horse	0.95	0.94	0.99	

Scatterplot ⓘ Category column: correct ⏺[Download](#)

Double-click to recenter. Shift-click and drag to lasso-select.

**FairVis** ⓘ

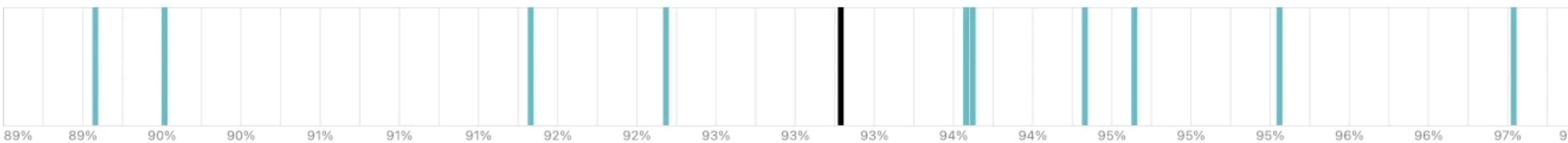
Label column: label

Prediction column: output

Plot type: Strip Plot

[Download](#)

Accuracy

**Settings**

Show unfiltered charts



Samples per page

150

Filter

|

[Apply](#)**Group**

label

[Clear](#)[Group](#)

