

Visual Analytics of the Machine Learning Process

Mengchen Liu, Tsinghua University

Supervisor: Shixia Liu



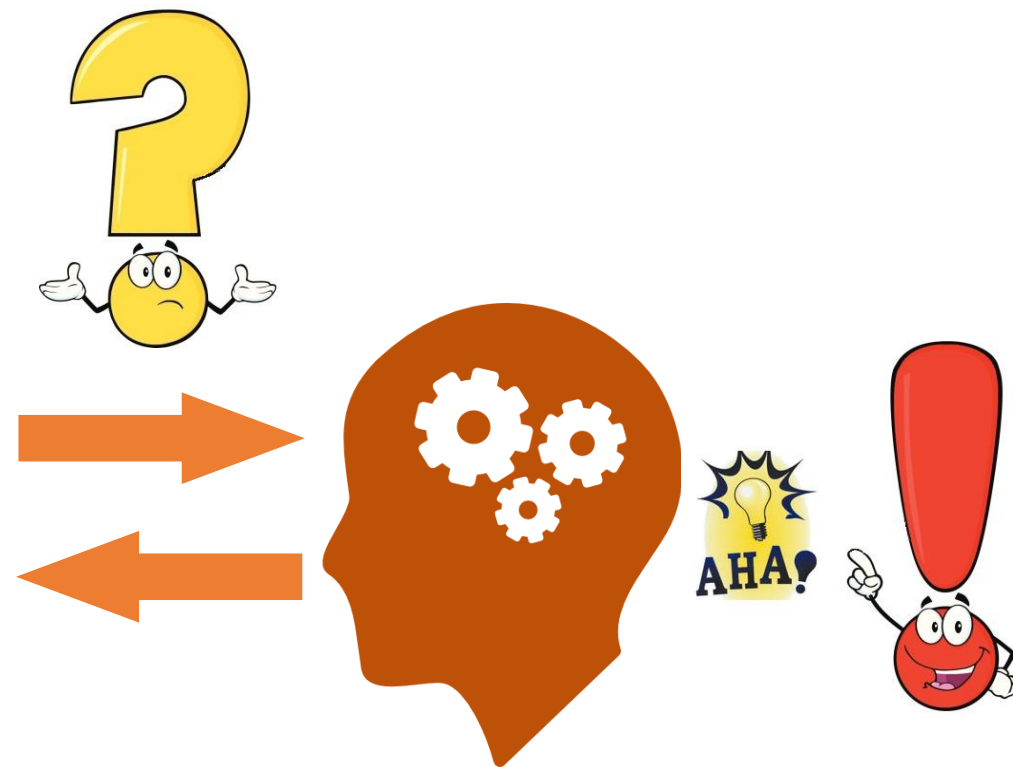
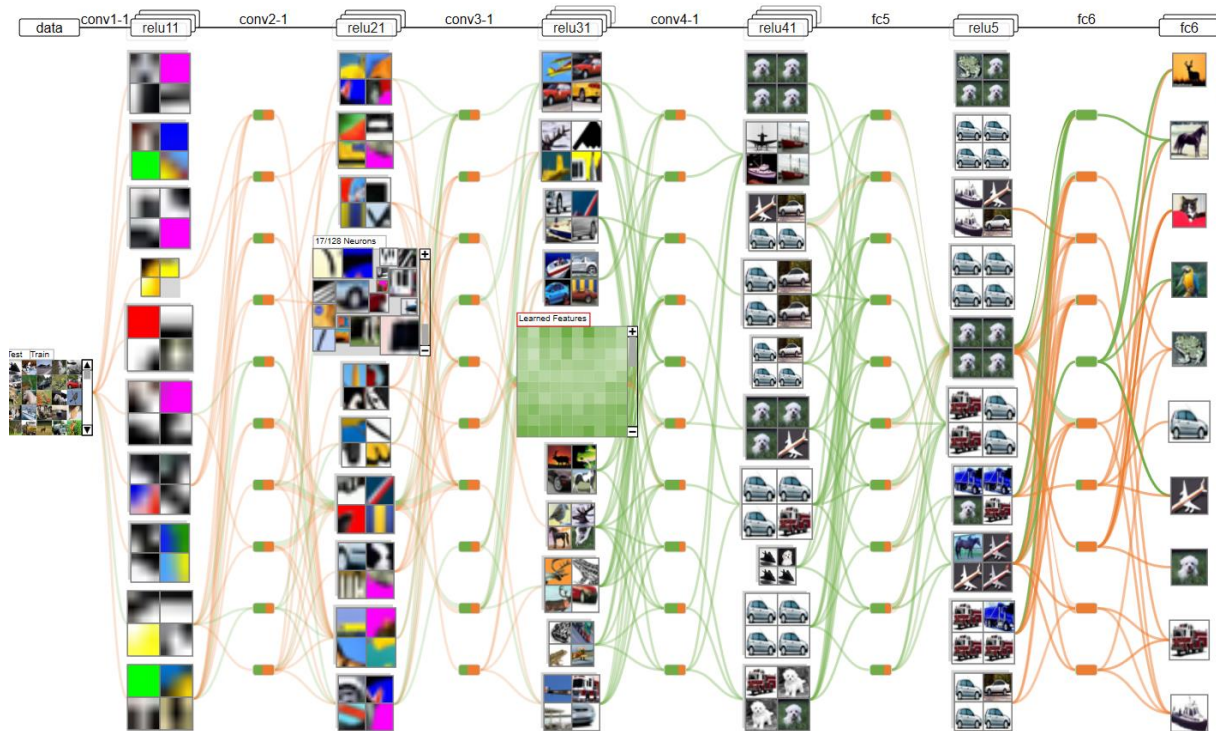
Machine learning is hot



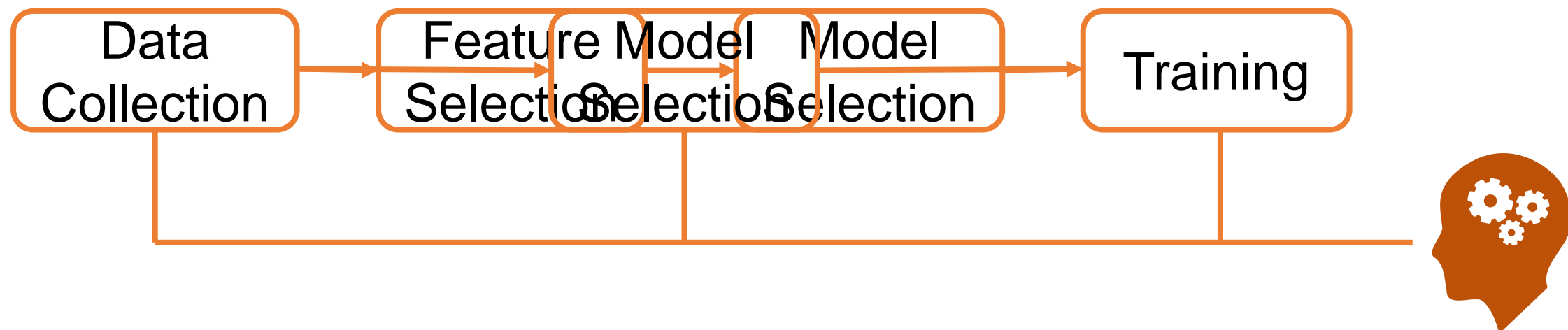
Machine Learning is Hard

- A large amount of data is needed
 - MSCOCO: 330K images
- Machine learning models are often treated as “black boxes”
 - Dozens of layers, millions of parameters
- Successful training a machine learning model needs time, skill.
 - Trial-and-error
 - Single trial
 - AlphaGo (2016): 1 week on 50 GPUs using Google Cloud
-

Machine Learning + Visualization

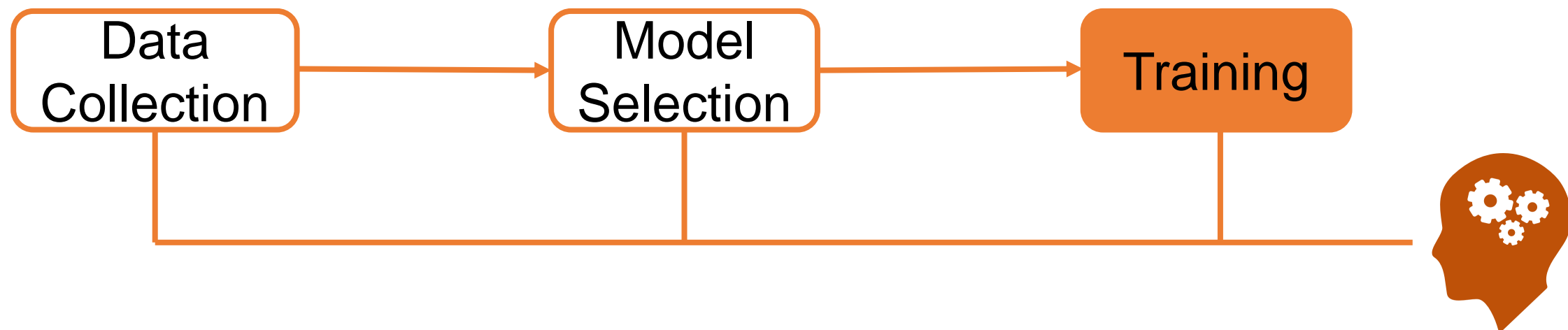


Machine Learning + Visualization



Deep Learning. Goodfellow et al.,
Pattern Classification. Duda et al.

Machine Learning + Visualization



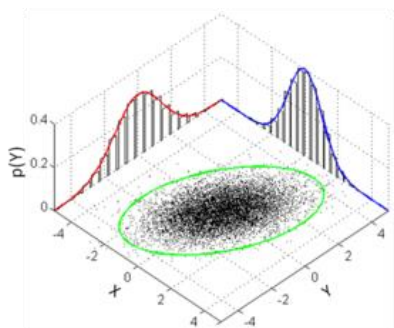
Deep Learning. Goodfellow et al.,
Pattern Classification. Duda et al.

Analyzing the Training Processes of Deep Generative Networks

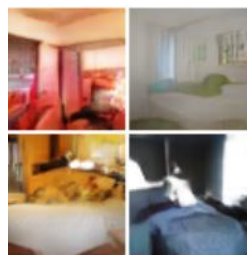
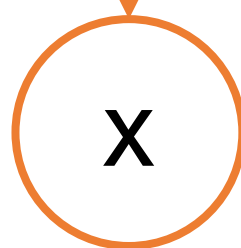
Mengchen Liu, Jiaxin Shi, Kelei Cao, Jun Zhu, Shixia Liu

TVCG 2018

Deep Generative Model (DGM)



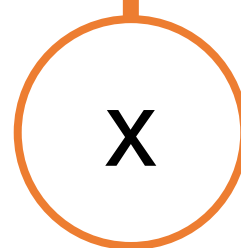
Deep neural network



Unsupervised /
semi-supervised

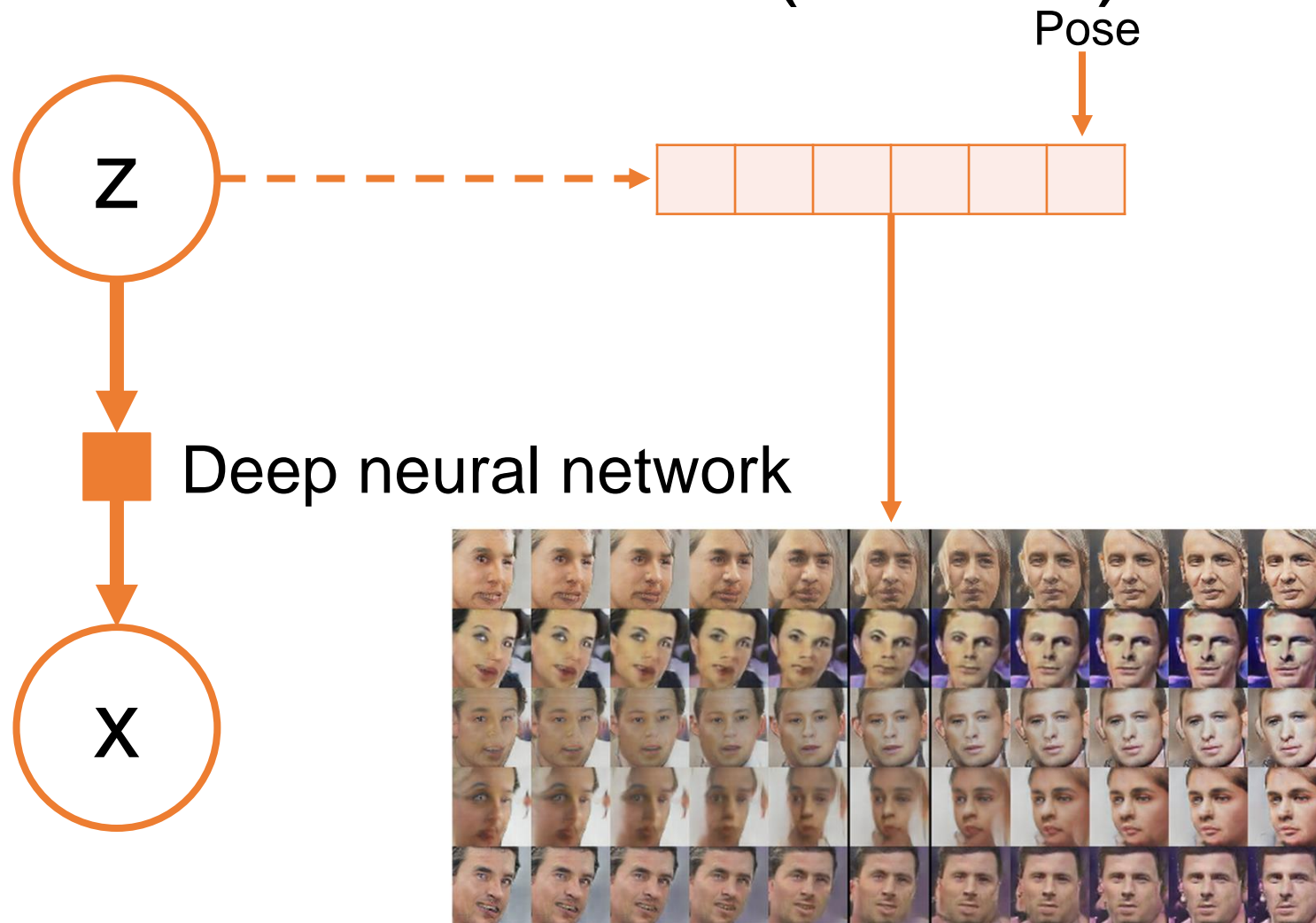


Smiling face



Supervised

Deep Generative Models (DGMs)

