

Lunch & Learn:

features in convolutional neural nets

presentation by: Denis Kazakov

Rough Notes (read before looking at slides)

Slides

Key points

I am going to refer to brain's analogies and how humans learn purely for metaphorical purposes. I am not making any claims whatsoever that the two are anyhow identical

0. What did I actually work on

Road Lane classifier - too small of a dataset - noticed that it learned form context, not from roads

Traffic Light detectors - pipeline -

got to learn DNN's architectures, spark

concept of fusion

1. Lack of data => retrain

- we have to learn how to communicate what we are looking for and in images it's hard

imaging being born again and seeing a bunch of images of the world without knowing what the world is. It's hard to learn anything

- introducing new information as a way to compensate for lack of data

* eventually we are trying to teach DNN about object's invariance (we learn it by seeing world in 3d)

multiple loss functions

augmentations

basic

noise addition (read that paper)

- features are real

analogy to regression at each layer

about how 2 nets have partially distributed clusters of neurons in each layer

- pre-training as a way to transfer those features, should be generic - can't always rely on truth generation

introduce deep dreams and how steepest ascent works

show what the filters catch

show filters of not pertained -> hmmm (suspish) -> compare to large data basis -> hmmm

show examples of pertained and not

Conclusion: training from scratch always isn't easy. Too much data dependency (especially if there is a way to avoid it, worth exploring) all those features - are a complex nonlinear basis. intuitively whatever basis was found for 5mil images is more uniform for "seeing" that the basis it would find fro 50000 images

2. Doubtfulness as a job necessity

- overfitting isn't easy to catch

potential addition: pipeline to visualize network before using it in production.

if your test/val set isn't independent enough, you won't see it in the plot.

Ex: pertained model tested worse than not pertained. that's uncommon but just looking at the number doesn't tell us anything.

segway into adversarial networks?

3. Sidetones & thoughts on the field:

- **adversarial networks - have smth \Leftrightarrow relation to manifold theory?**

- **NN have many equivalent local minimas**

Agenda

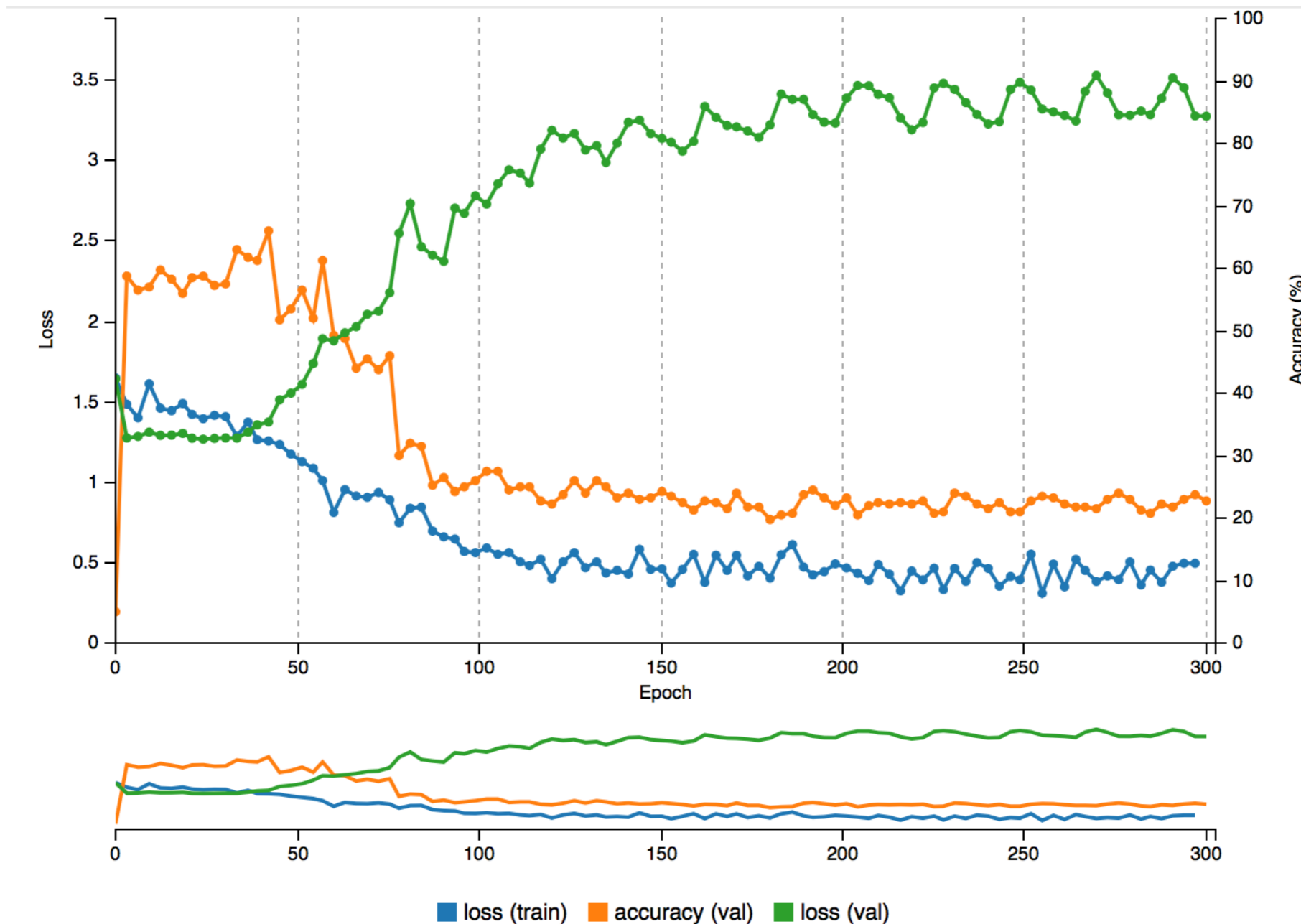
- 1) My summer experience, things I did
- 2) Addressing the problem of insufficient data
- 3) Transferring features in DNN
- 4) Visualizing layer maximization