ship (6,000) ×

deer (6,000) ×

truck (6,000) ×

cat (6,000) ×

dog (6,000) ×

automobile (6,000) ×

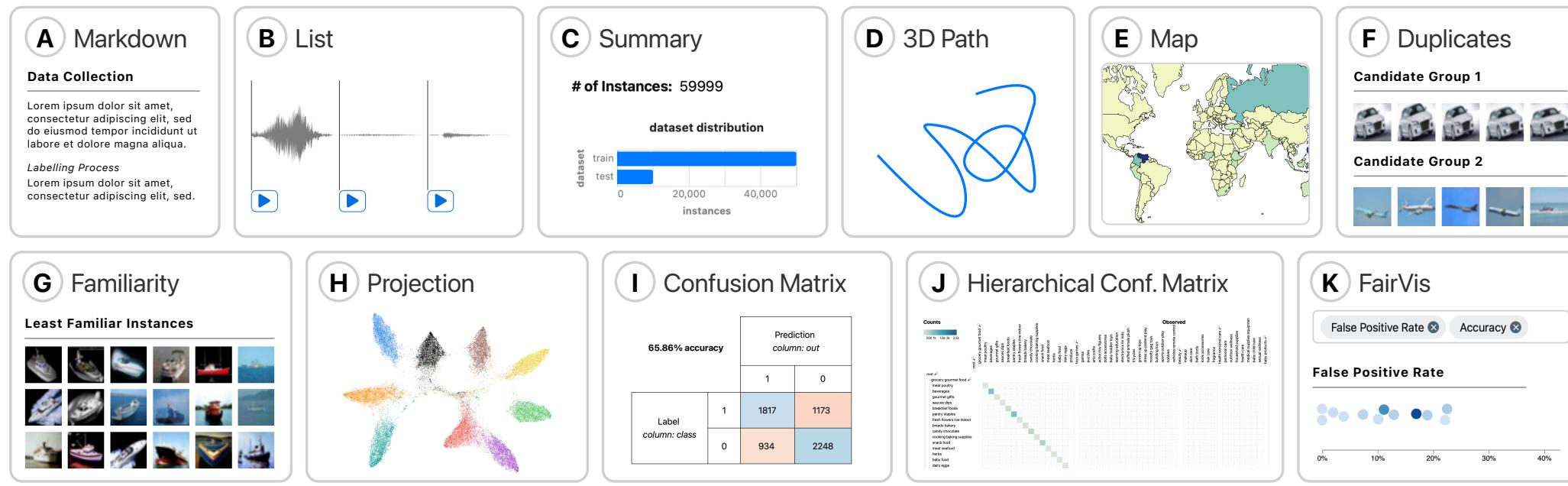
bird (6,000) ×

frog (6,000) ×

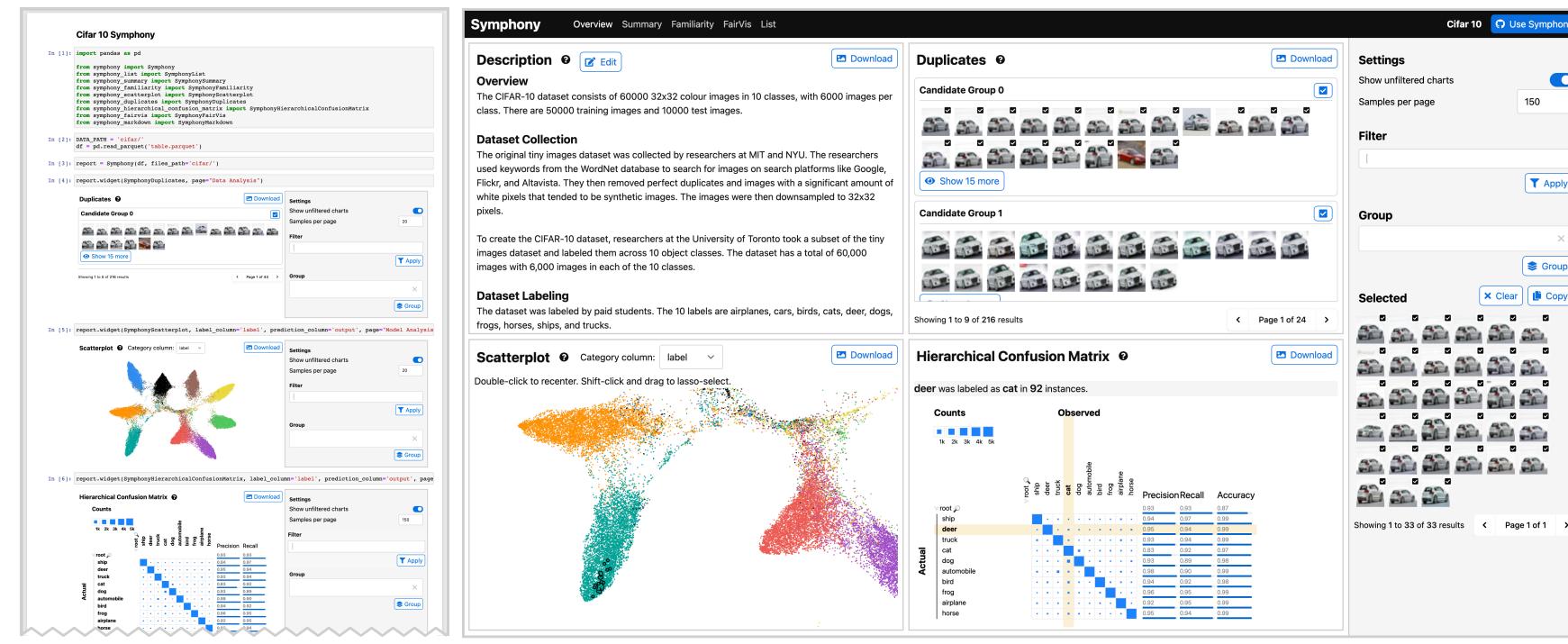
airplane (6,000) ×