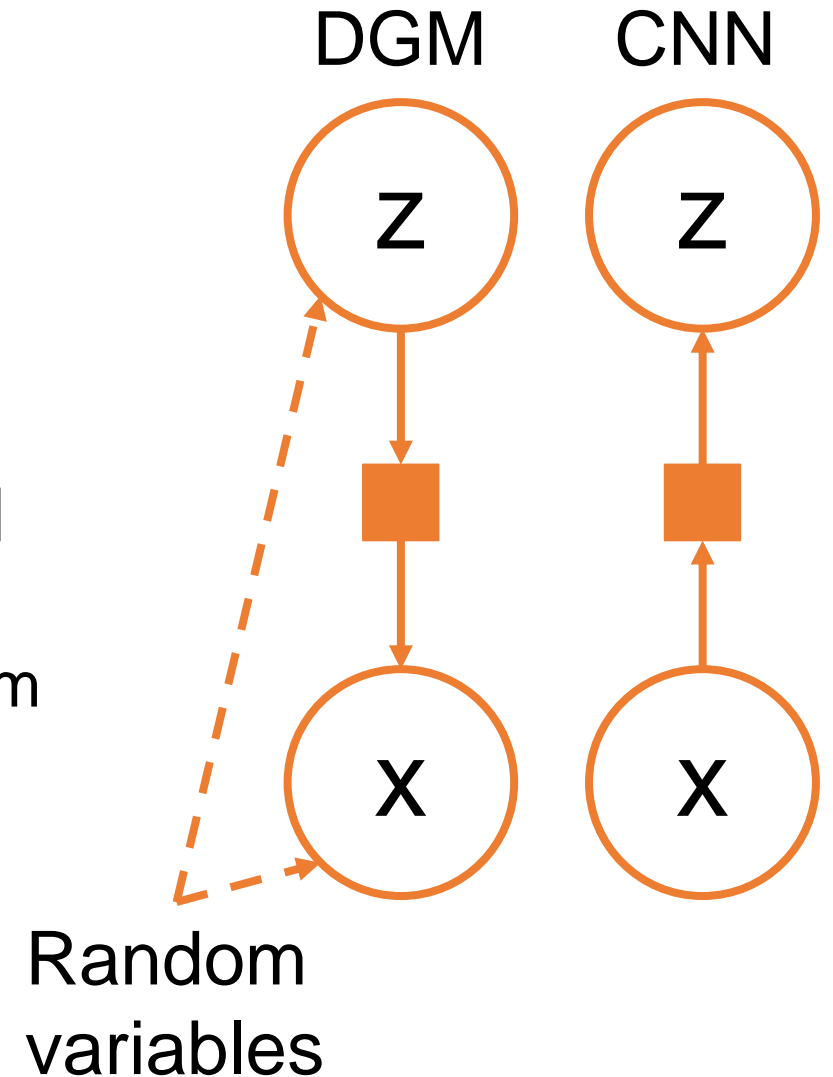


By generative adversarial network (GAN)

Training a DGM is Hard

- DGM: both deterministic functions and random variables (x, z)
 - Convolutional neural network (CNN): deterministic functions (e.g., convolution and pooling)
- DGM: a top-down generative process and a bottom-up Bayesian inference process
 - CNN: a bottom-up process: input at the bottom layer \rightarrow high-level features \rightarrow outputs

Visualization



Challenge 1

- Handle a large amount of time series data
 - Typical time series data: activation/gradient/weight changes over time
 - Millions of activations/gradients/weights in a DGM

Contribution 1

- Handle a large amount of time series data
 - Activation/gradient/weight changes over time
 - Millions of activations/gradients/weights in a DGM
- Line chart + **blue noise polyline sampling** algorithm
 - Each time series as a polyline
 - Select polyline samples with blue-noise properties
 - Preserve outliers and reduce visual clutter

Challenge 2

- Identify the root cause of a failed training process
 - It's often difficult to locate the specific neurons leading to the training failure
 - Neurons influence each other

Contribution 2

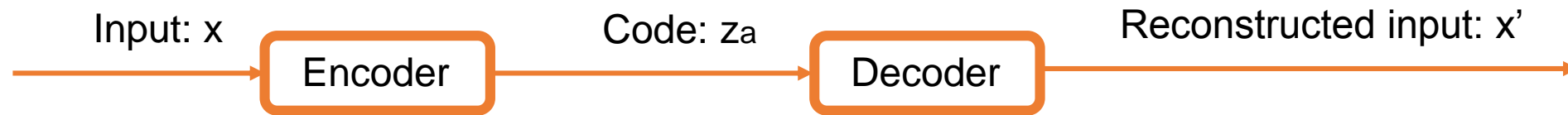
- Identify the root cause of a failed training process
 - It's often difficult to locate the specific neurons leading to the training failure
 - Neurons influence each other
- A **credit assignment** algorithm
 - Explain how other neurons contribute to the output of the neuron causing a training problem

DGMTracker

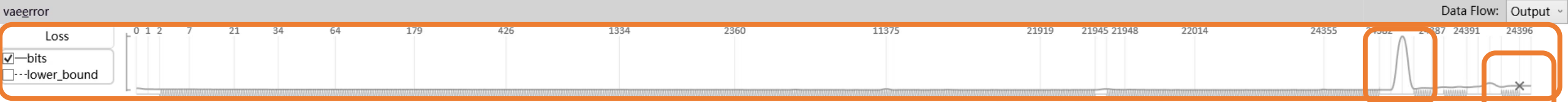
- Better understand and diagnose the training process of a DGM

Case Study: Debugging a Failed Training Process of a Variational Autoencoder (VAE)

- Autoencoder
 - Reconstruct their input with minimum information loss

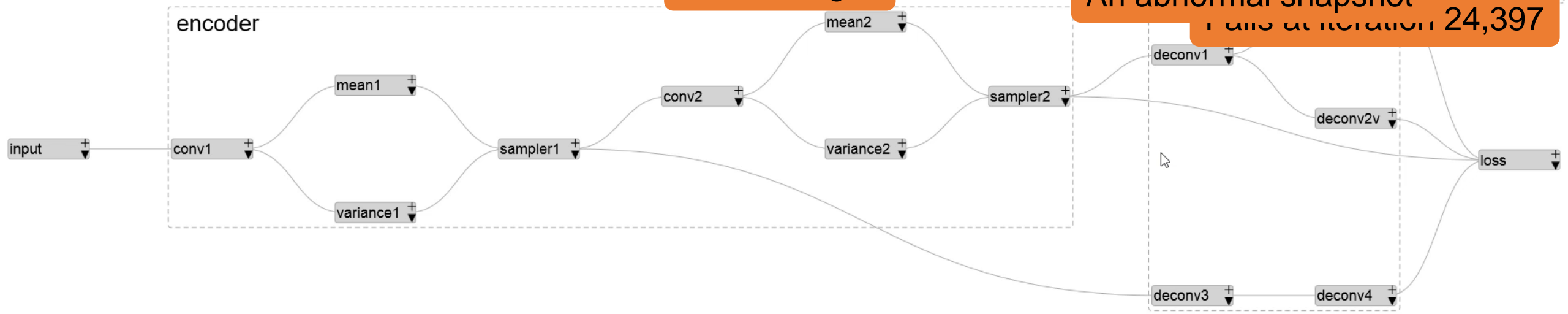


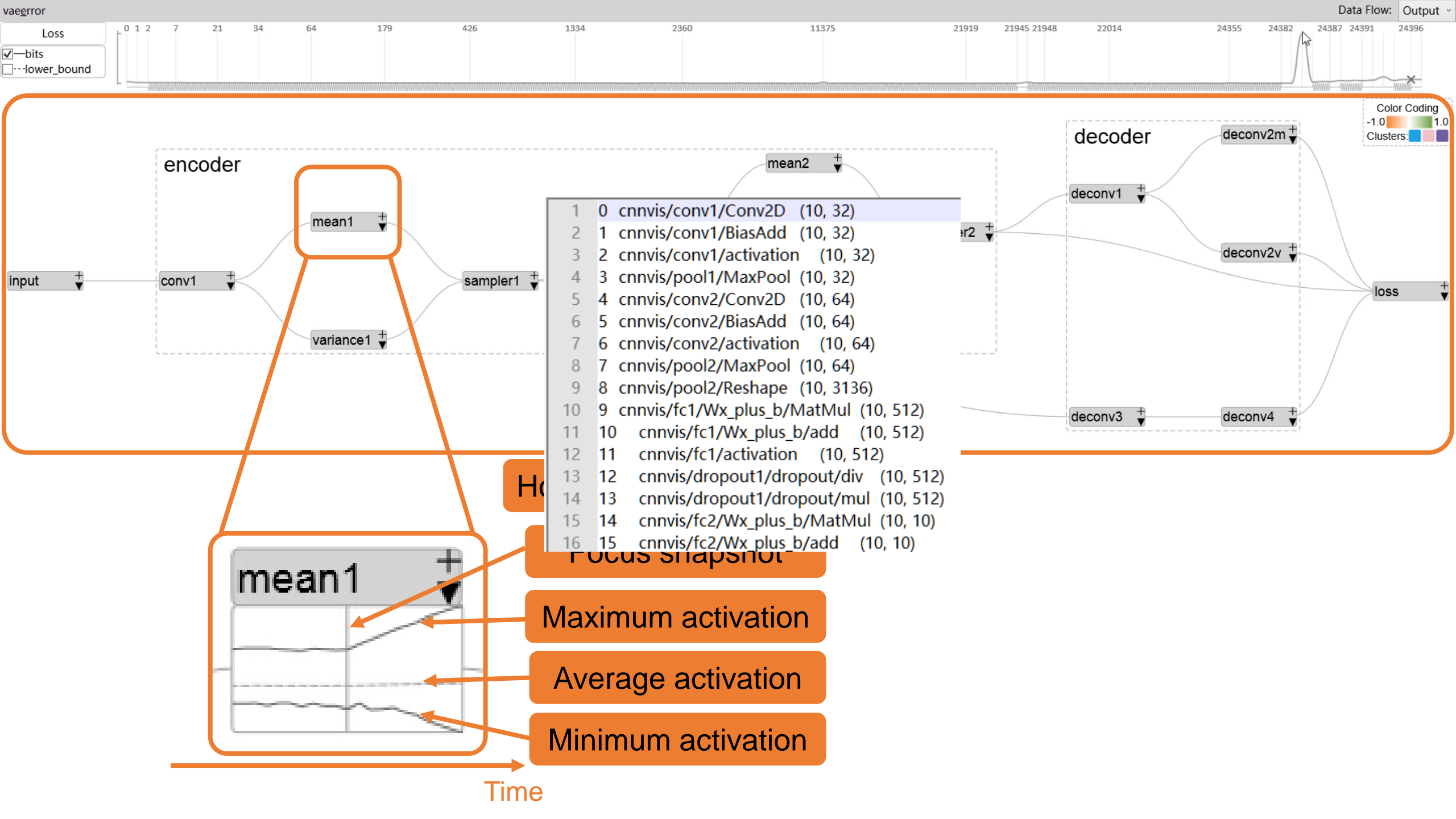
- Variational autoencoder
 - Probabilistic version of an autoencoder
 - z_v : a vector of random variables
 - z_a : a vector of real numbers
- Dataset: CIFAR10 dataset
Loss = NaN (10k-30k iterations)
An example case: fails at 24,397