Summer Internship

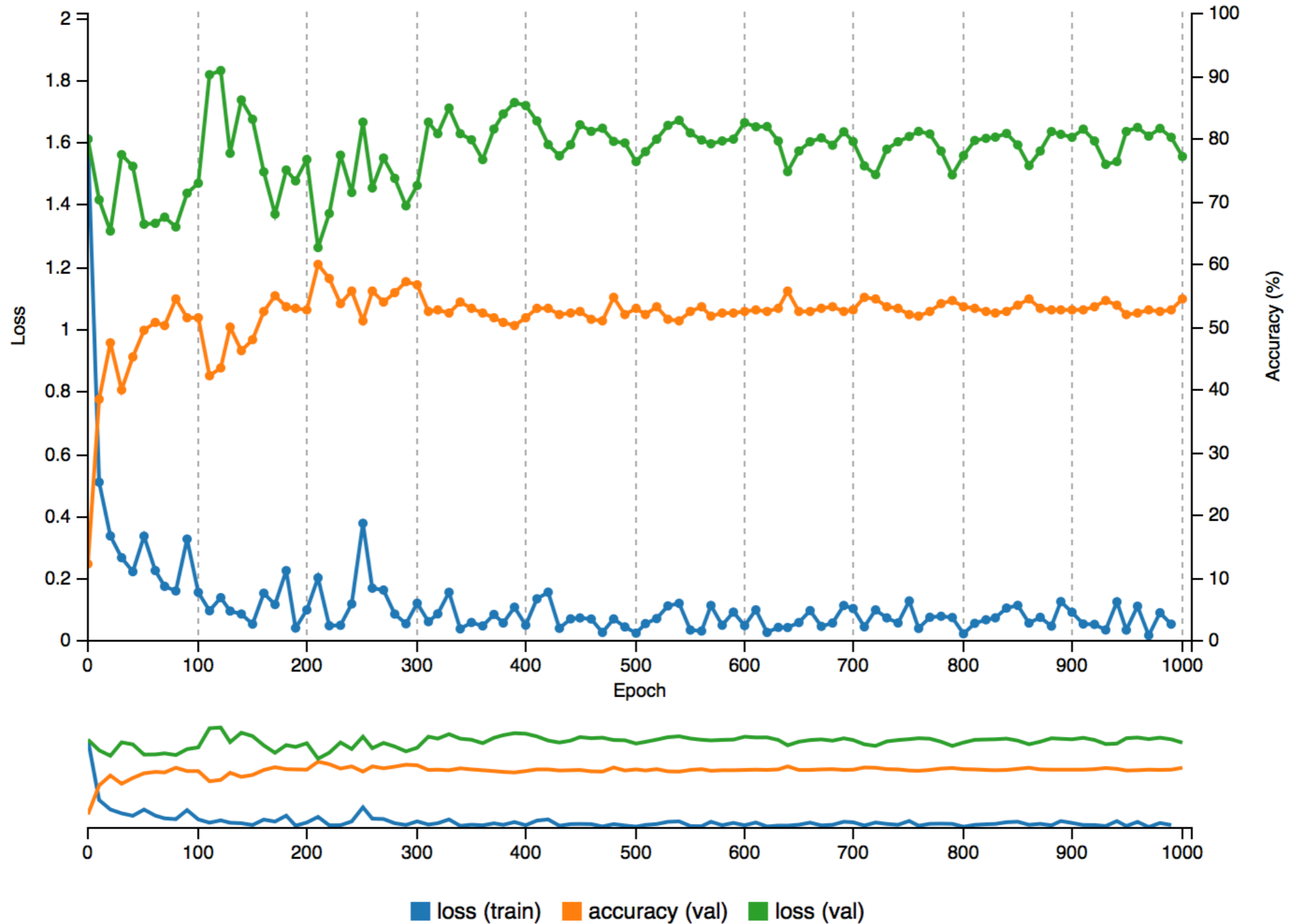
- 1) Road lane classifier
- 2) Traffic light object detector