horse (6,000) ×

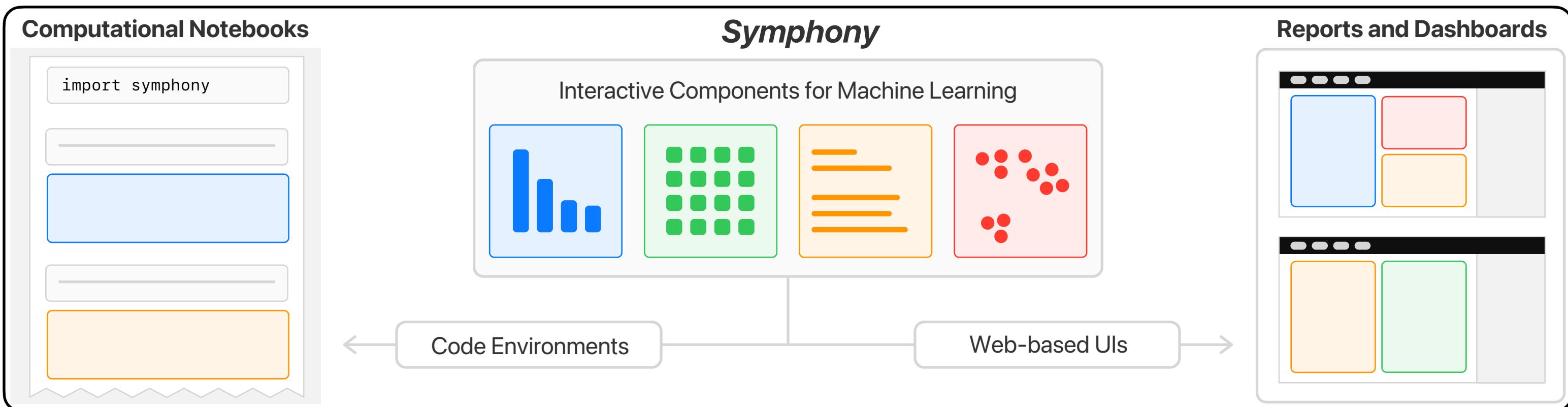
Reusable visualization components



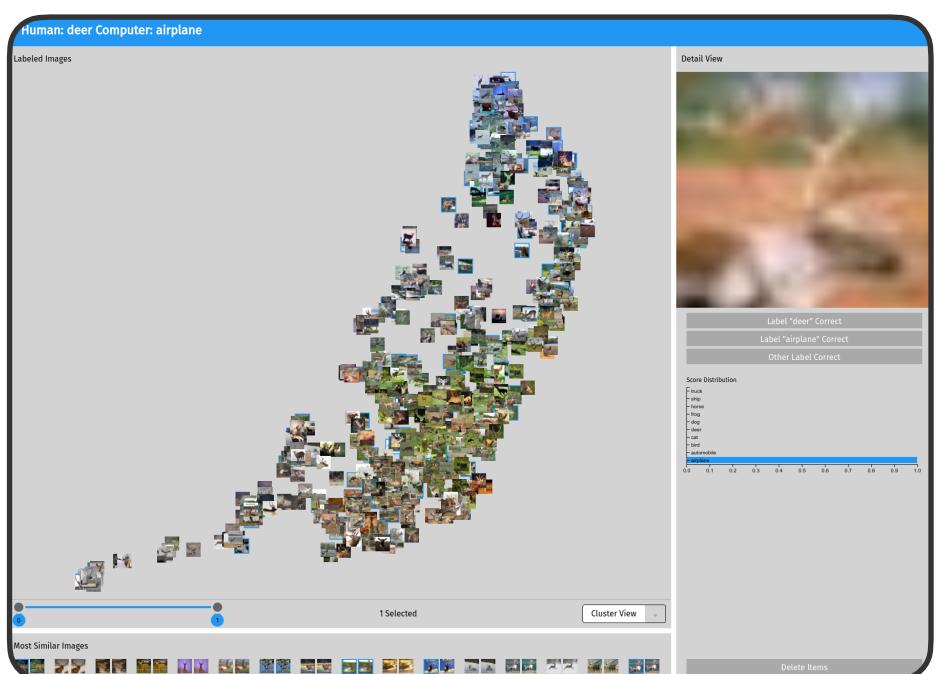
Across multiple environments



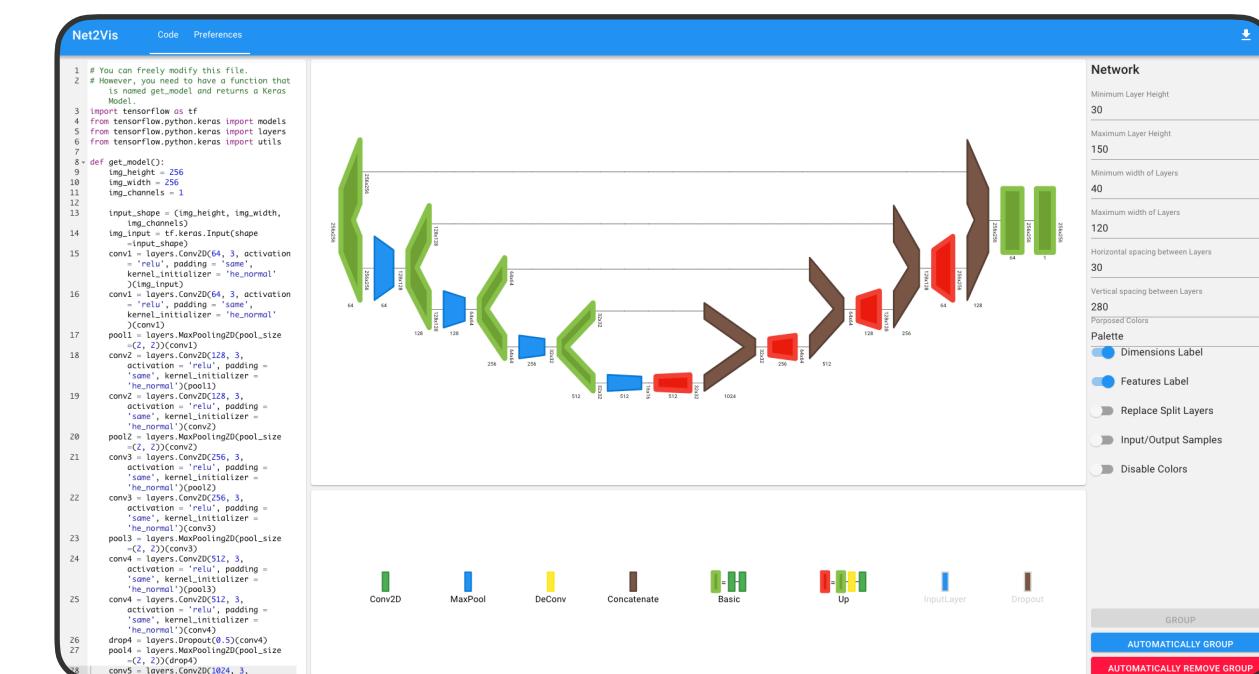
Shared organizational understanding