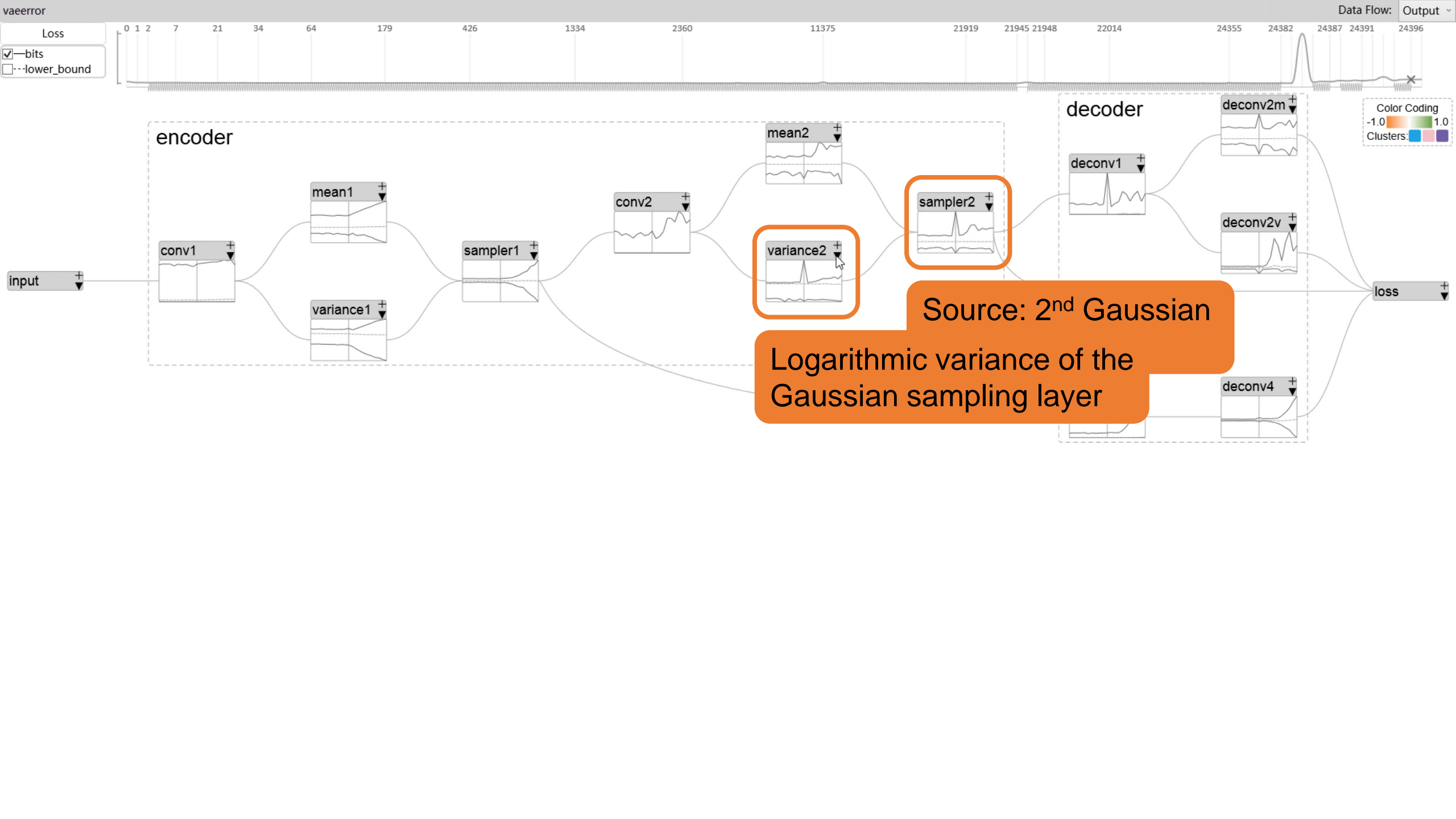


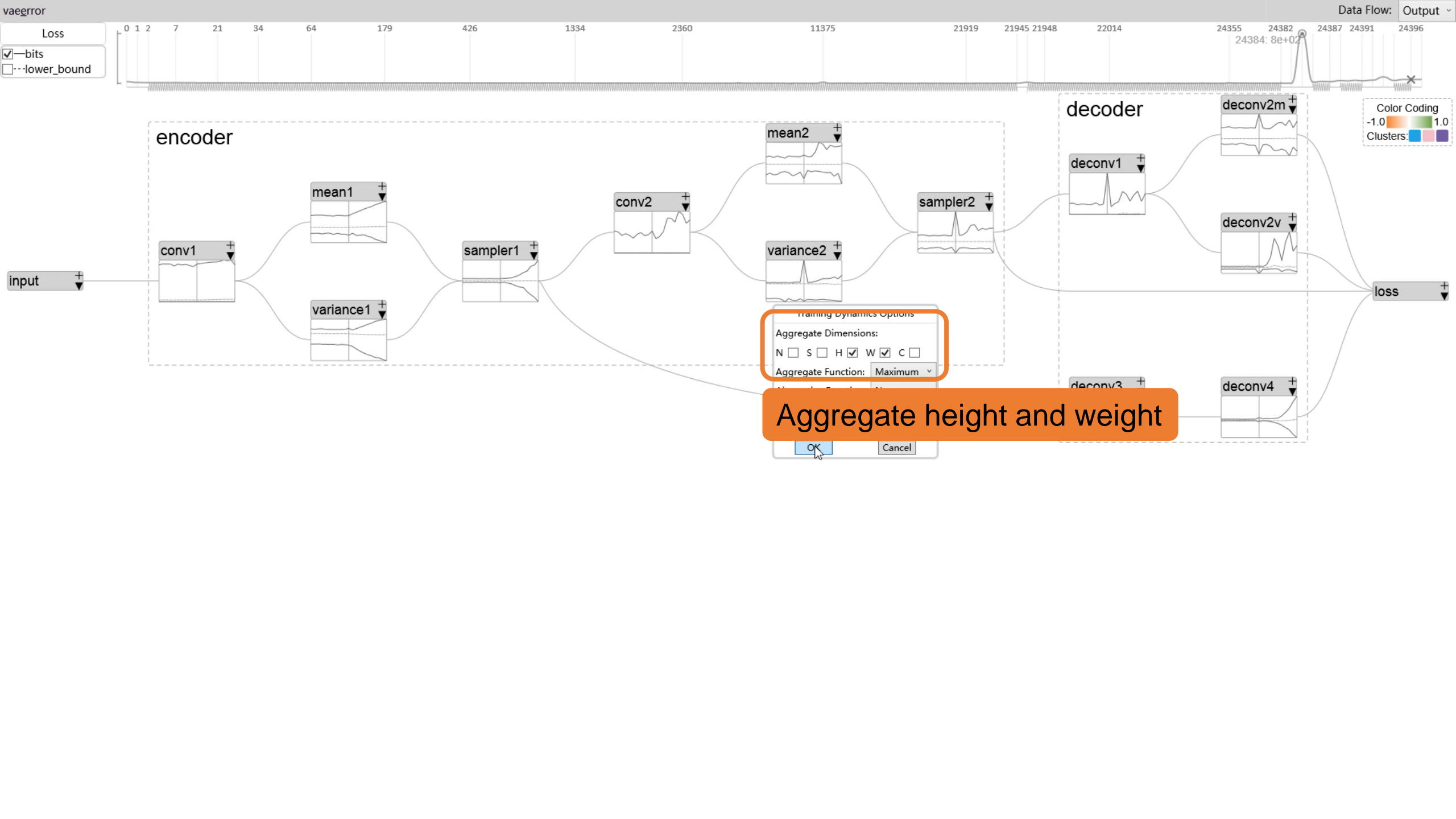
Loss changes

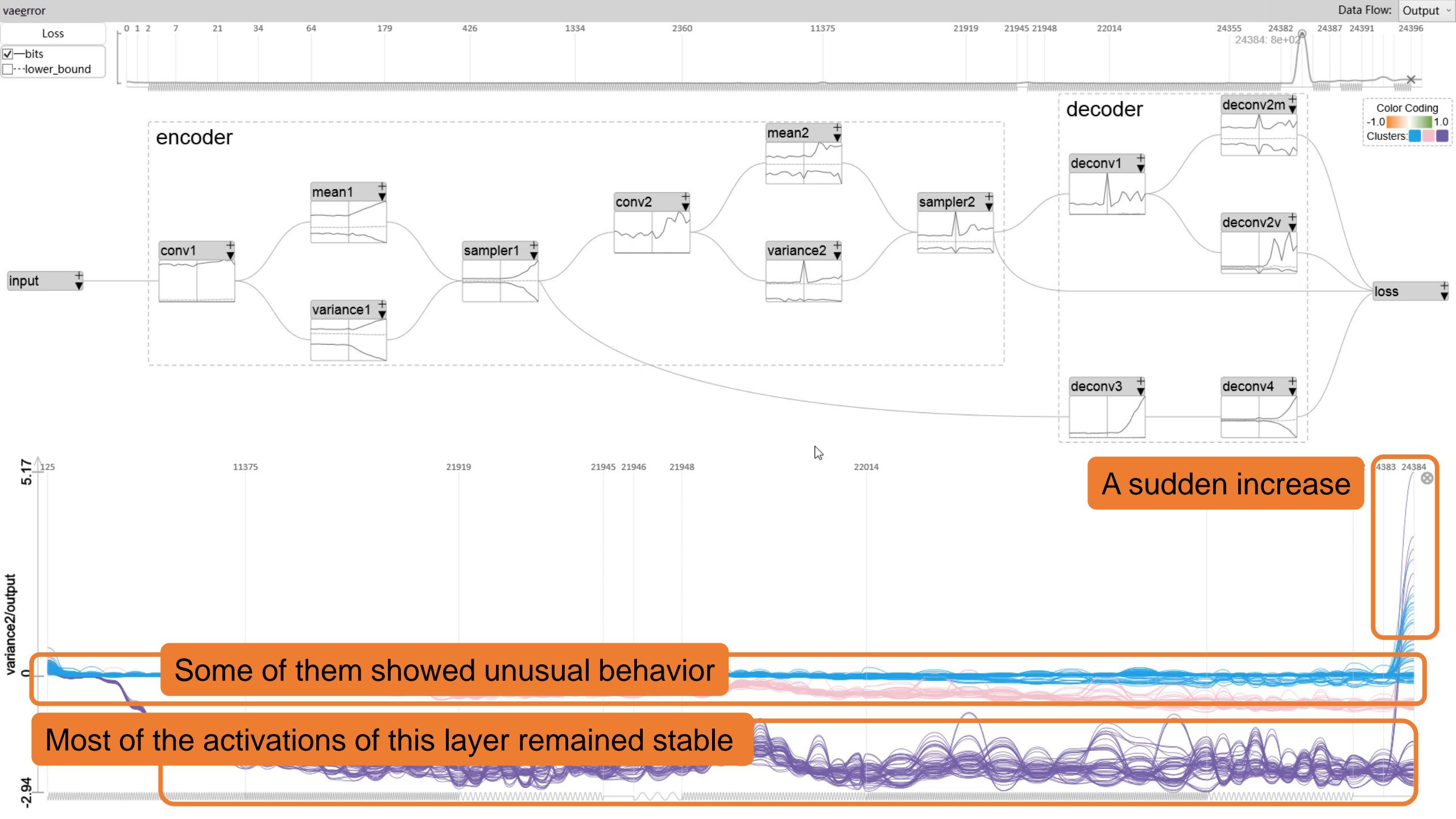
An abnormal snapshot fails at iteration 24,397

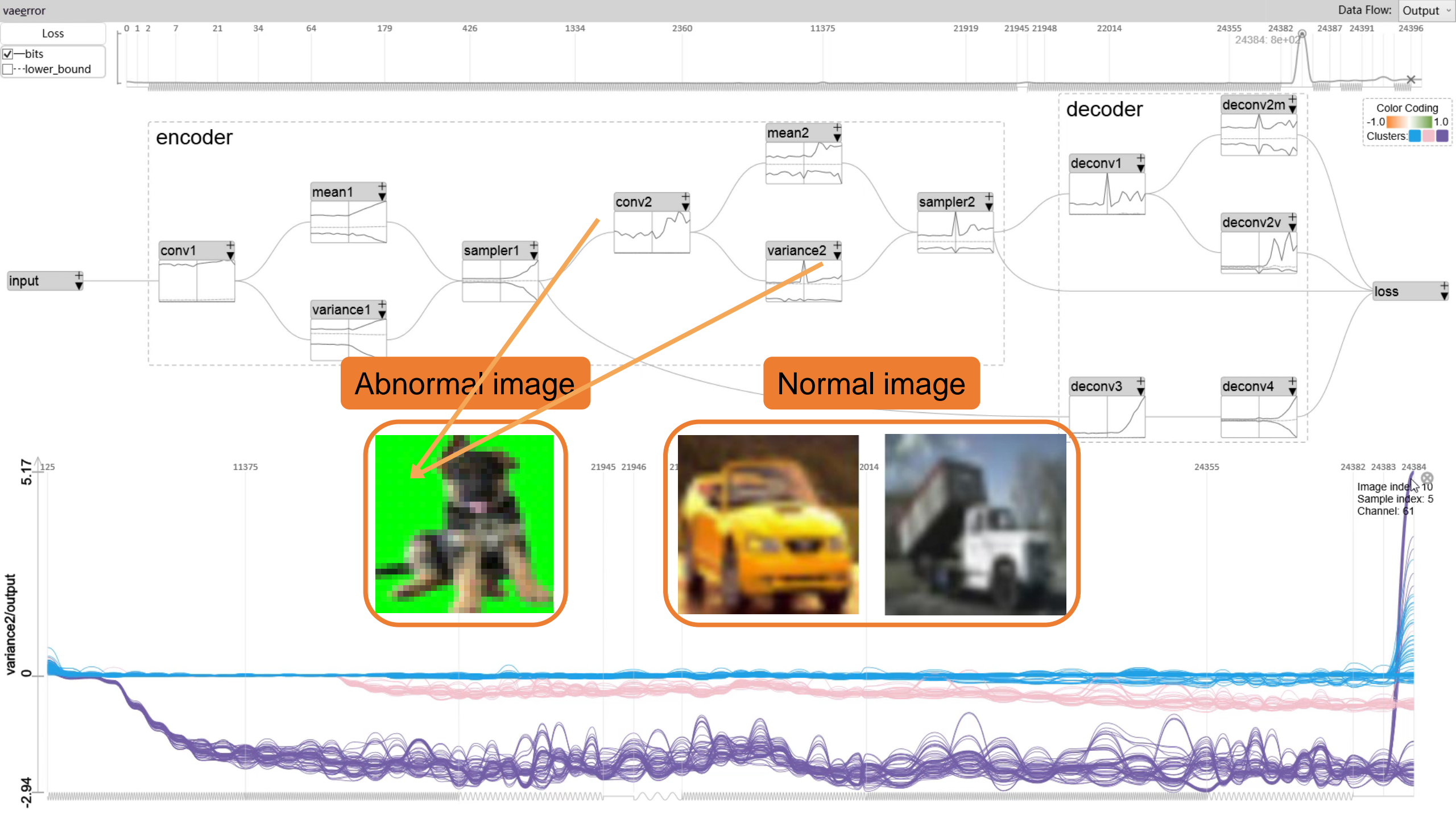








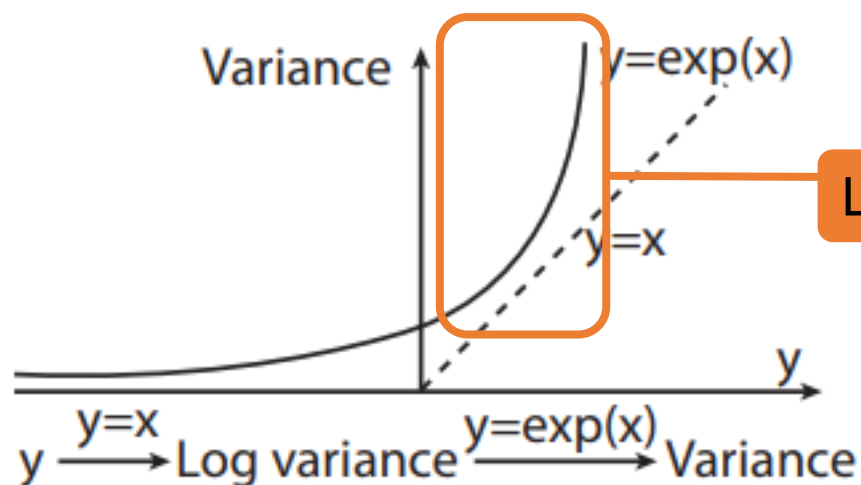




Solution

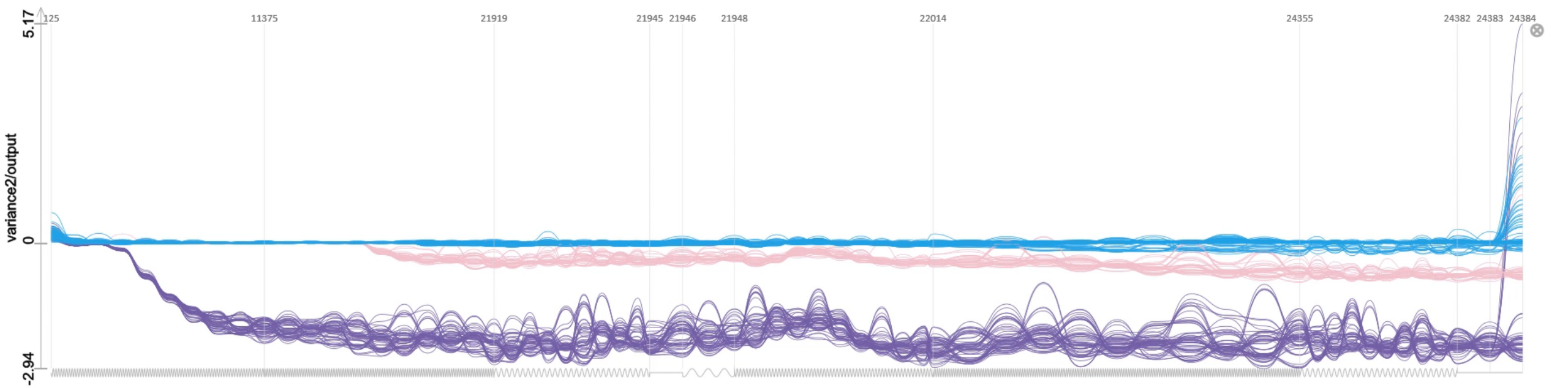
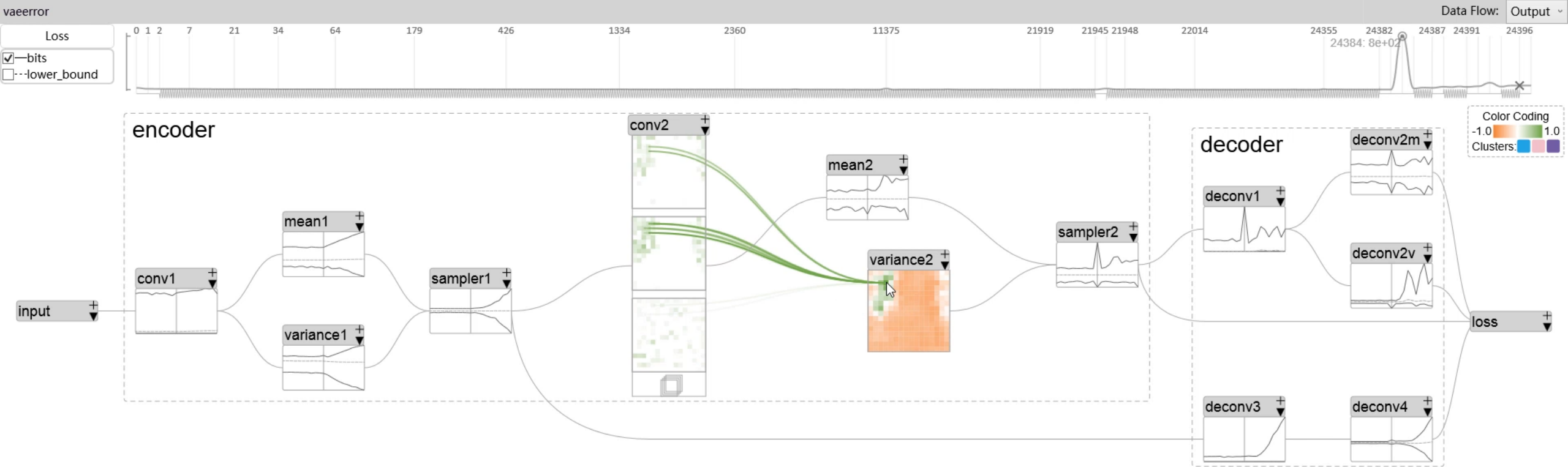
- Trial 1:
 - Replacing  with , but the network failed again
 - By the same analysis, we find another “bad” image 

- Trial 2:

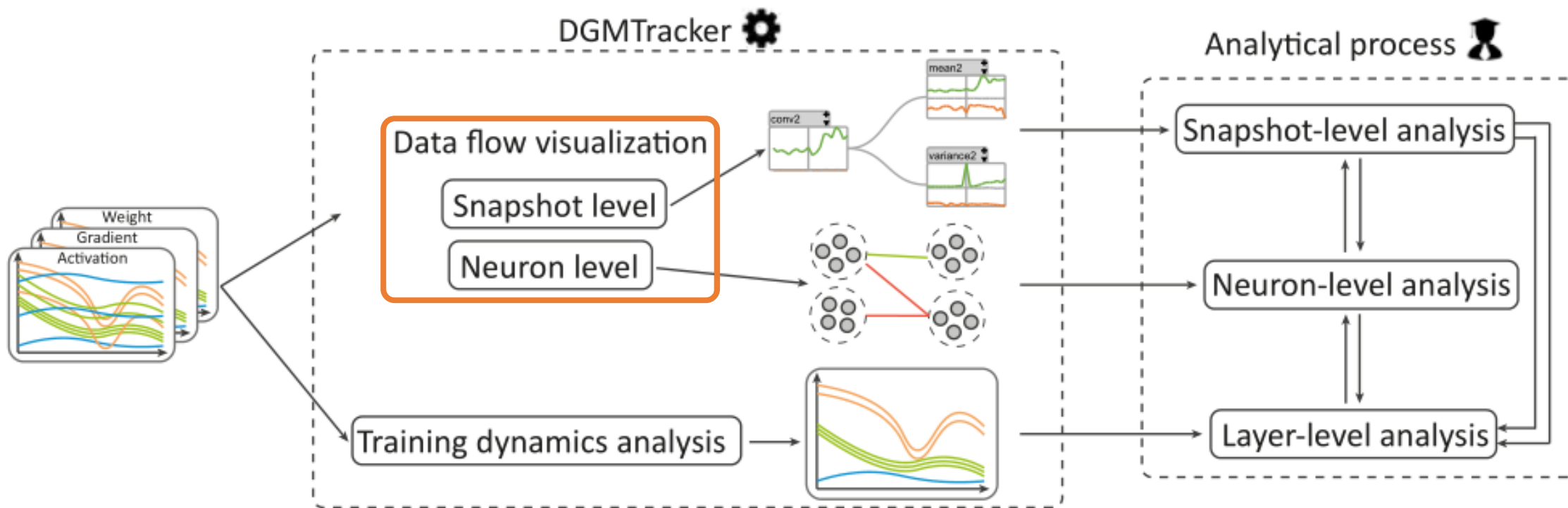


Large variance \rightarrow large samples \rightarrow large increase in loss

Much smoother

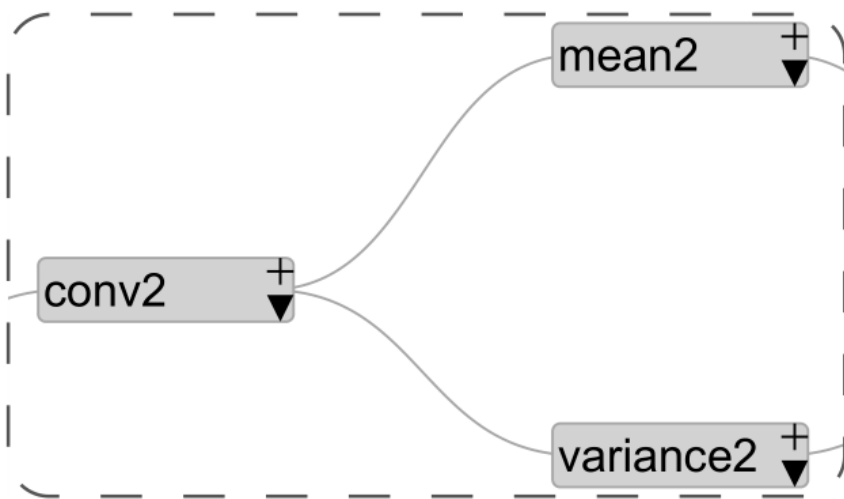


DGMTracker



Snapshot-Level Visualization

- How data flows through a network

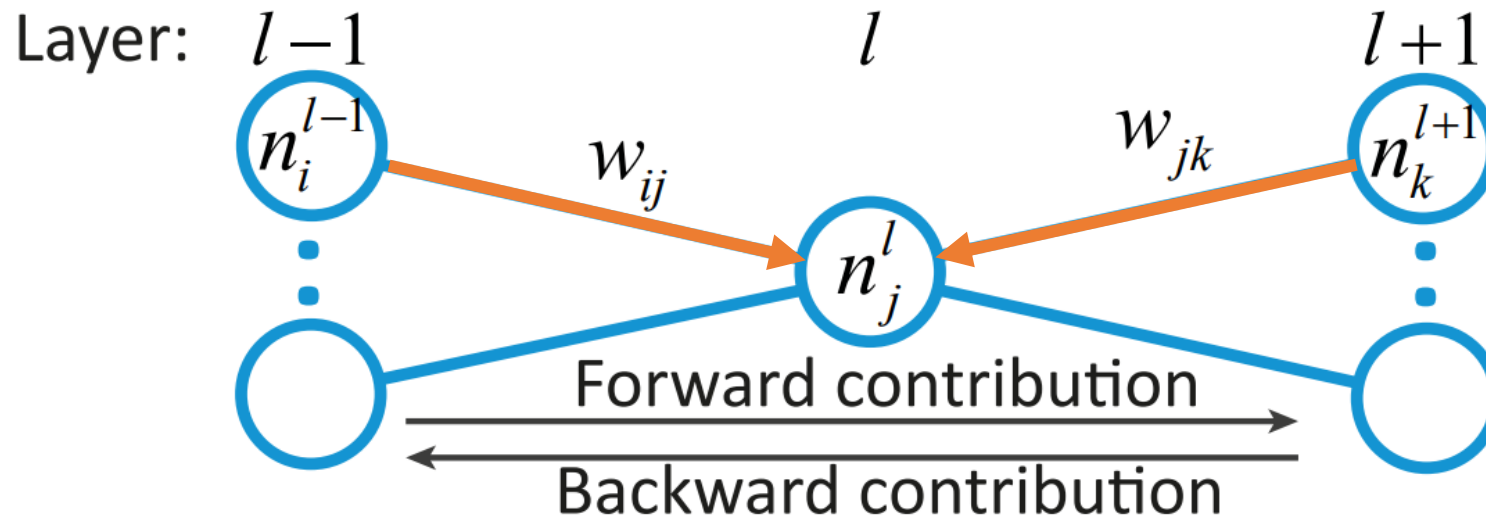


DAG visualization for DGM structure

Line chart for representing data flow

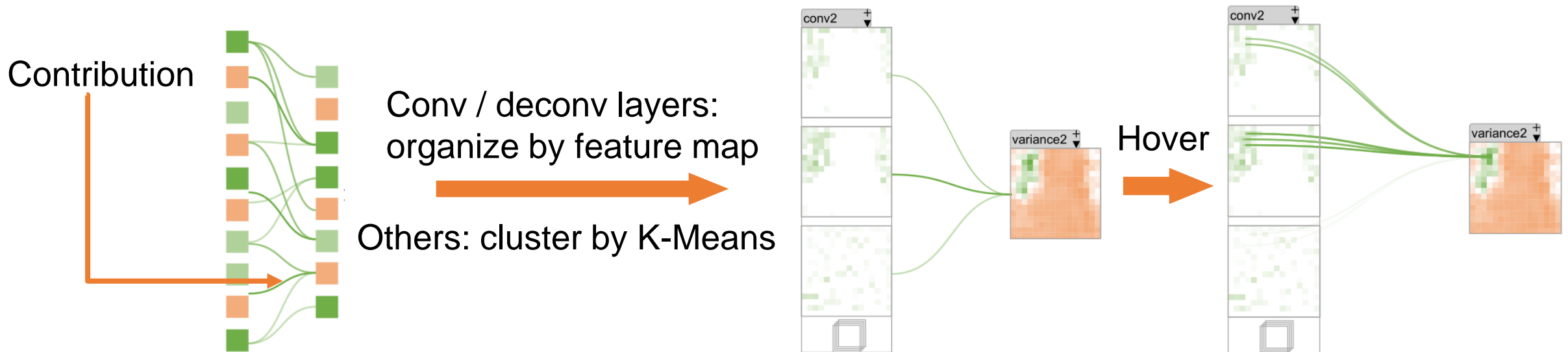
Neuron Level Visualization

- **Computing** and presenting how other neurons contribute to the output of the neuron being explored
 - Forward contribution: based on Layer-wise Relevance Propagation (LRP)
 - Backward contribution: based on backpropagation (BP)

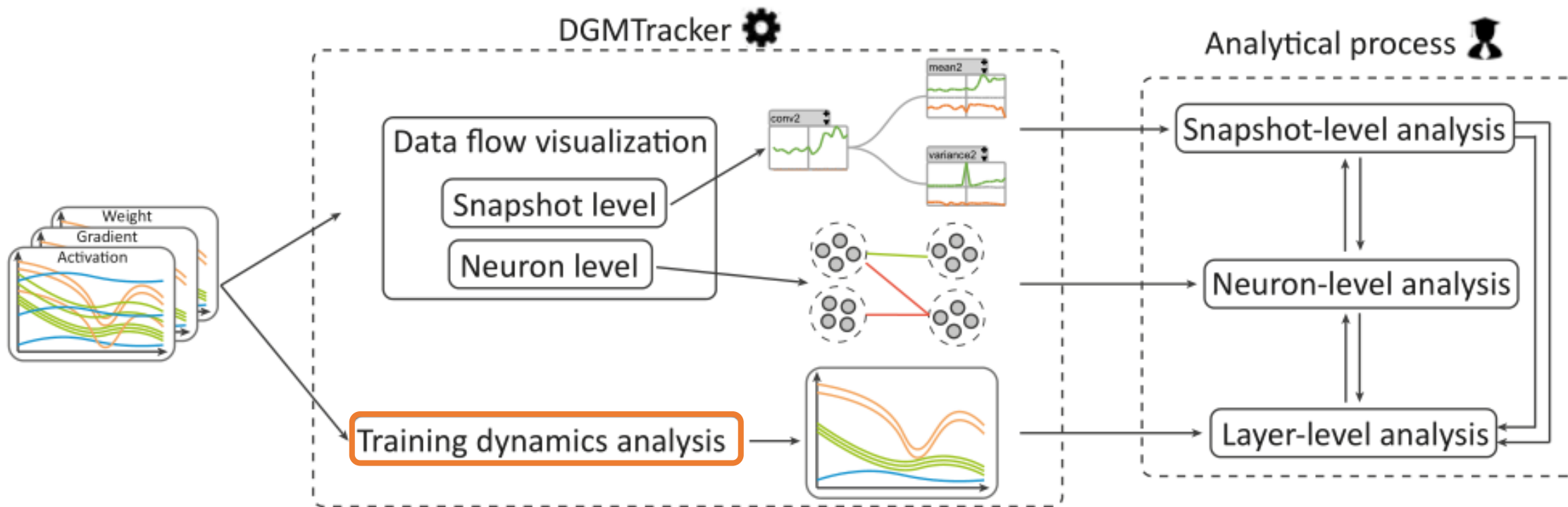


Neuron Level Visualization


- Computing and **presenting** how other neurons contribute to the output of the neuron being explored



DGMTracker

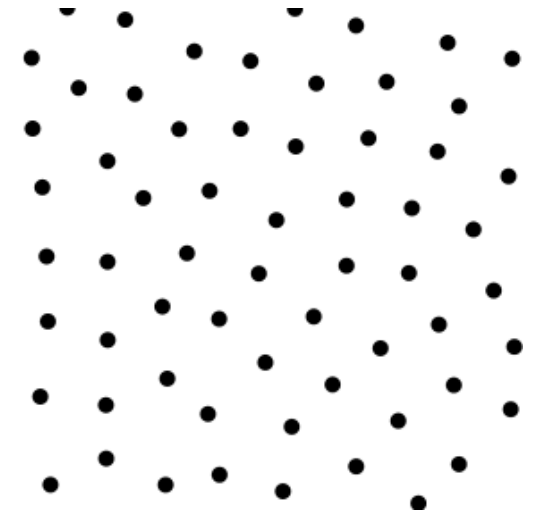


Training Dynamics Analysis

- Employ a line chart to visually convey the training dynamics
 - Training dynamics: activation/gradient/weight changes over time
- Challenge
 - Visual clutter caused by a large amount of time series data
- Solution:  polyline sampling
 - Both reduce visual clutter and preserve outliers

Blue Noise Sampling

- The selected samples have blue-noise properties
 - The selected samples are located **randomly** and **uniformly** in the space
- Compared with traditional random sampling
 - Low sampling rate in the high-density regions
 - Reduce visual clutter
 - High sampling rate in the low-density regions
 - Preserve outliers



Blue noise sampling example

Blue-Noise Polyline Sampling

- State-of-the-art: blue-noise line segment sampling [Sun et al., 2013]
 - Sample a line segment → if the distance with others is large enough, accept, or reject
- Intuitive method
 - Selecting line segment samples with blue-noise properties
 - Selecting polylines that contain the selected line segments as samples
- Solution: select “complete” polylines
 - How to select a polyline sample
 - How to compute the distance between two polylines

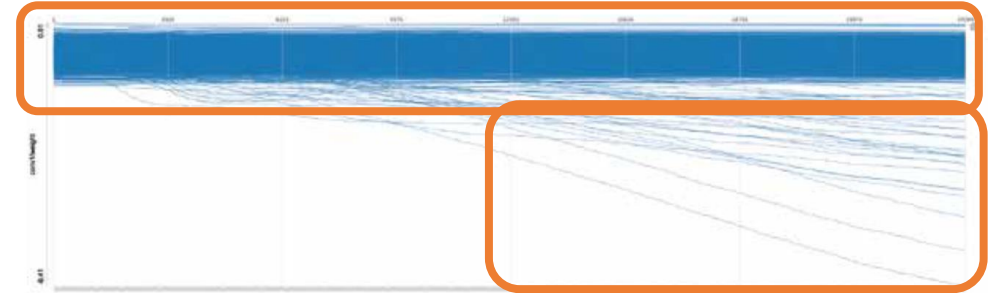


Blue-Noise Polyline Sampling (cont'd)