Road Lane Classifier

Using standard AlexNet



Using pre-trained DNN



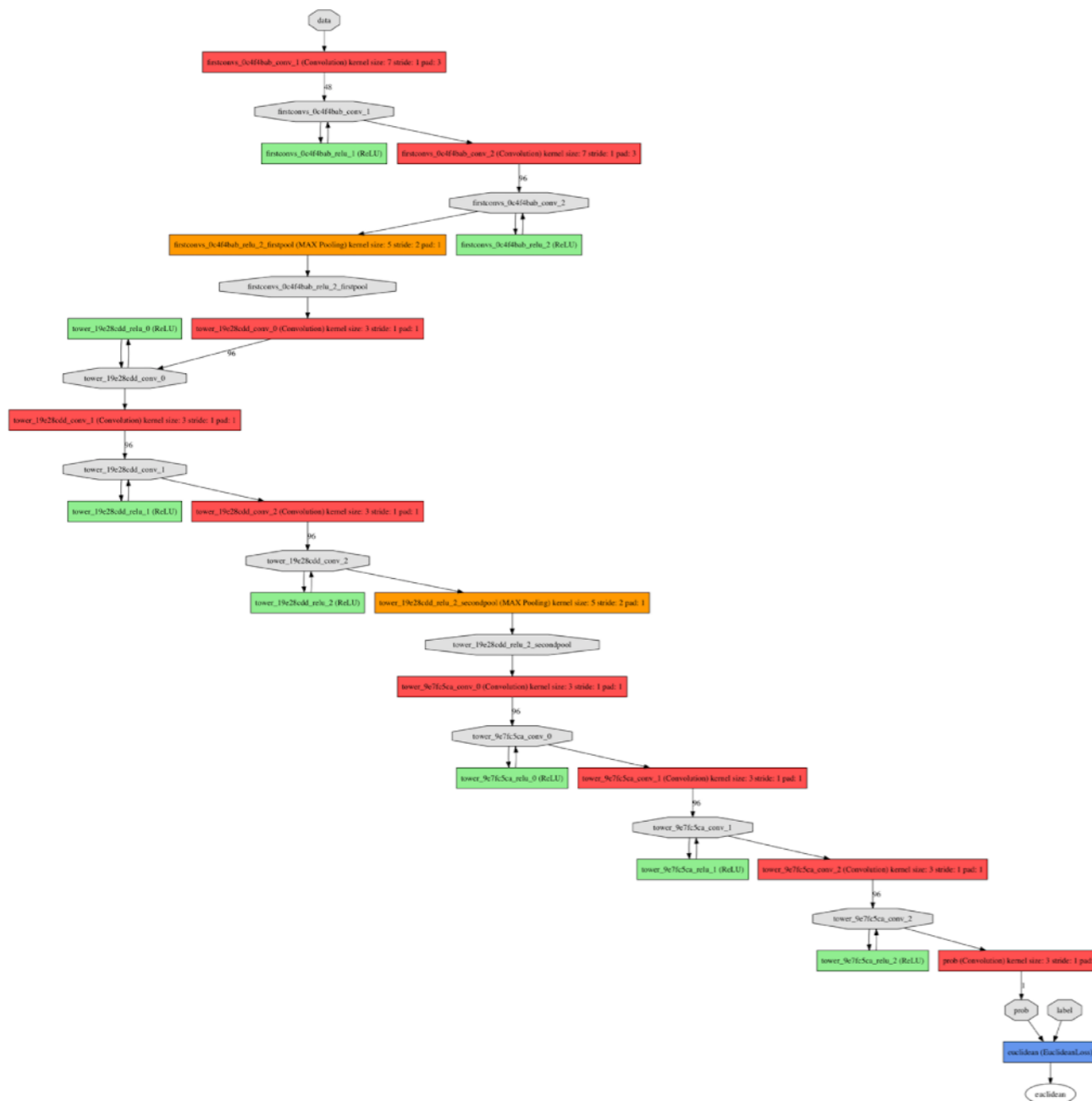
Conclusion on Road Lanes



Learned context, but didn't care about the road lanes themselves

Traffic Light Detector

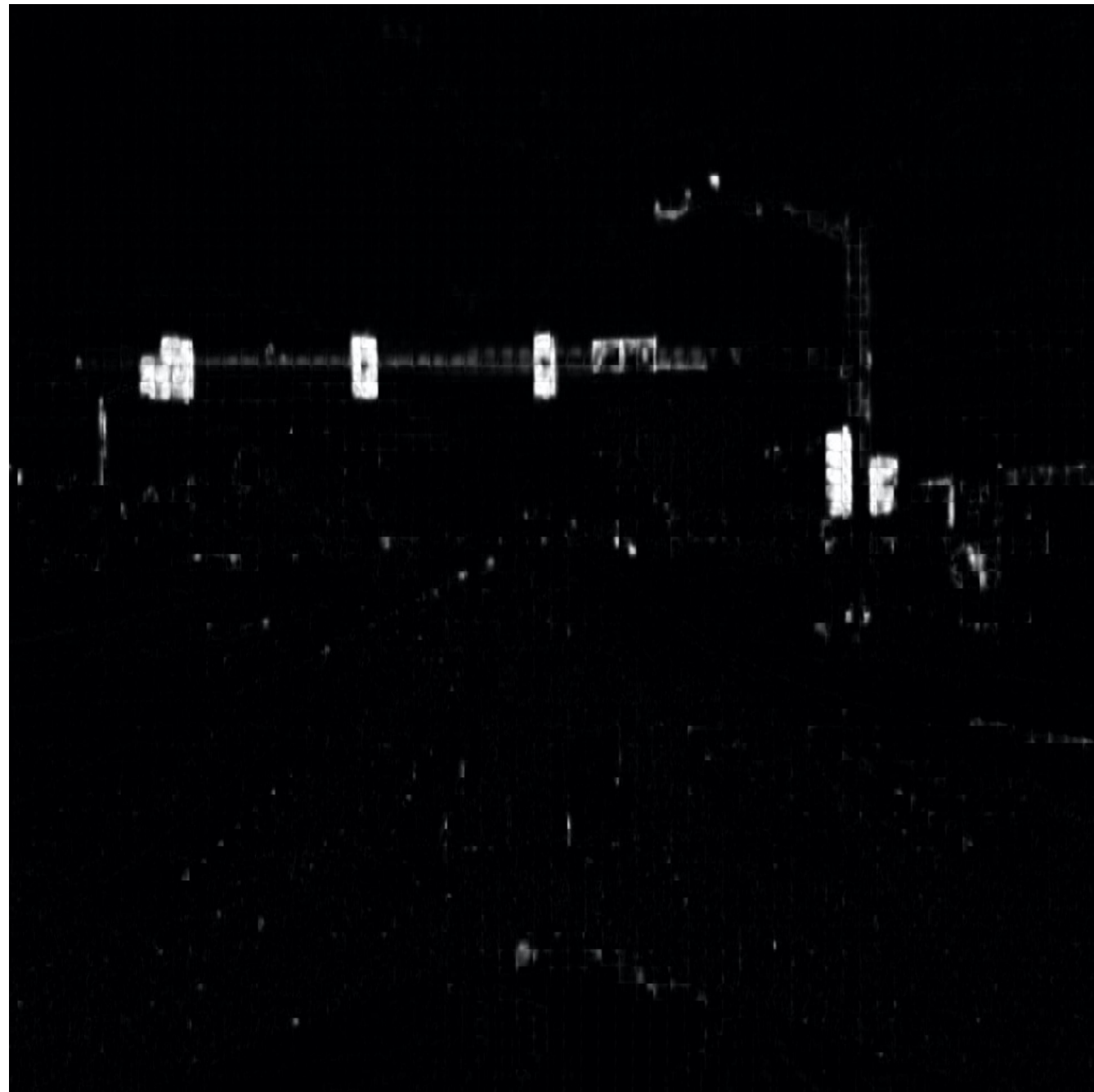
Simple detector: architecture



input image \longrightarrow heat map

- A. 2 CONV (7x7, stride 1)
- B. 1 POOL (5x5, stride 2)
- C. 3 CONV (3x3, stride 1)
- D. 1 POOL (5x5, stride 2)
- E. 3 CONV (3x3, stride 1)
- F. 1 heat map

Sample response (close distance)



Sample response (medium distance)



Sample response (far distance)



(deleted for IP protection)

Summer Conclusion

- 1) Learned Spark
- 2) Learned industry deep learning practices
- 3) Read a ton of papers and learned a ton about neural nets