of data and models encourages the creation of accurate, responsible, and robust AI products.



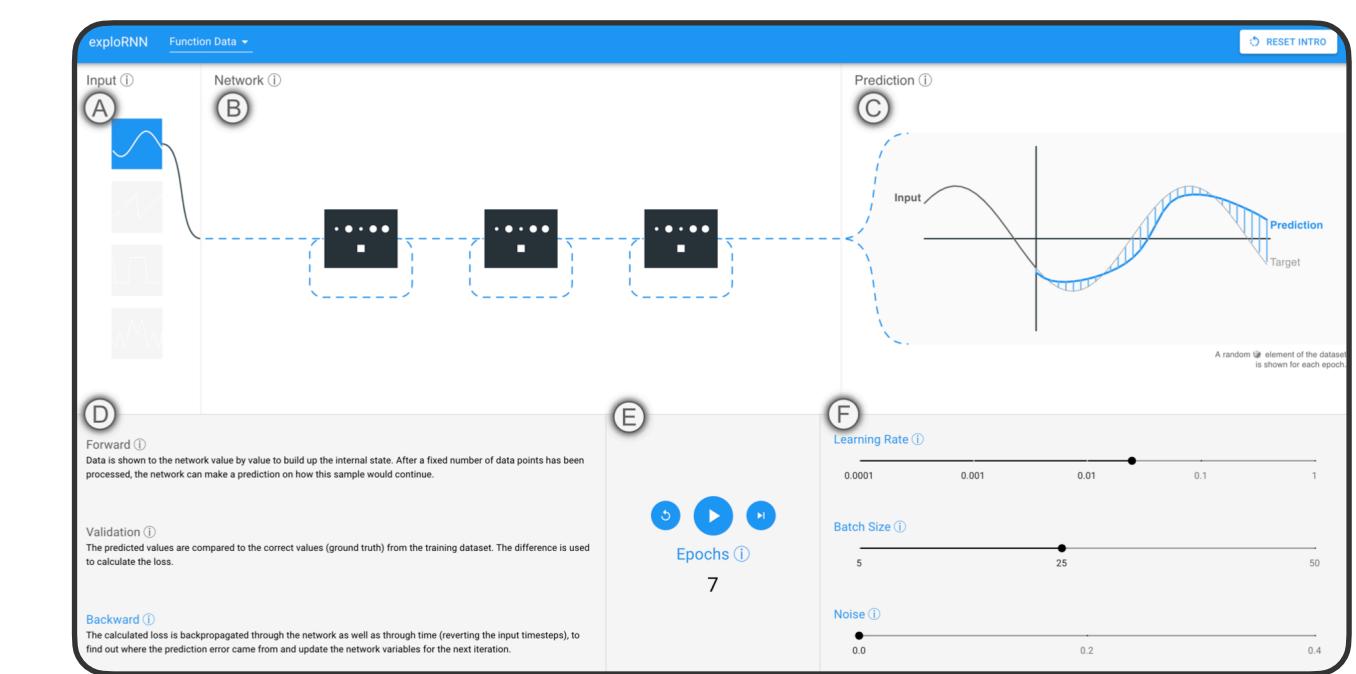
Bäuerle and Cabrera et al. 2022, CHI
Symphony: Composing Interactive Interfaces for Machine Learning