- How to select a polyline sample
 - Solution: select the polyline that can make the samples the most balanced in direction
- How to compute the distance between two polylines
 - Solution: sum of distances between corresponding line segments

$$d(L_1, L_2) = \frac{1}{N_S} \sum_{i=1}^{N_S} d_C(s_1^i, s_2^i)$$

Without sampling



Random sampling

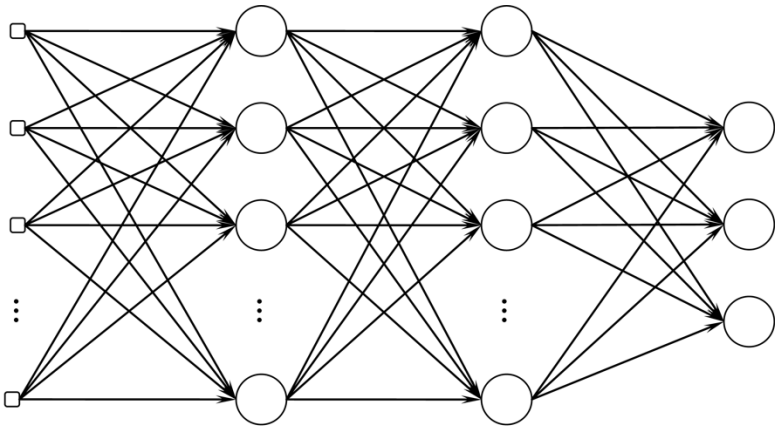


Blue-noise polyline sampling

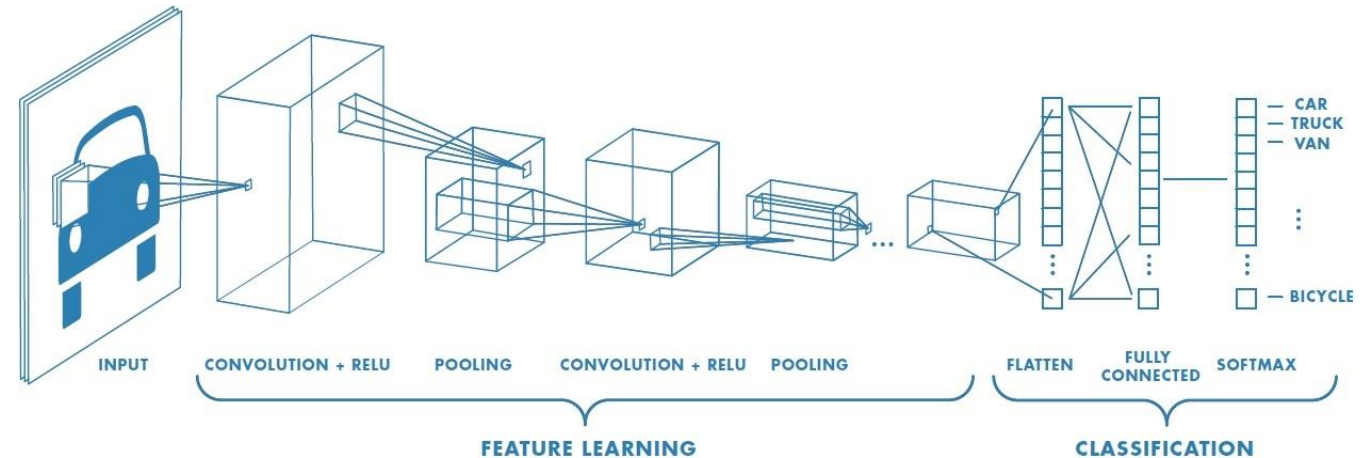


Generalization

- DGMTracker can be directly extended to other models, such as CNNs and MLPs
 - Often a base component of a DGM



Multiplayer perceptrons (MLPs)

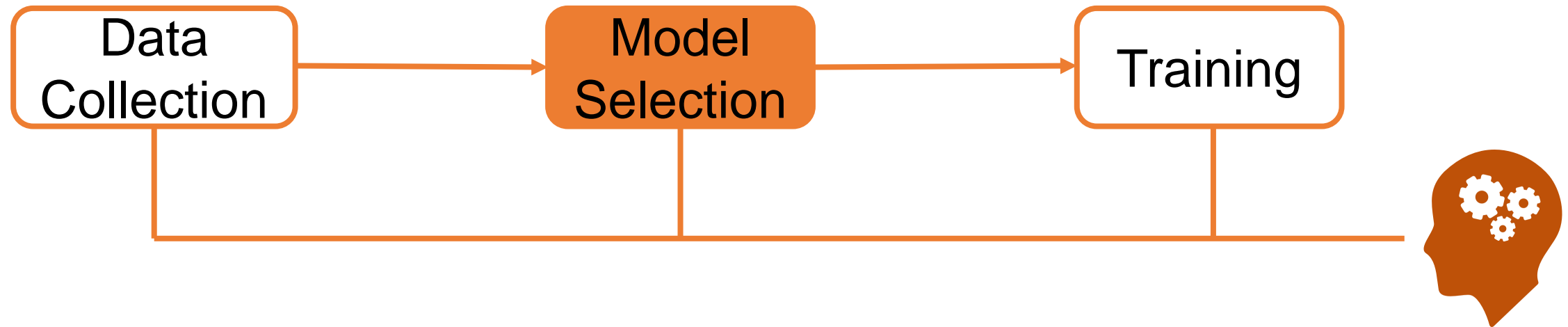


Convolutional neural networks (CNNs)

Conclusion and Future Work

- We have developed a visual analytics tool, DGMTracker, to facilitate machine learning experts in better understanding and diagnosing DGMs.
- Future work
 - Offline analysis → online analysis
 - Employ pattern mining techniques to disclose interesting patterns
 - Further reduce the amount of data needed for analysis: information theory
 - Originally 6TB → currently 6GB

Machine Learning + Visualization



Deep Learning. Goodfellow et al.,
Pattern Classification. Duda et al.

Towards Better Analysis of Deep Convolutional Neural Networks

Mengchen Liu, Jiaxin Shi, Zhen Li, Chongxuan Li, Jun Zhu, Shixia Liu

TVCG 2017

Challenges

DeepMind challenge match

One key factor: two deep CNNs →



May 2017



Why do the neural networks work?

AlphaGo



Beats

Nature match



Fan Hui (2p)

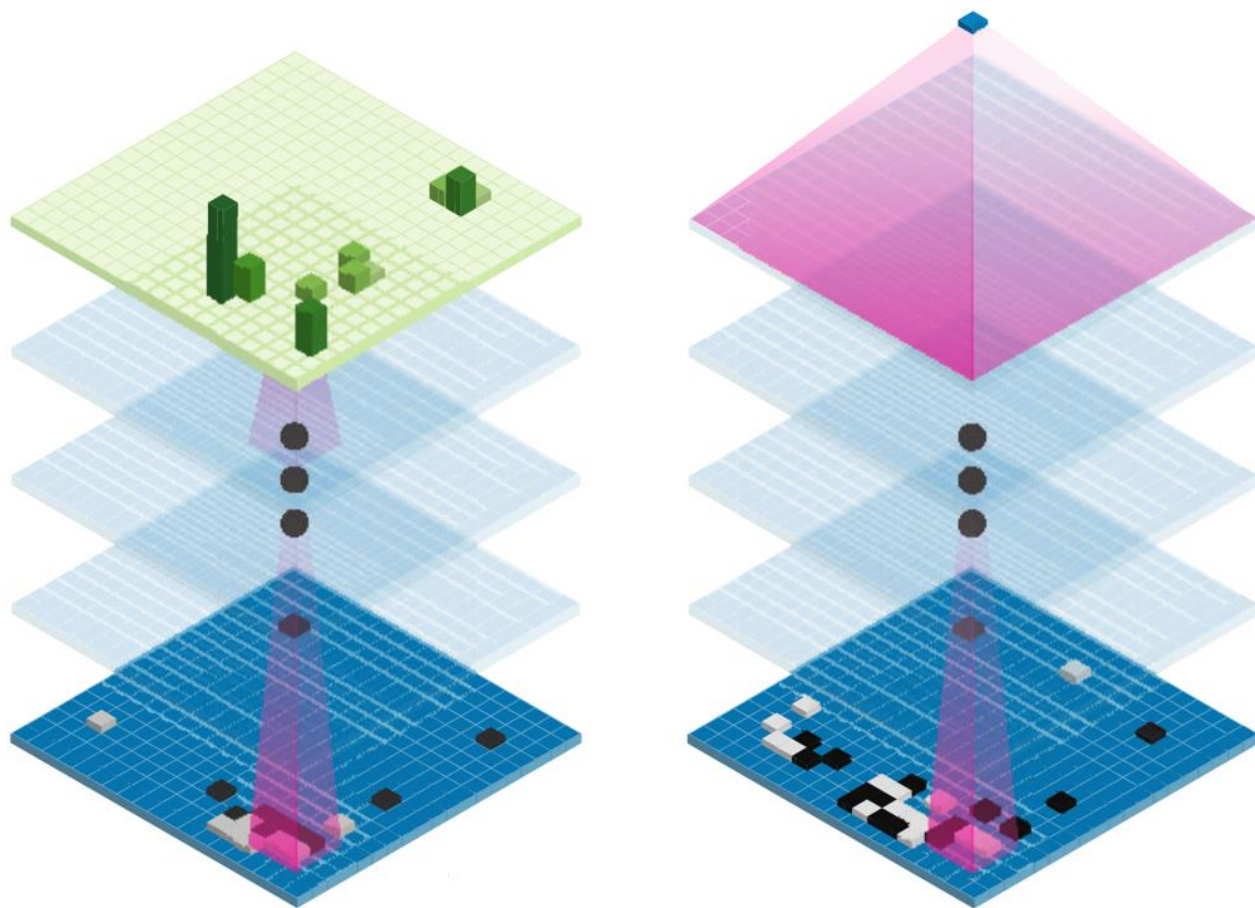


Oct 2015

Policy network

Value network

Challenges



- The size of a CNN is large
 - Tens or hundreds of layers (depth)
 - Thousands of neurons in each layer (width)
 - Millions of connections between neurons
- Many functional components
 - Their values and roles are not well understood

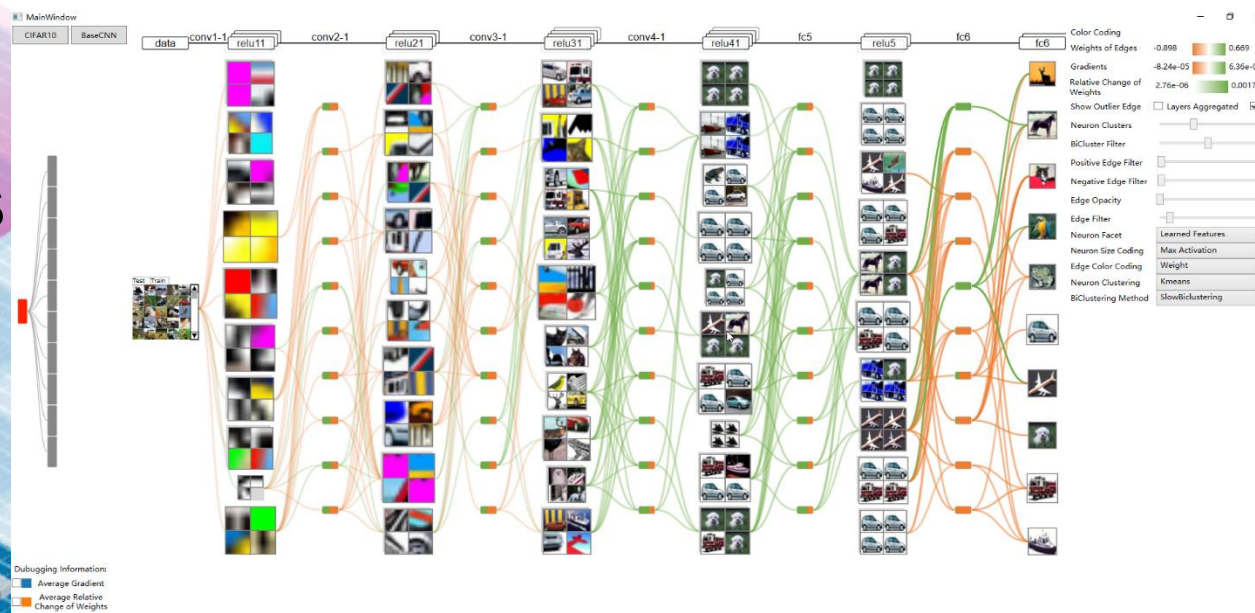
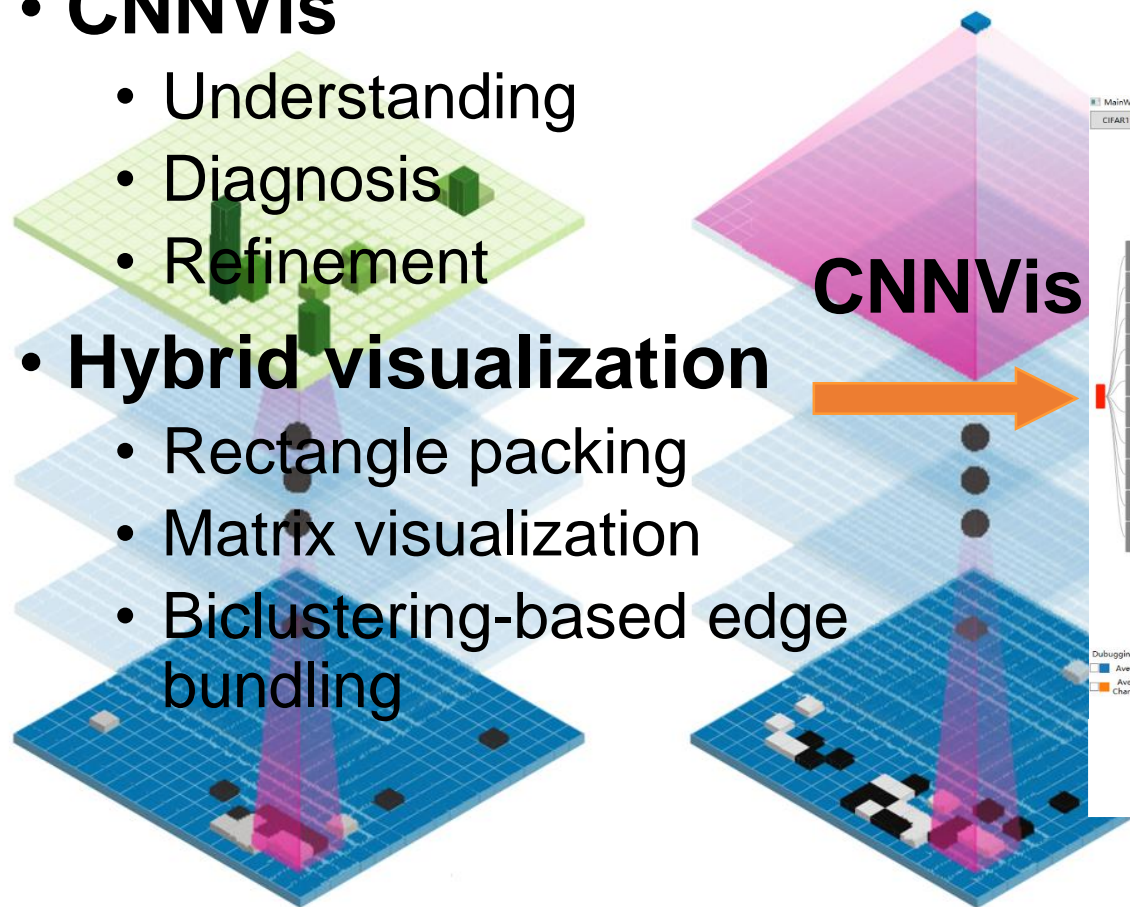
Contributions

- **CNNVis**

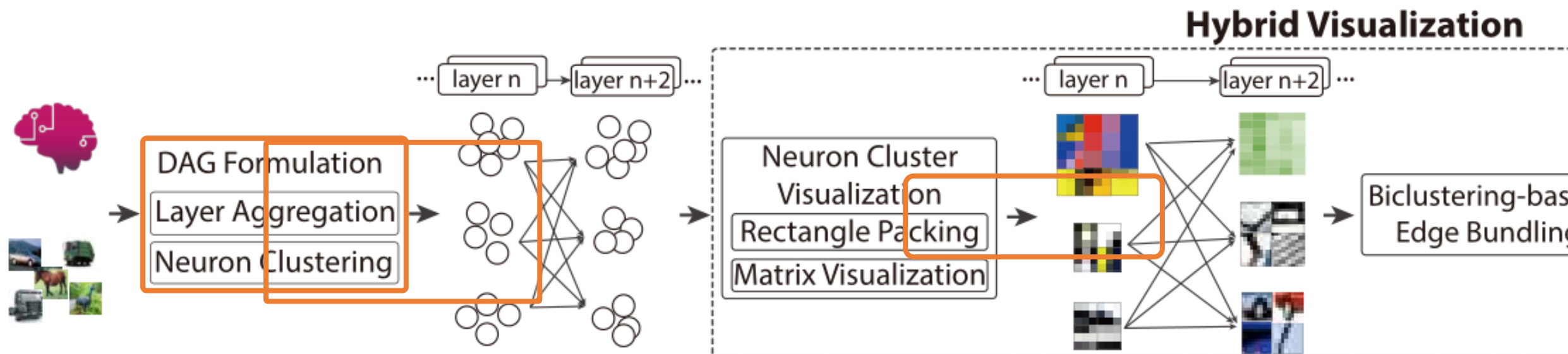
- Understanding
- Diagnosis
- Refinement