Lunch & Learn



Training DNN is hard

randomly
initialized
filters

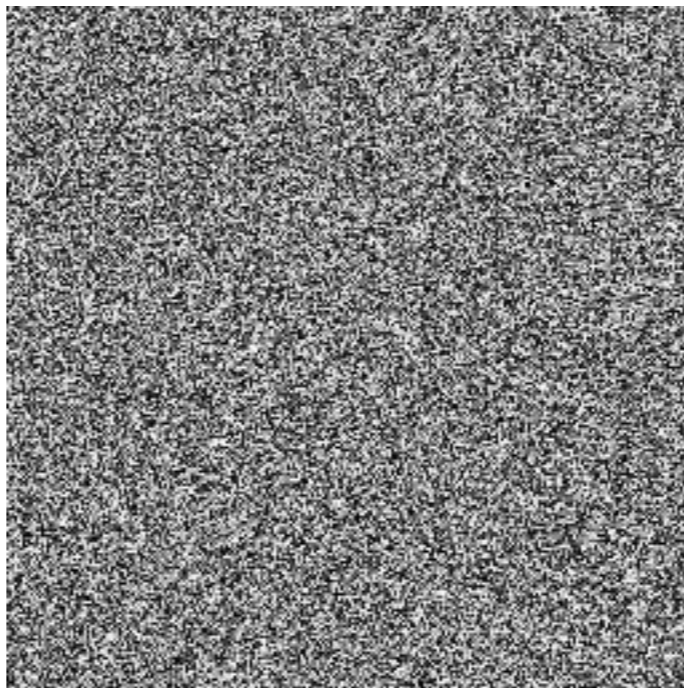


sensible
interpretation
of images



Training DNN is hard

Here, you see a traffic light



Okay...



Invariance in objects



Creating more information

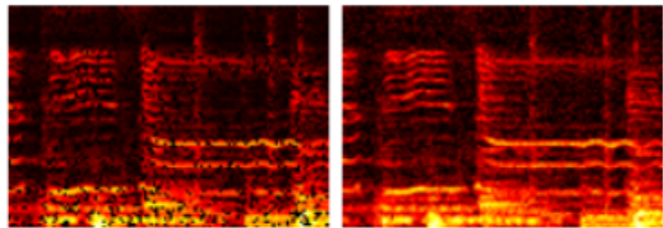
Two ways to teach invariance:

- 1) show examples from different contexts where something we are interested in stays constant
- 2) artificially 'augment' data to simulate different context

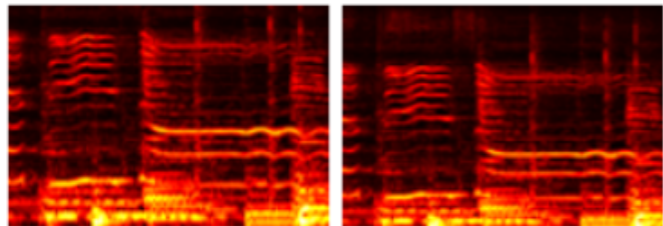
Data augmentation in images

Random(
Crop,
Scale,
Sheer,
Stretch,
Flip,
Noise)

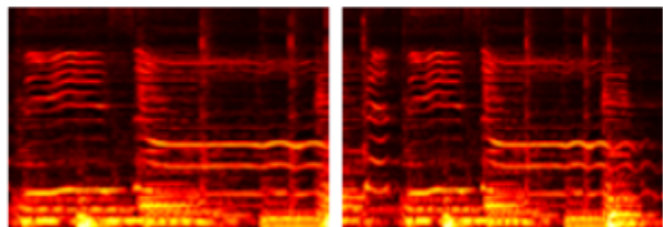




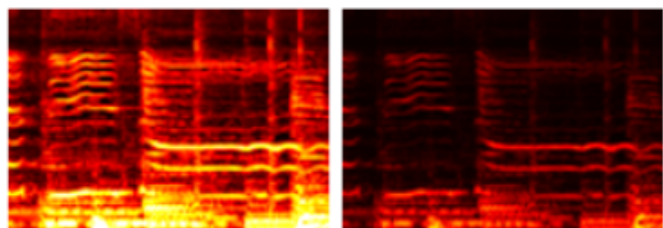
(c) Dropout and Gaussian noise.



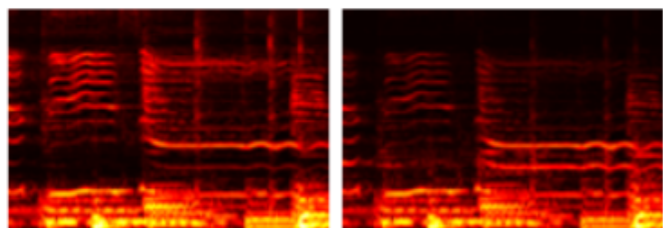
(d) Pitch shift of $\pm 20\%$.



(e) Time stretch of $\pm 20\%$.



(f) Loudness of ± 10 dB.



(g) Random frequency filters.

Data augmentation in sound

CONVERGENT LEARNING: DO DIFFERENT NEURAL NETWORKS LEARN THE SAME REPRESENTATIONS?