Bäuerle et al. 2020, CGF
Classifier-Guided Visual Correction of Noisy Labels for Image Classification Tasks



Bäuerle et al. 2021, TVCG
Net2Vis - A Visual Grammar for Automatically Generating Publication-Tailored CNN Architecture Visualizations



Bäuerle et al. 2022, TVC
exploRNN: Understanding Recurrent Neural Networks through Visual Exploration

Introspection Technique

Visualization Interface

1

Bäuerle et al. 2020, CGF

Classifier-Guided Visual Correction of Noisy Labels for Image Classification Tasks

2

Bäuerle et al. 2021, TVCG

Net2Vis - A Visual Grammar for Automatically Generating Publication-Tailored CNN Architecture Visualizations

3

Bäuerle et al. 2022, TVC

exploRNN: Understanding Recurrent Neural Networks through Visual Exploration

4

Bäuerle and Cabrera et al. 2022, CHI

Symphony: Composing Interactive Interfaces for Machine Learning

Introspection Technique

Visualization Interface

1

Bäuerle et al. 2020, CGF
Classifier-Guided Visual Correction of Noisy Labels for Image Classification Tasks

Bäuerle et al. 2021, TVCG
Net2Vis - A Visual Grammar for Automatically Generating Publication-Tailored CNN Architecture Visualizations

2

Bäuerle et al. 2022, TVC
exploRNN: Understanding Recurrent Neural Networks through Visual Exploration

3



Bäuerle and Cabrera et al. 2022, CHI
Symphony: Composing Interactive Interfaces for Machine Learning

4

5 Aka et al. 2021, AIES
Measuring Model Biases in the Absence of Ground Truth 

 Bäuerle et al. 2022, arXiv
Visual Identification of Problematic Bias in Large Label Spaces

6

7 Bäuerle and Wexler 2020, VISxAI 
What does BERT dream of?

8 Weber et. al. 2020, Journal of Microscopy
Automatic identification of crossovers in cryo-EM images of murine amyloid protein A fibrils with machine learning

10 Bäuerle and Jönsson et. al. 2022, arXiv
Neural Activation Patterns (NAPs): Visual Interpretability of Learned Concepts

Bäuerle and van Onzenoodt et. al. 2022, CGF
Where did my Lines go? Visualizing Missing Data in Parallel Coordinates

9

Model attributes

Modern architectures

Introspection

Develop methods for ML introspection and make them accessible through visualization.

Visualization

Data types, tasks & users

Beyond research projects



Visualization-Based Neural Network Introspection

Alex Bäuerle

Defense, 21.12.2022