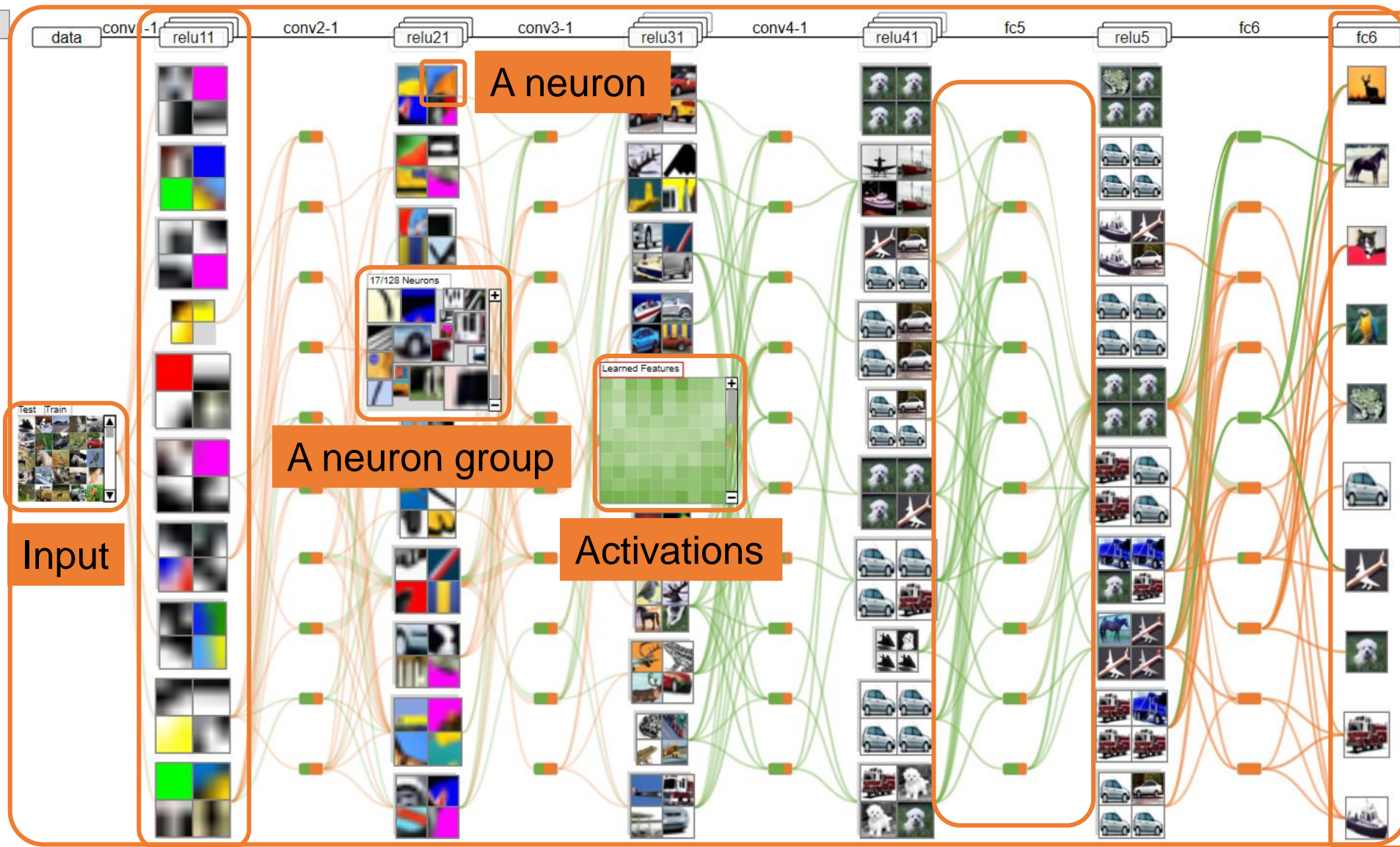
- **Hybrid visualization**

- Rectangle packing
- Matrix visualization
- Biclustering-based edge bundling



CNNVis Overview





Color Coding

Weights of Edges: -0.898 to 0.669

Gradients: -8.24e-05 to 6.36e-05

Relative Change of Weights: 2.76e-06 to 0.00179

Show Outlier Edge: Layers Aggregated:

Neuron Clusters: [Slider]

BiCluster Filter: [Slider]

Positive Edge Filter: [Slider]

Negative Edge Filter: [Slider]

Edge Opacity: [Slider]

Edge Filter: [Slider]

Neuron Facet: Learned Features

Neuron Size Coding: Max Activation

Edge Color Coding: Weight

Neuron Clustering: Kmeans

BiClustering Method: SlowBiclustering

Debugging Information:

Average Gradient

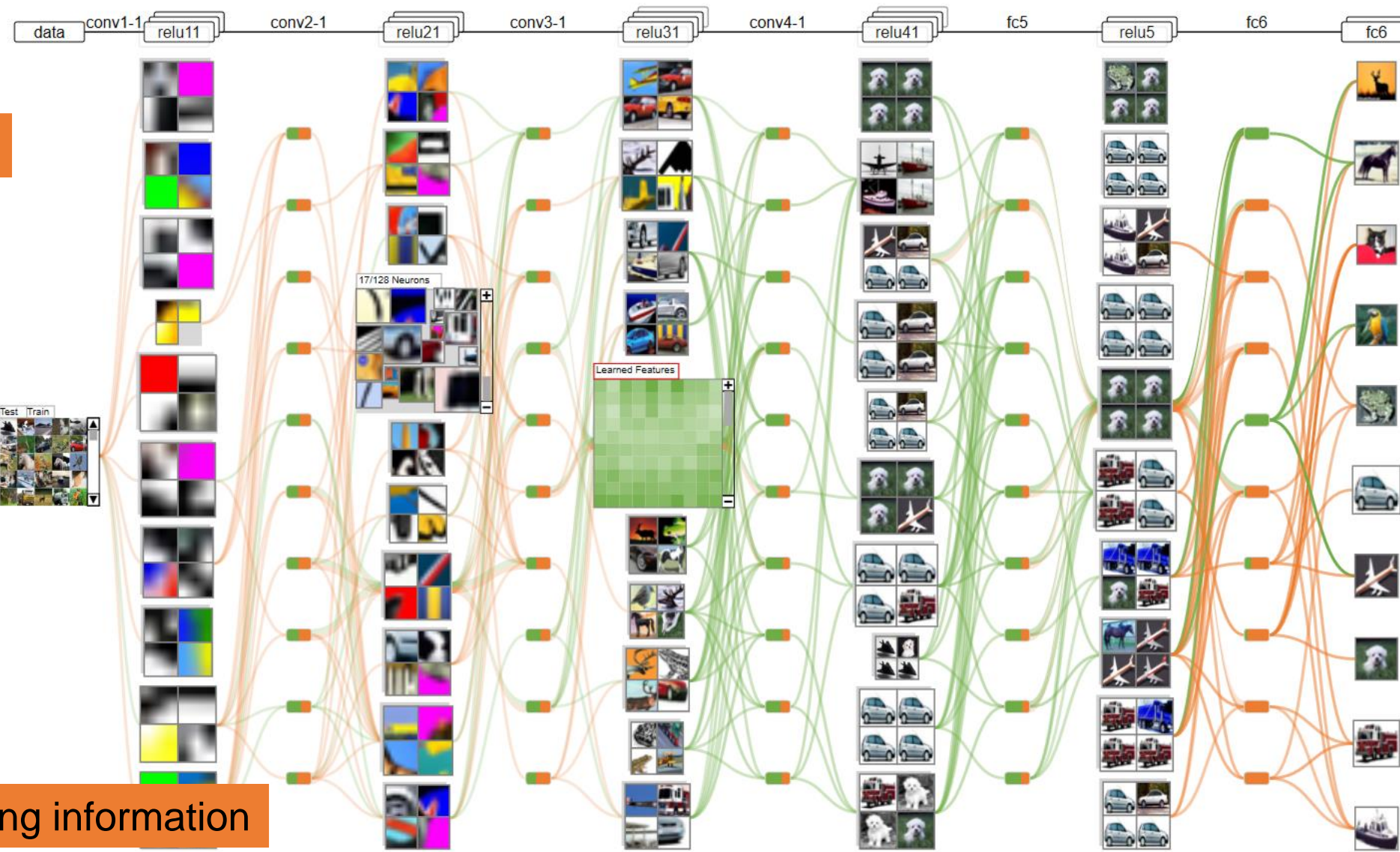
Average Relative Change of Weights

A layer group

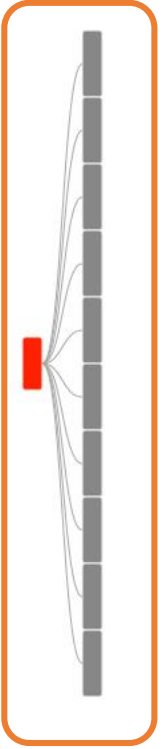
Connections

Output

2.37e-05



Selection



Color Coding

Weights of Edges: -0.898 to 0.669

Gradients: -8.24e-05 to 6.36e-05

Relative Change of Weights: 2.76e-06 to 0.00179

Show Outlier Edge: Layers Aggregated:

Neuron Clusters:

BiCluster Filter:

Positive Edge Filter:

Negative Edge Filter:

Edge Opacity:

Edge Filter:

Neuron Facet: Learned Features

Neuron Size Coding: Max Activation

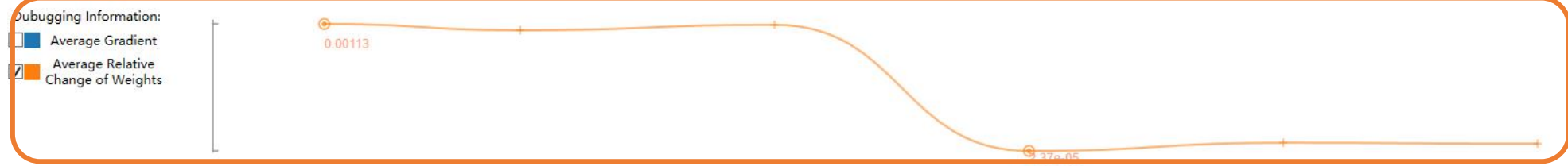
Edge Color Coding: Weight

Neuron Clustering: Kmeans

BiClustering Method: SlowBiclustering

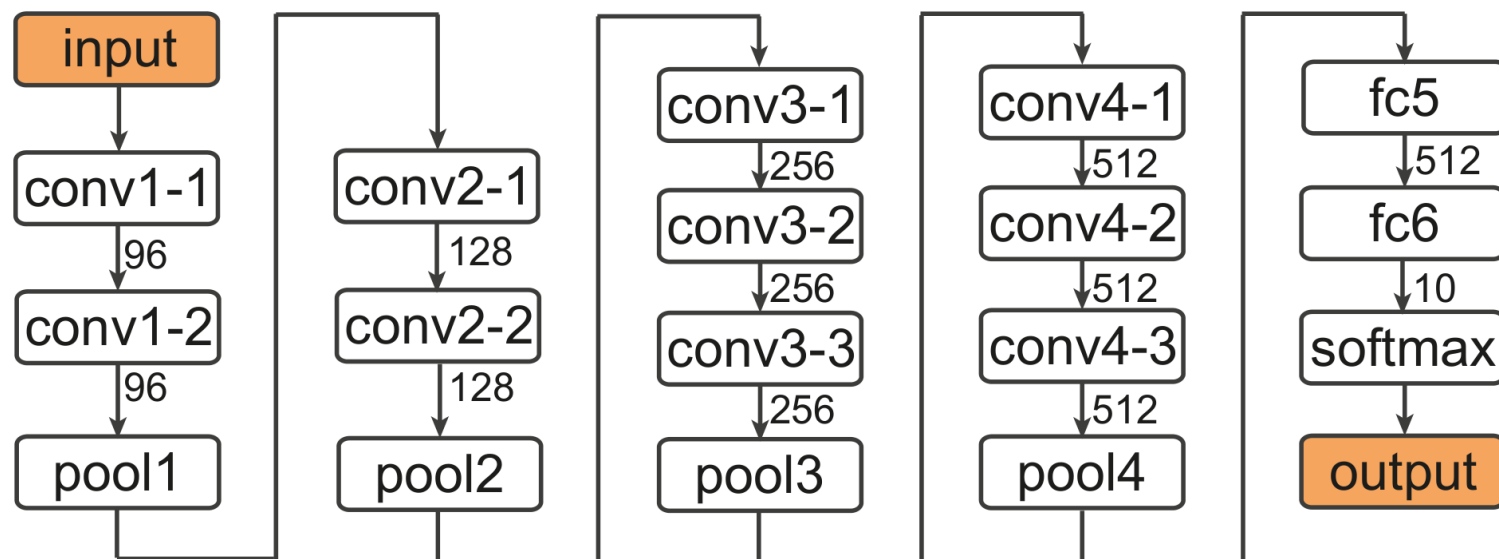
Control panel

Debugging information

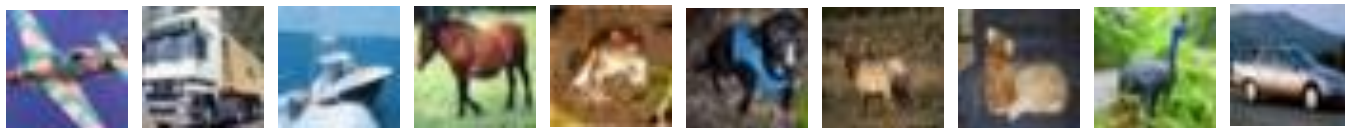


Case Study: Understanding

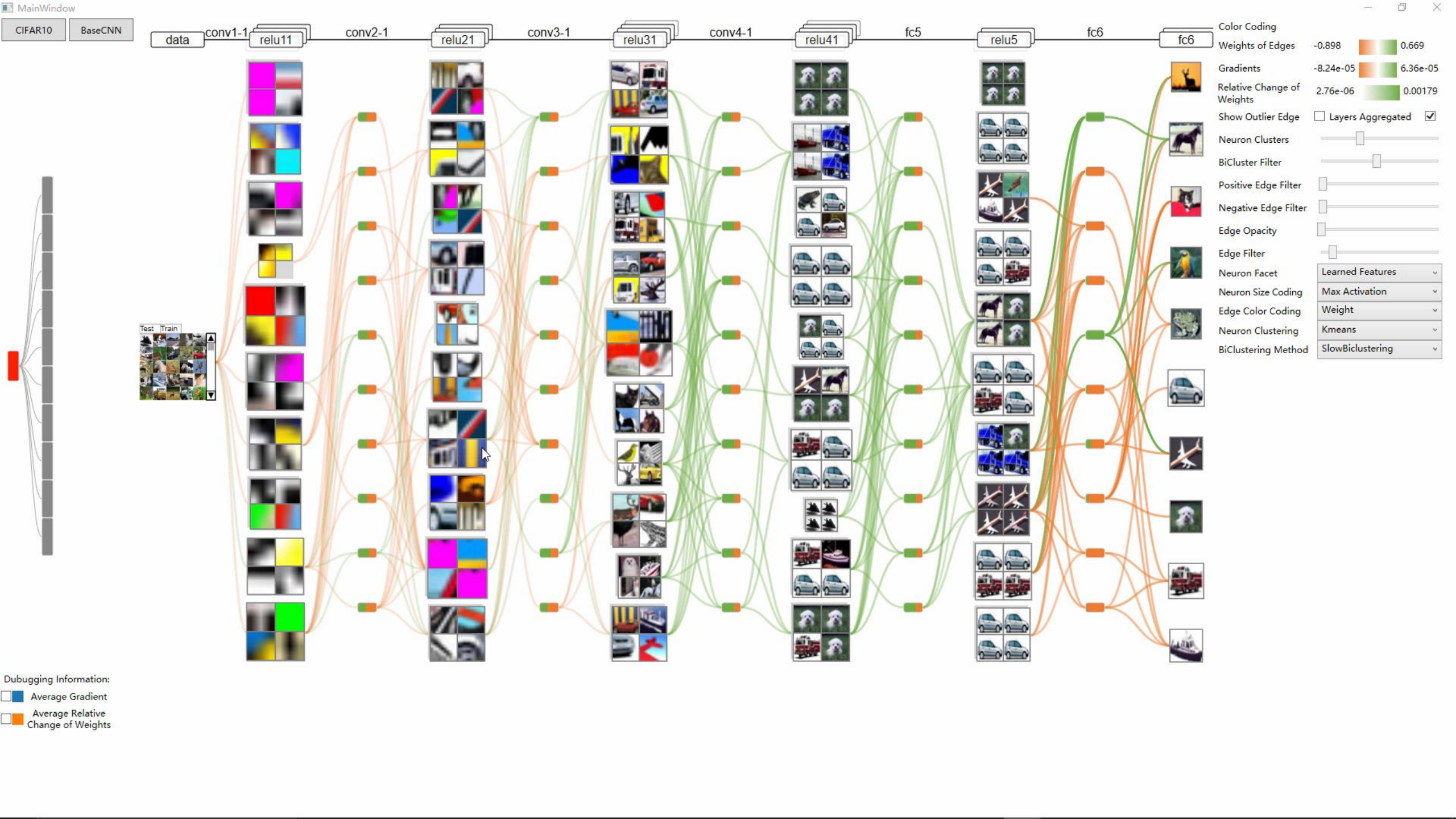
- Network: BaseCNN (Inspired by VGG)