Yixuan Li^{1*}, Jason Yosinski^{1*}, Jeff Clune², Hod Lipson³, & John Hopcroft¹

¹Cornell University

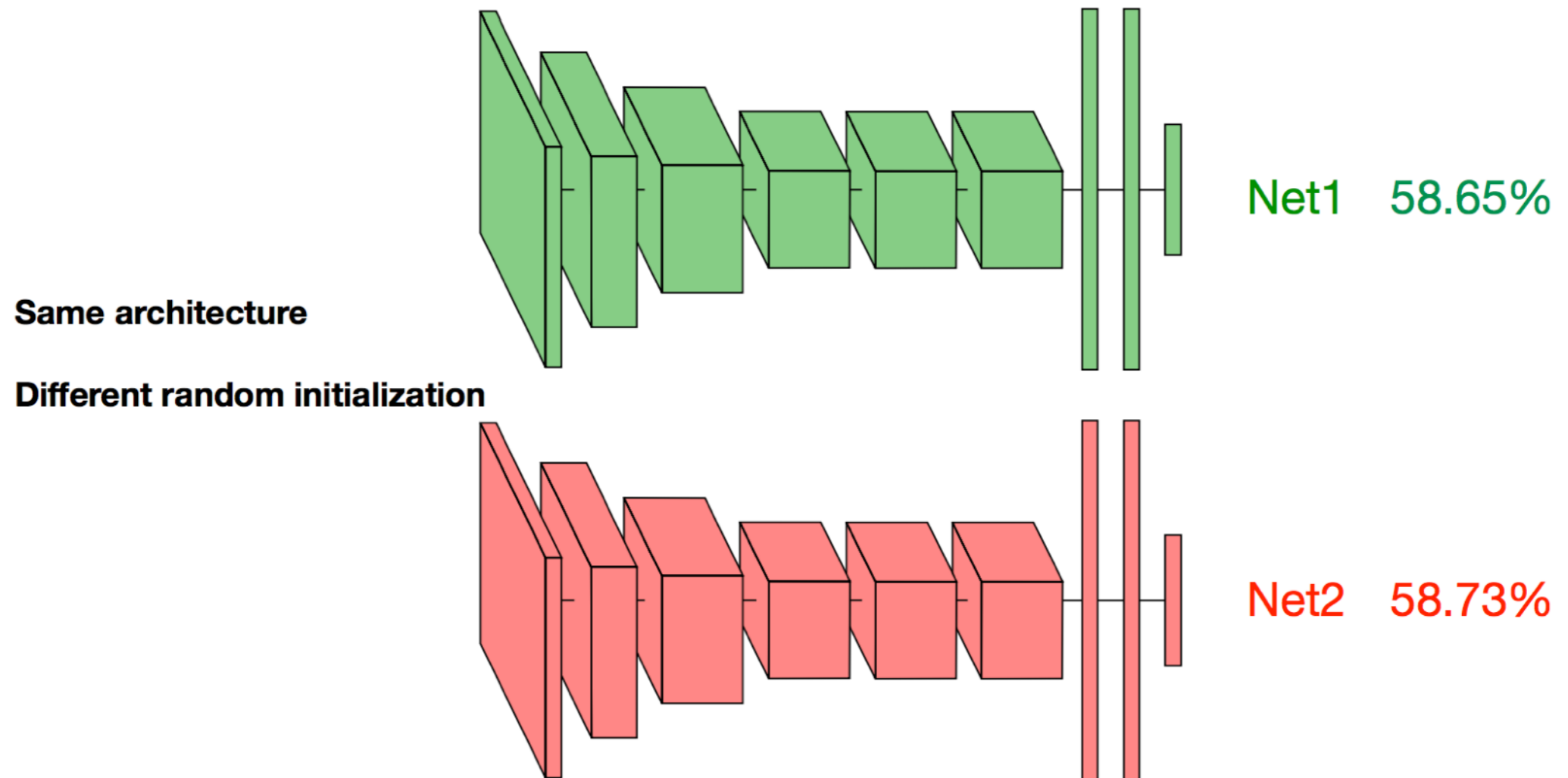
²University of Wyoming

³Columbia University

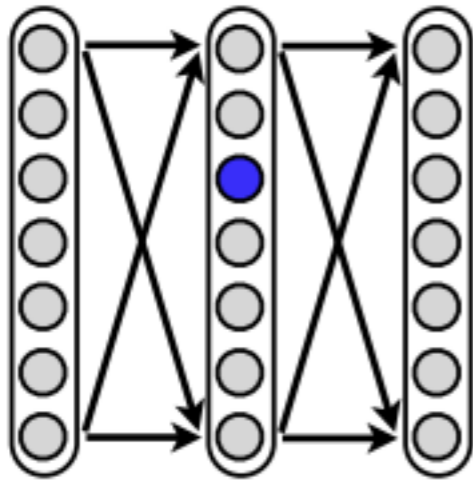
{yli, yosinski, jeh}@cs.cornell.edu

jeffclune@uwyo.edu, hod.lipson@columbia.edu

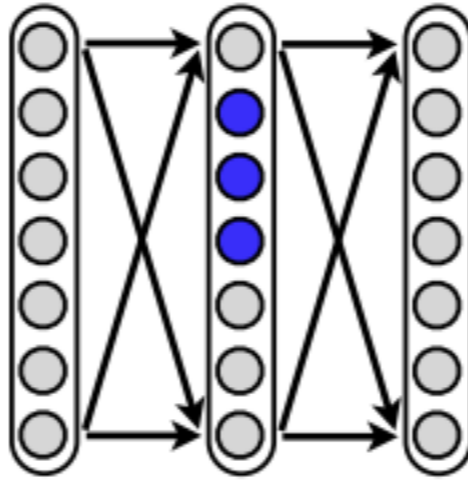
Features



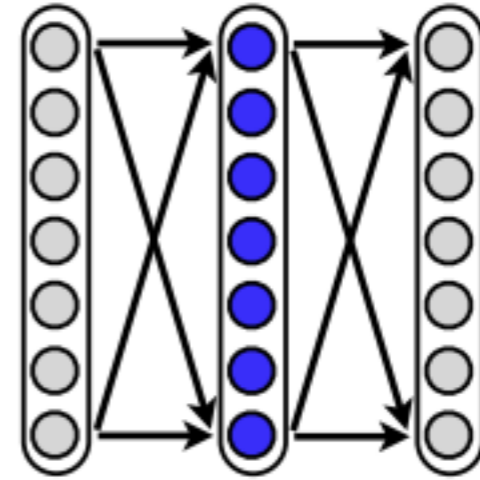
Which one is it?



Local?



Partially-
Distributed?



Distributed?

Measuring correlation & mutual information

i, j - units, l - layer, X - series of activation values

Mean: $\mu_{l,i}^{(n)} = \mathbb{E}[X_{l,i}^{(n)}]$

Standard deviation: $\sigma_{l,i}^{(n)} = \sqrt{\mathbb{E}[(X_{l,i}^{(n)} - \mu_{l,i}^{(n)})^2]}$

Within-net correlation: $c_{l,i,j}^{(n)} = \mathbb{E}[(X_{l,i}^{(n)} - \mu_{l,i}^{(n)})(X_{l,j}^{(n)} - \mu_{l,j}^{(n)})] / \sigma_{l,i}^{(n)} \sigma_{l,j}^{(n)}$

Between-net correlation: $c_{l,i,j}^{(n,m)} = \mathbb{E}[(X_{l,i}^{(n)} - \mu_{l,i}^{(n)})(X_{l,j}^{(m)} - \mu_{l,j}^{(m)})] / \sigma_{l,i}^{(n)} \sigma_{l,j}^{(m)}$