- Dataset: CIFAR10



- Performance: 11.33% error rate





Debugging Information:

- Average Gradient
- Average Relative Change of Weights

Color Coding

Weights of Edges

Gradients

Relative Change of Weights

Show Outlier Edge Layers Aggregated

Neuron Clusters

BiCluster Filter

Positive Edge Filter

Negative Edge Filter

Edge Opacity

Edge Filter

Neuron Facet

Neuron Size Coding

Edge Color Coding

Neuron Clustering

BiClustering Method

Learned Features

Max Activation

Weight

Kmeans

SlowBiclustering

 Layers Aggregated

Neuron Clusters

BiCluster Filter

Positive Edge Filter

Negative Edge Filter

Edge Opacity

Edge Filter

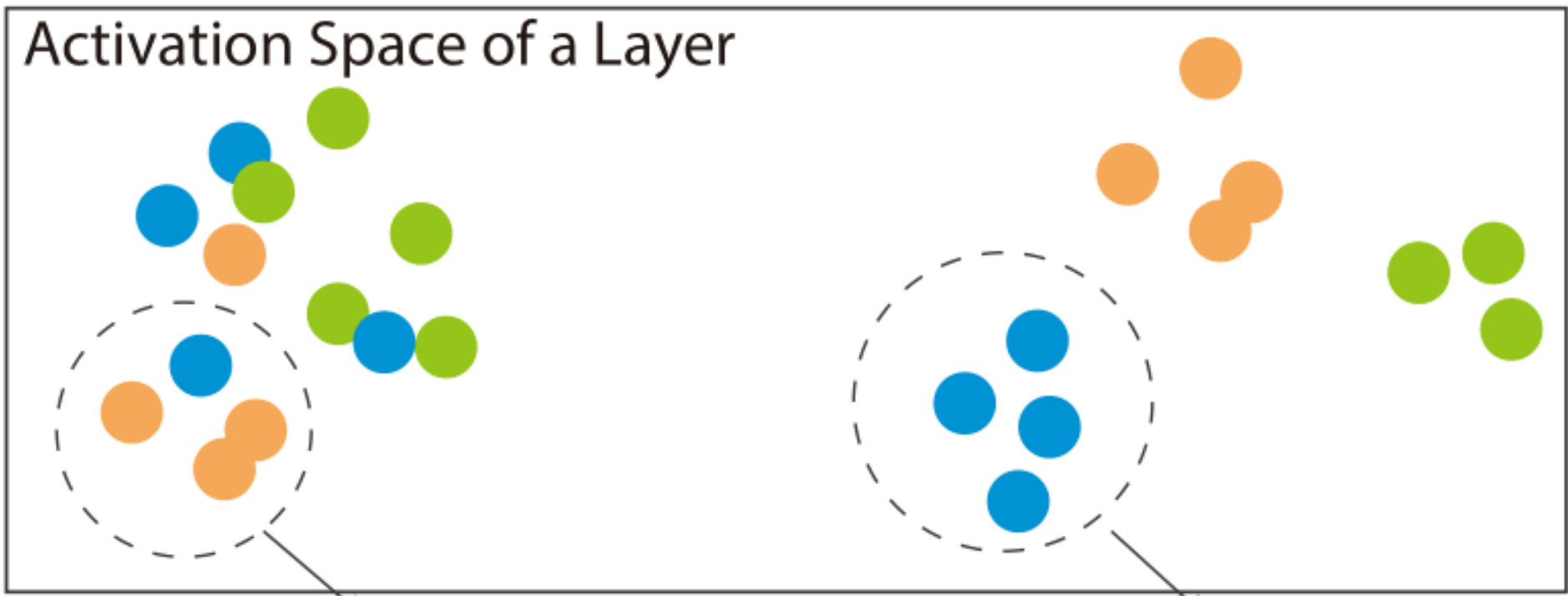
Neuron Facet

Neuron Size Coding

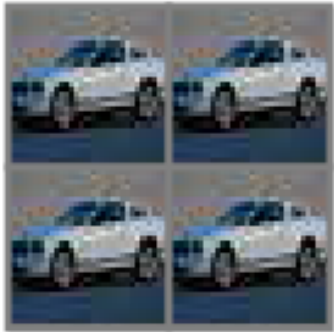
Edge Color Coding

Neuron Clustering

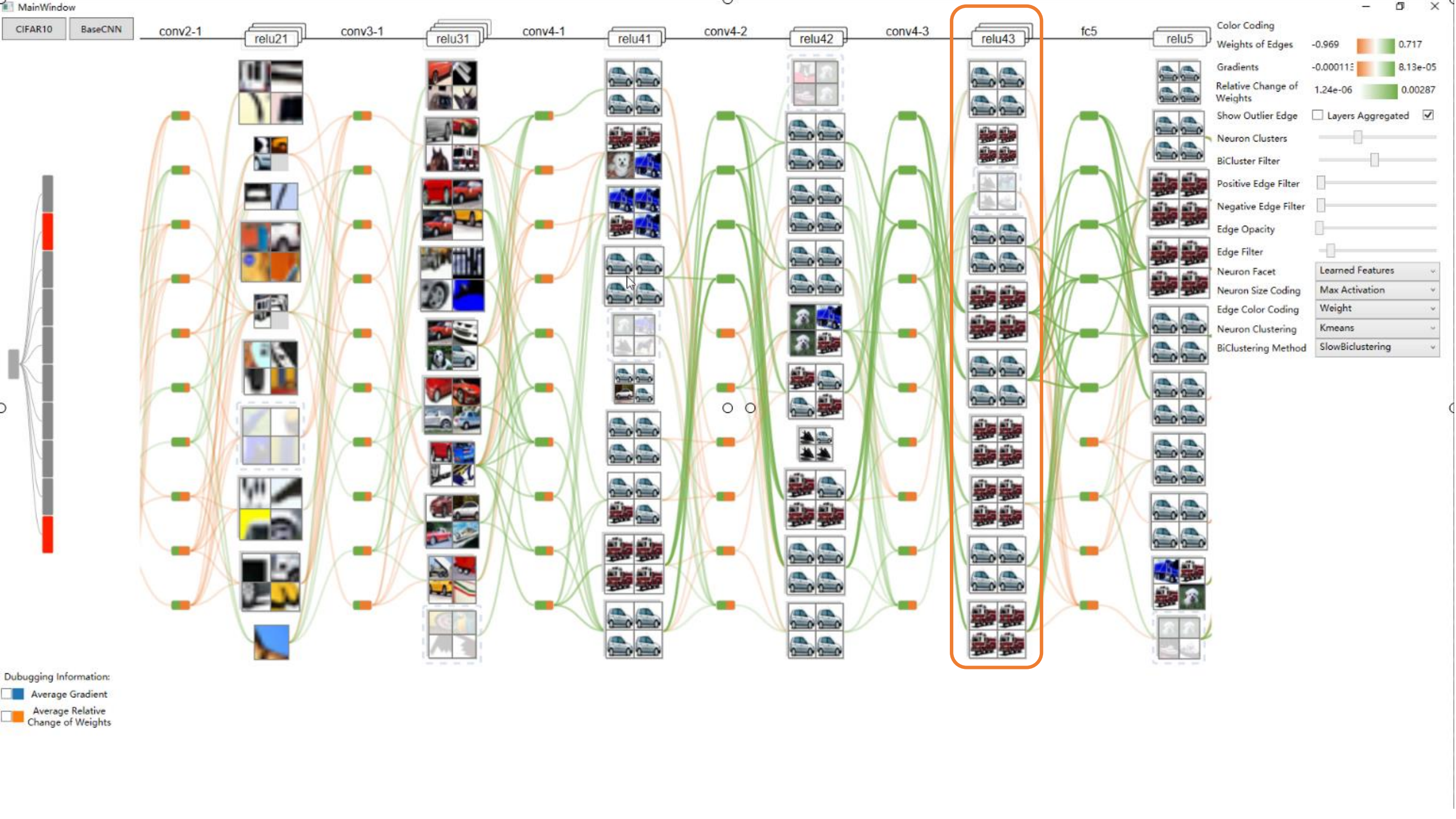
BiClustering Method



Impure Cluster



Pure Cluster

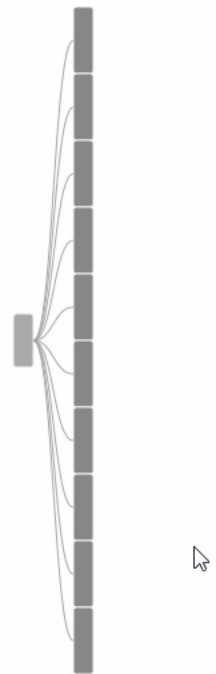


Network Depth

	Error	#ConvLayers	#Layers
→ ShallowCNN	11.94%	7	30
BaseCNN	11.33%	10	40
→ DeepCNN	14.77%	20	70

CIFAR10


ShallowCNN





Debugging Information:

- Average Gradient
- Average Relative Change of Weights

Color Coding

Weights of Edges 

Gradients 

Relative Change of Weights 

Show Outlier Edge Layers Aggregated

Neuron Clusters

BiCluster Filter

Positive Edge Filter

Negative Edge Filter

Edge Opacity

Edge Filter

Neuron Facet Learned Features ▾

Neuron Size Coding Max Activation ▾

Edge Color Coding Weight ▾

Neuron Clustering Kmeans ▾

BiClustering Method SlowBiclustering ▾



Color Coding

Weights of Edges

Gradients

Relative Change of Weights

Show Outlier Edge Layers Aggregated

Neuron Clusters

BiCluster Filter

Positive Edge Filter

Negative Edge Filter

Edge Opacity

Edge Filter

Neuron Facet Learned Features ▾

Neuron Size Coding Max Activation ▾

Edge Color Coding Weight ▾

Neuron Clustering Kmeans ▾

BiClustering Method SlowBiclustering ▾

Debugging Information:

Average Gradient

Average Relative Change of Weights

Network Width

	Error	#Params	Training loss	Testing loss	
→ BaseCNN×4	12.33%	4.22M	0.04	0.51	↑ Overfitting ↓ Underfitting
BaseCNN×2	11.47%	2.11M	0.07	0.43	
BaseCNN	11.33%	1.05M	0.16	0.40	
BaseCNN×0.5	12.61%	0.53M	0.34	0.40	
→ BaseCNN×0.25	17.39%	0.26M	0.65	0.53	

CIFAR10

BaseCNNx0.25



Color Coding

Weights of Edges

Gradients

Relative Change of Weights

Show Outlier Edge Layers Aggregated

Neuron Clusters

BiCluster Filter

Positive Edge Filter

Negative Edge Filter

Edge Opacity

Edge Filter

Neuron Facet

Neuron Size Coding

Edge Color Coding

Neuron Clustering

BiClustering Method

Learned Features

Max Activation

Weight

Kmeans


SlowBiclustering


Debugging Information:


 Average Gradient Average Relative Change of Weights




Color Coding


Weights of Edges 


Gradients 


Relative Change of Weights 

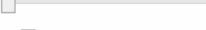
Show Outlier Edge Layers Aggregated

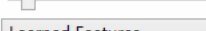
Neuron Clusters 

BiCluster Filter 

Positive Edge Filter 

Negative Edge Filter 

Edge Opacity 

Edge Filter 

Neuron Facet Learned Features ▾

Neuron Size Coding Max Activation ▾

Edge Color Coding Weight ▾

Neuron Clustering Kmeans ▾

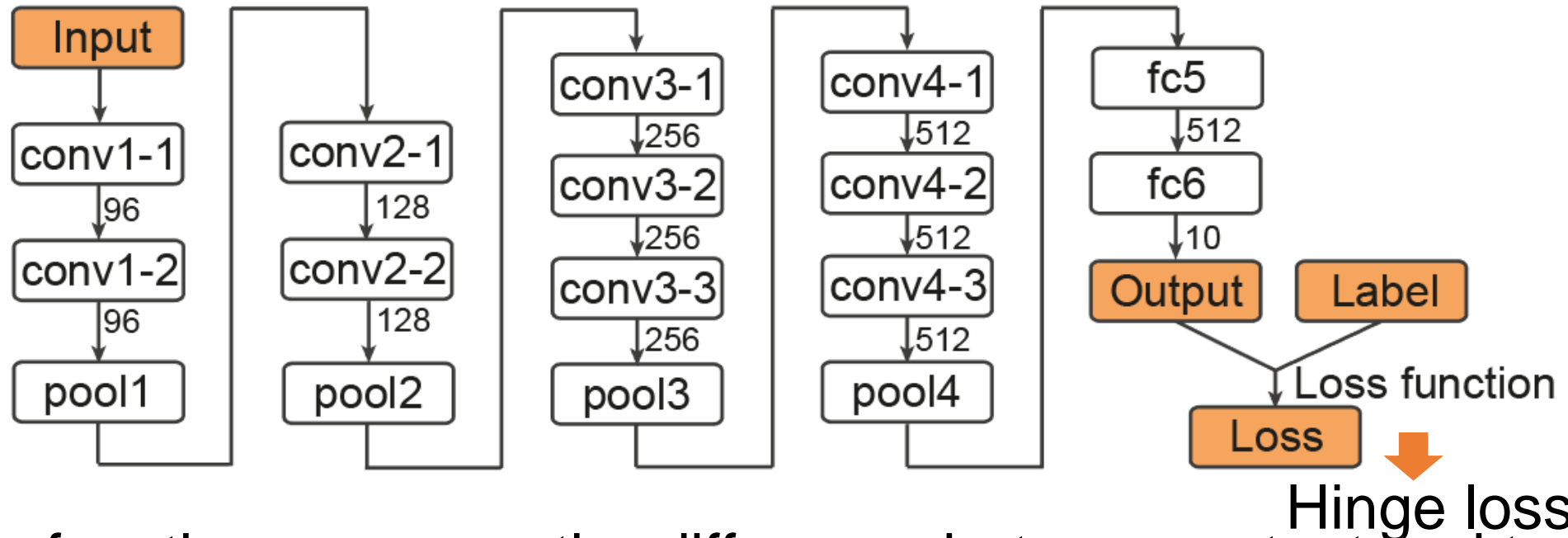
BiClustering Method SlowBiClustering ▾

Debugging Information:

- Average Gradient
- Average Relative Change of Weights

Case Study: Training Diagnosis

- Network:





- Loss function: measure the difference between output and true labels
- Training stuck at loss=2.0 (far from achieving good accuracy)


cifar10

lhinge7_4


Color Coding


Weights of Edges 


Gradients 


Relative Change of Weights 


Show Outlier Edge Layers Aggregated


Neuron Clusters 

BiCluster Filter 

Positive Edge Filter 

Negative Edge Filter 

Edge Opacity 

Edge Filter 

Neuron Facet Learned Features ▾

Neuron Size Coding Max Activation ▾

Edge Color Coding Weight ▾

Neuron Clustering Kmeans ▾

BiClustering Method SlowBiclustering ▾



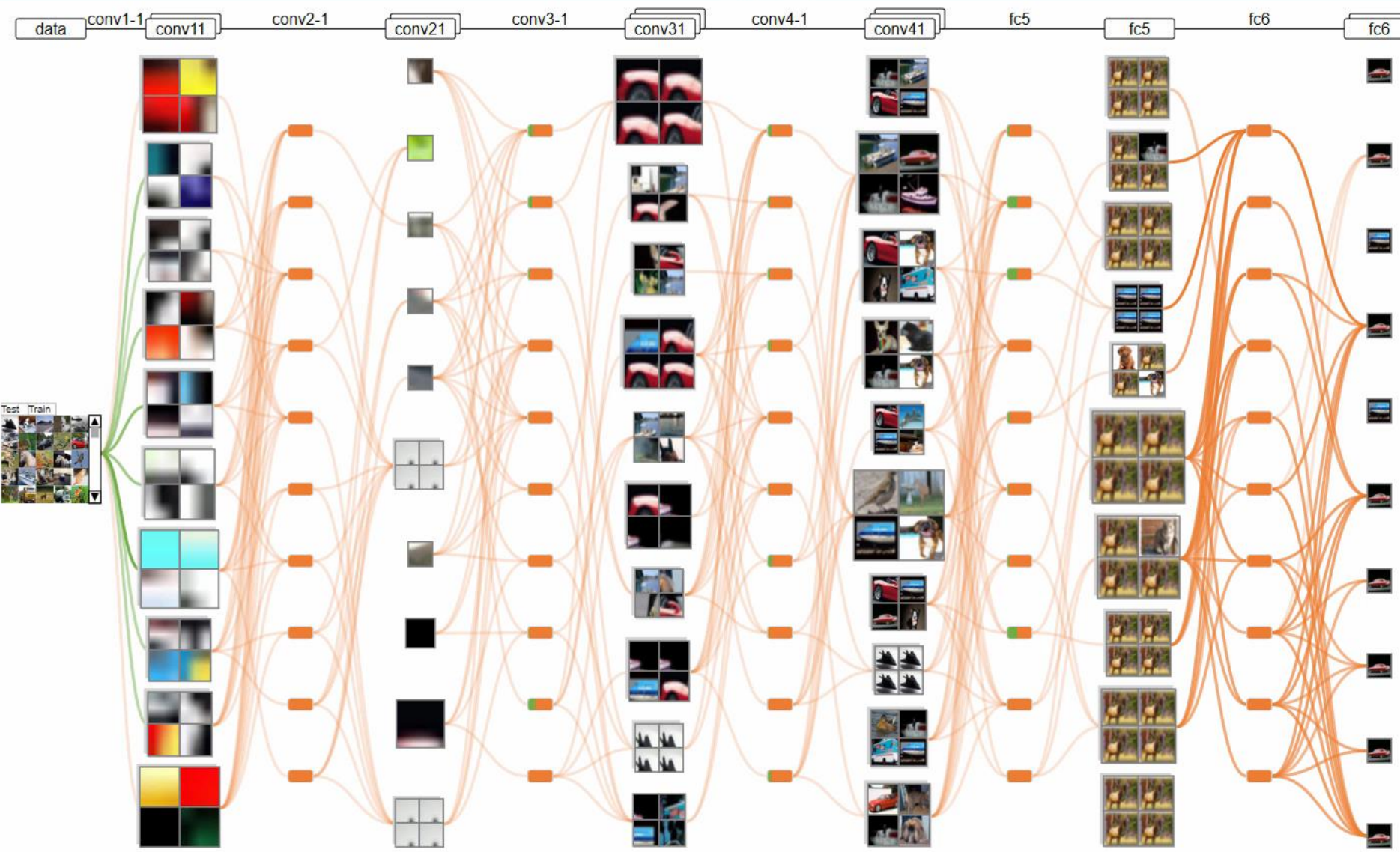
Debugging Information:

Average Gradient

Average Relative Change of Weights



cifar10 | l1hinge7_4



Color Coding

Weights of Edges: -1.61 0.195

Gradients: -6.65e-06 1.66e-06

Relative Change of Weights: 2.39e-06 3.12e-05

Show Outlier Edge Layers Aggregated

Neuron Clusters

BiCluster Filter

Positive Edge Filter

Negative Edge Filter

Edge Opacity

Edge Filter

Neuron Facet: Learned Features

Neuron Size Coding: Max Activation

Edge Color Coding: Weight

Neuron Clustering: Kmeans

BiClustering Method: SlowBiclustering

Debugging Information:

Average Gradient

Average Relative Change of Weights



Color Coding |

Weights of Edges -1.61 0.195

Gradients -6.65e-06 1.66e-06

Relative Change of Weights 2.39e-06 3.12e-05

Show Outlier Edge Layers Aggregated

Neuron Clusters

BiCluster Filter

Positive Edge Filter

Negative Edge Filter

Edge Opacity

Edge Filter

Neuron Facet Activation Matrix

Neuron Size Coding Max Activation

Edge Color Coding Weight

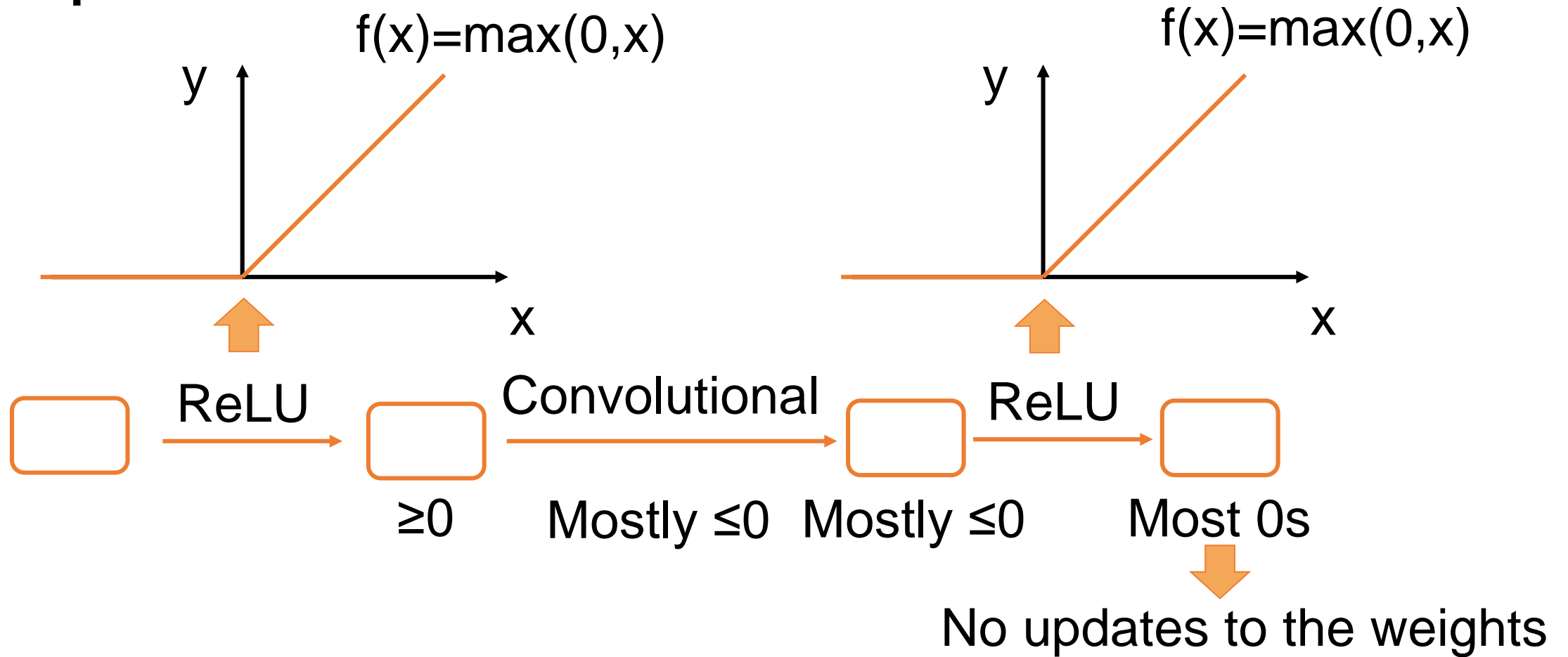
Neuron Clustering Kmeans

BiClustering Method SlowBiclustering

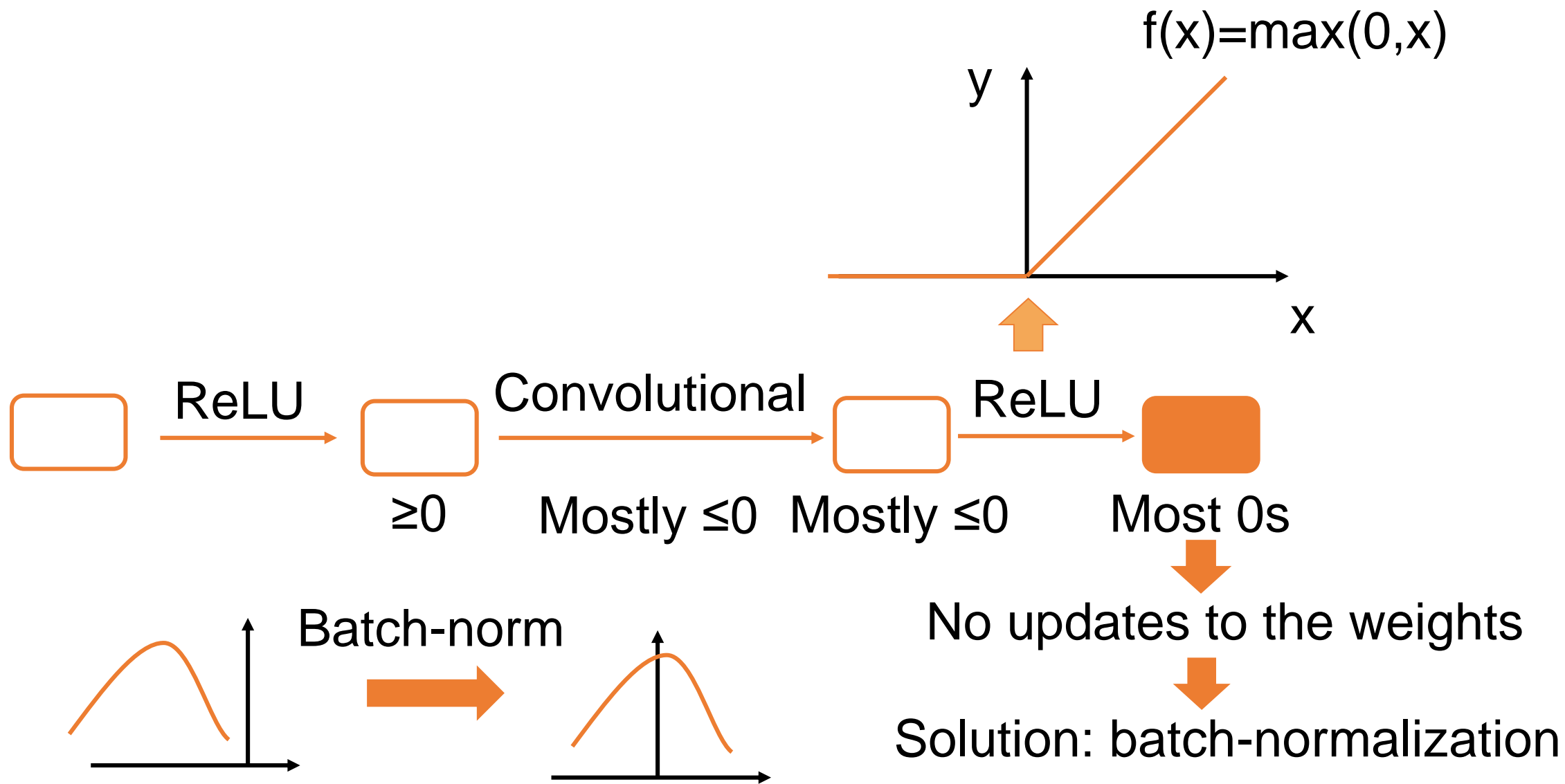
Debugging Information:

- Average Gradient
- Average Relative Change of Weights

Explanation



Solution



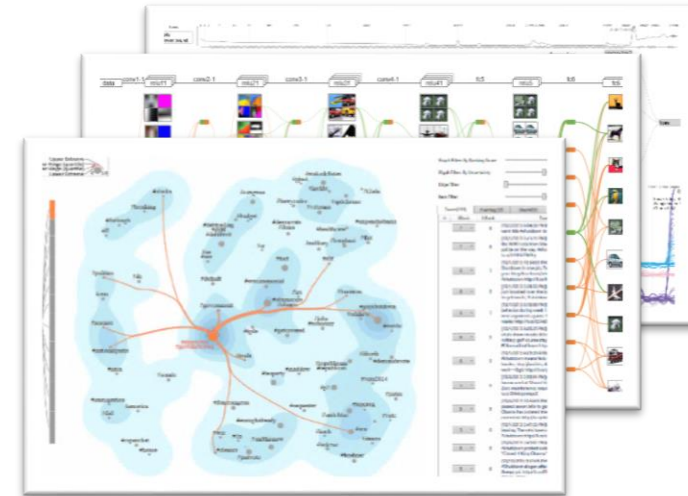
Discussion

- CNNVis cannot visualize deep models that cannot be formulated as DAGs
 - RNNs
- The scalability of activation matrix is limited
 - Number of columns = Number of classes
- There is a learning curve associated with the system
 - Neuron cluster and neuron

Future Work

- Provide an integrated system
 - End-to-end development

- Interactive feature selection
 - Manual feature construction / selection is still needed (tabular data)



Name	Thread pitch (mm)	Minor diameter tolerance	Nominal diameter (mm)	Head shape	Price for 50 screws	Available at factory outlet?	Number in stock	Flat or Phillips head?
M4	0.7	4g	4	Pan	\$10.08	Yes	276	Flat
M5	0.8	4g	5	Round	\$13.89	Yes	183	Both
M6	1	5g	6	Button	\$10.42	Yes	1043	Flat
M8	1.25	5g	8	Pan	\$11.98	No	298	Phillips
M10	1.5	6g	10	Round	\$16.74	Yes	488	Phillips
M12	1.75	7g	12	Pan	\$18.26	No	998	Flat
M14	2	7g	14	Round	\$21.19	No	235	Phillips
M16	2	8g	16	Button	\$23.57	Yes	292	Both
M18	2.1	8g	18	Button	\$25.87	No	664	Both
M20	2.4	8g	20	Pan	\$29.09	Yes	486	Both
M24	2.55	9g	24	Round	\$33.01	Yes	982	Phillips
M28	2.7	10g	28	Button	\$35.66	No	1067	Phillips
M36	3.2	12g	36	Pan	\$41.32	No	434	Both
M50	4.5	15g	50	Pan	\$44.72	No	740	Flat

Thank You!

Q & A