Mutual information: $I(X_{l,i}^{(n)}; X_{l,j}^{(m)}) = \sum_{a \in X_{l,i}^{(n)}} \sum_{b \in X_{l,j}^{(m)}} p(a, b) \log \left(\frac{p(a, b)}{p(a)p(b)} \right),$

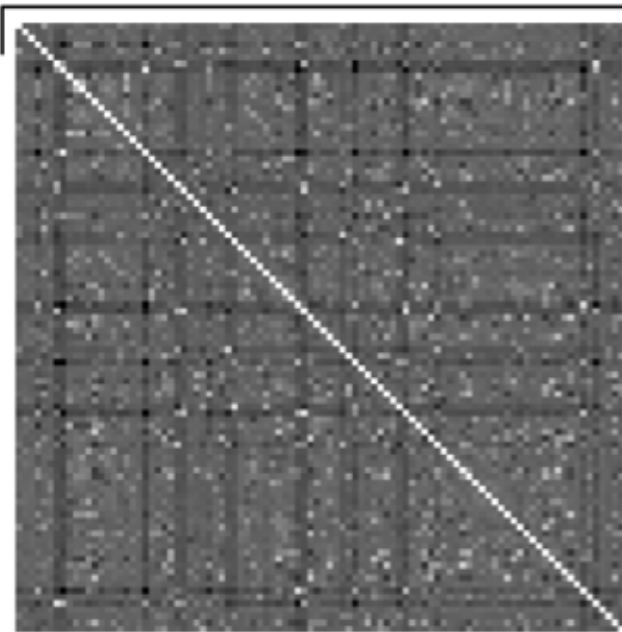
Neuron - Neuron Correlation

$c_{\text{conv1}}^{(\text{Net1})}$

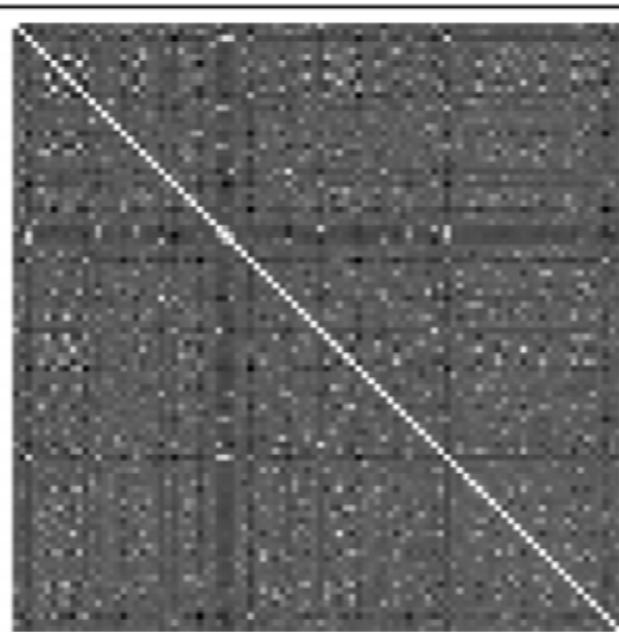
within-net

$c_{\text{conv1}}^{(\text{Net2})}$

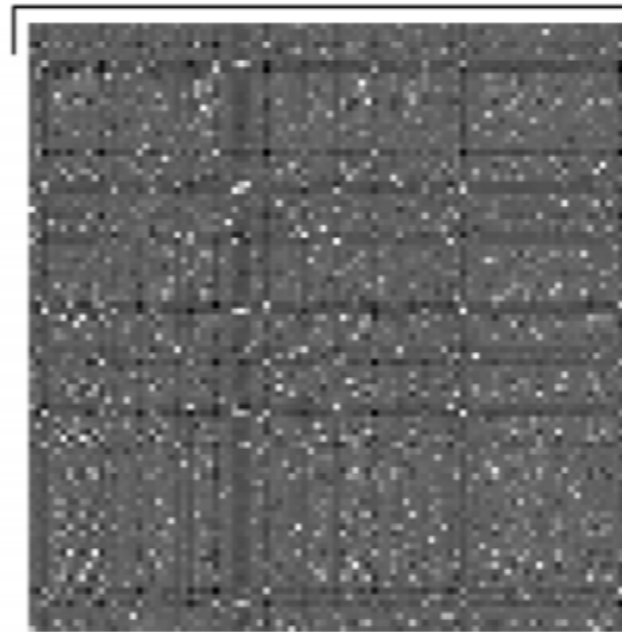
between-net $c_{\text{conv1}}^{(\text{Net1,Net2})}$



(a)

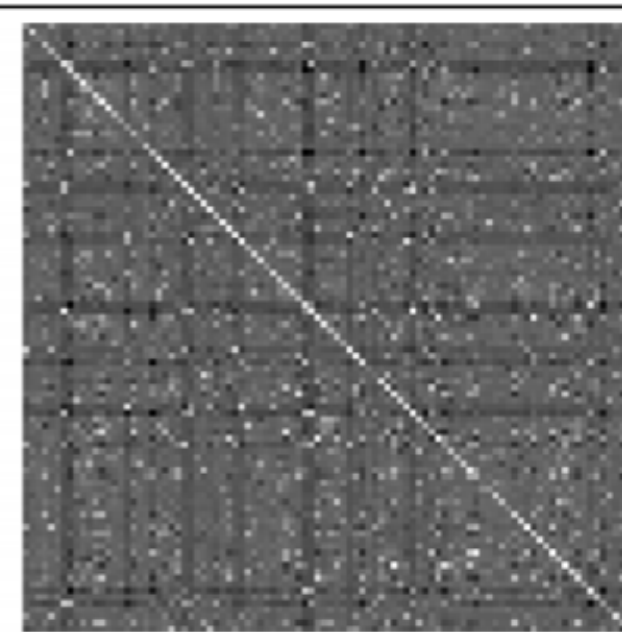


(b)



natural (unaligned) order

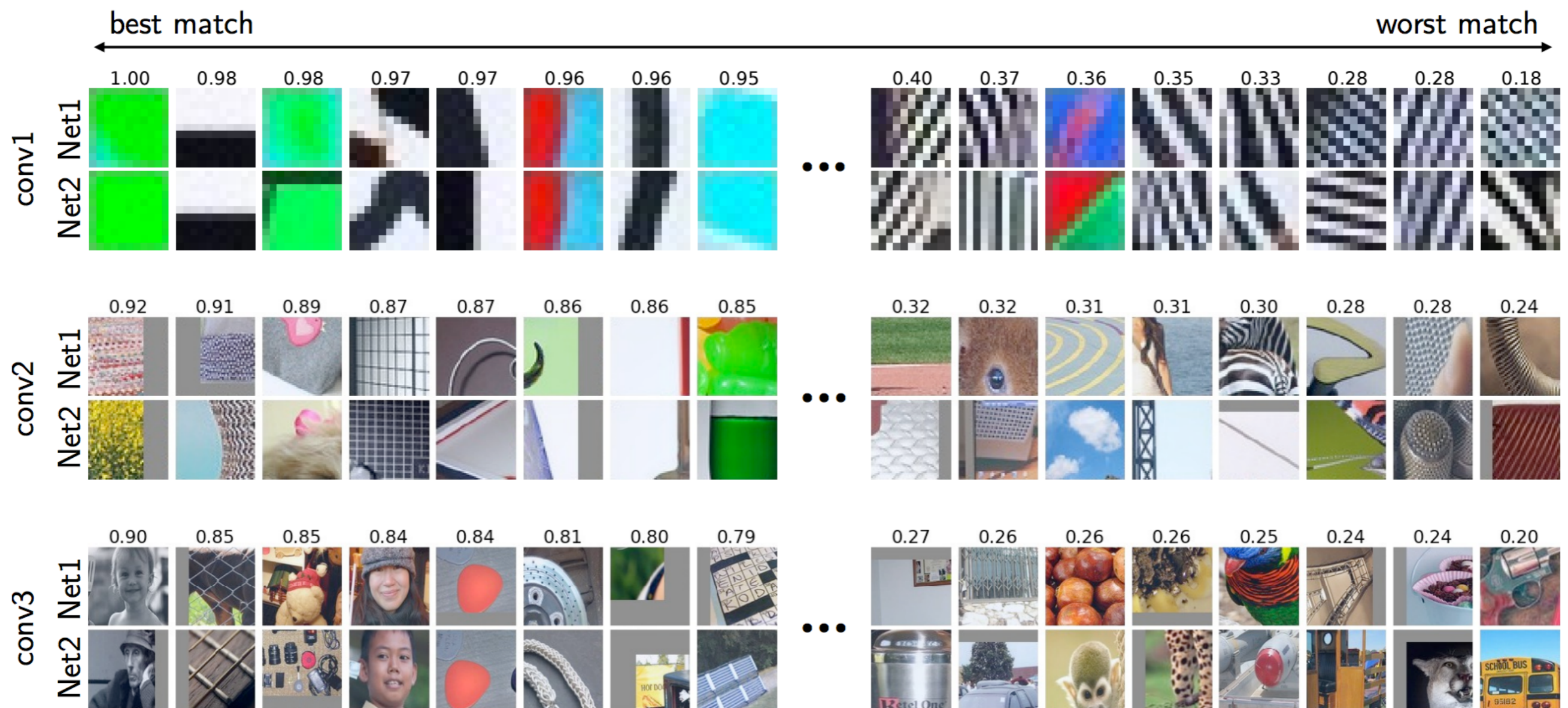
(c)



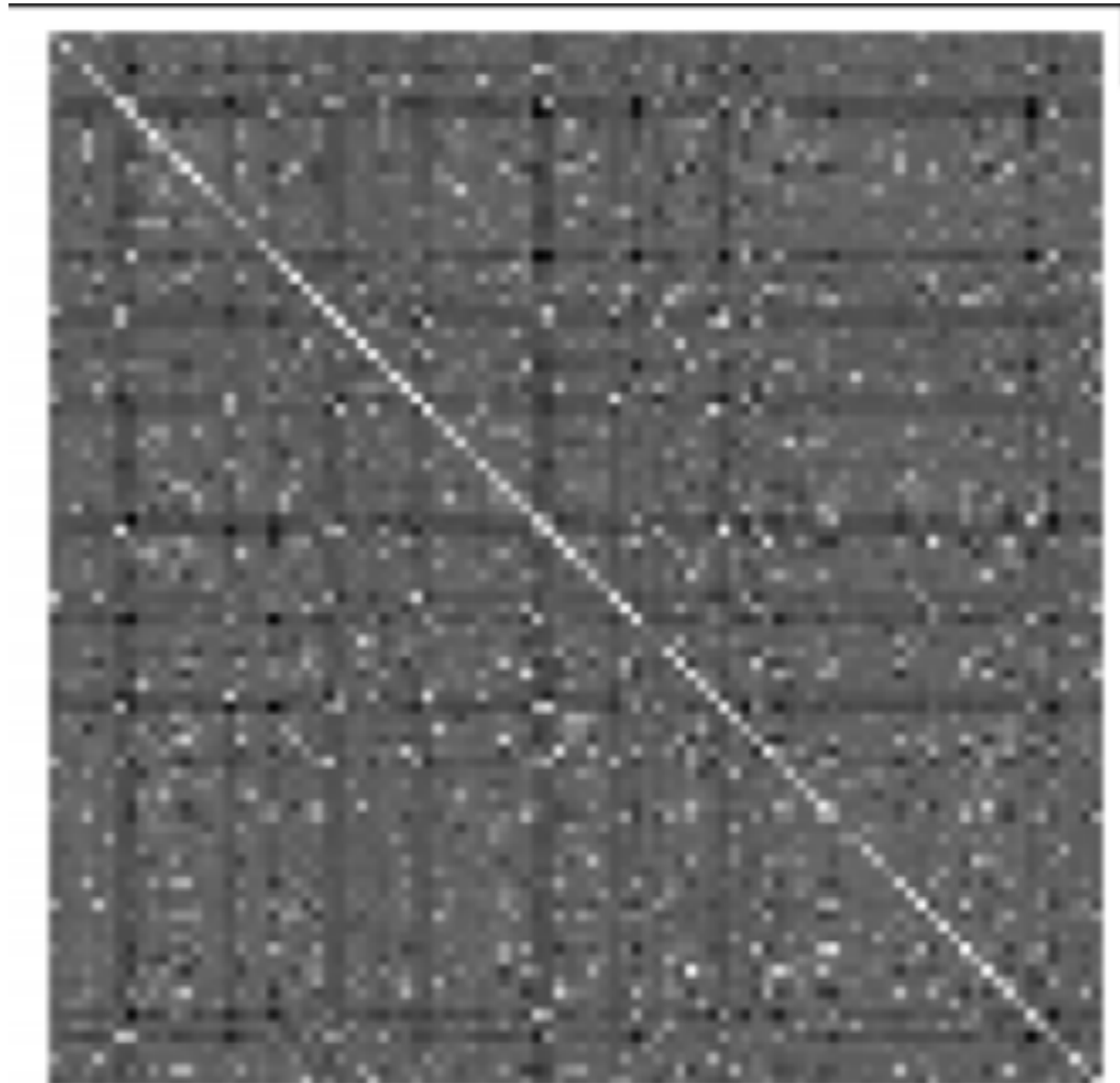
permuted (aligned) order

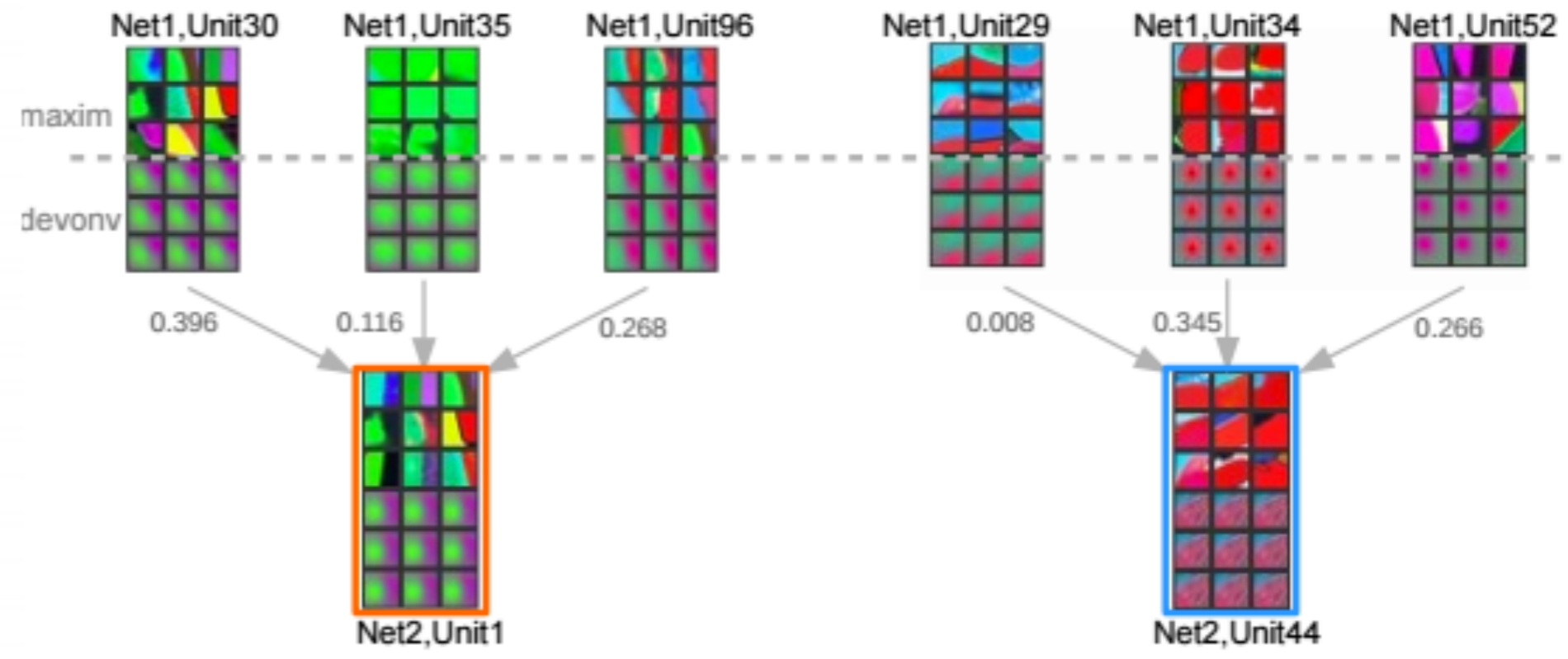
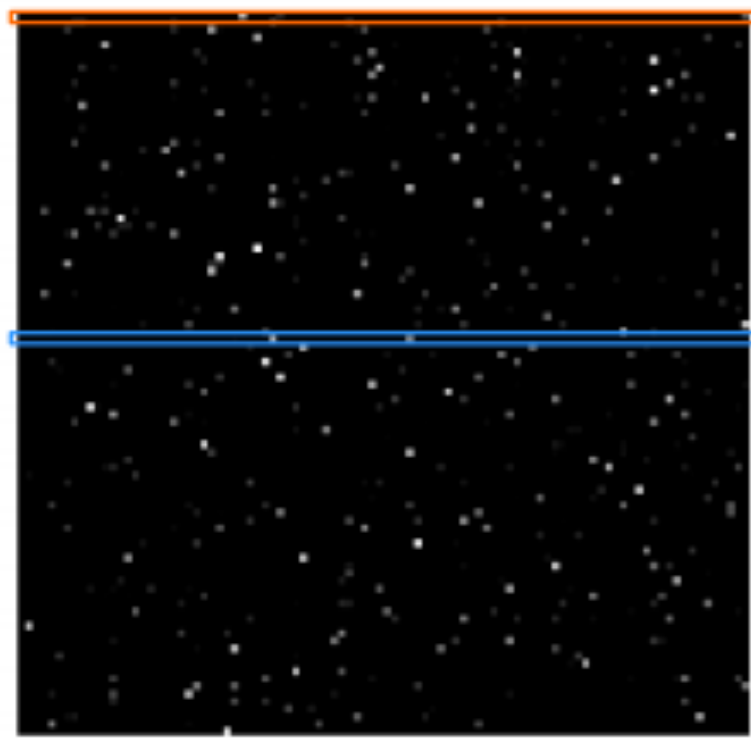
(d)

Common filters - chosen greedily

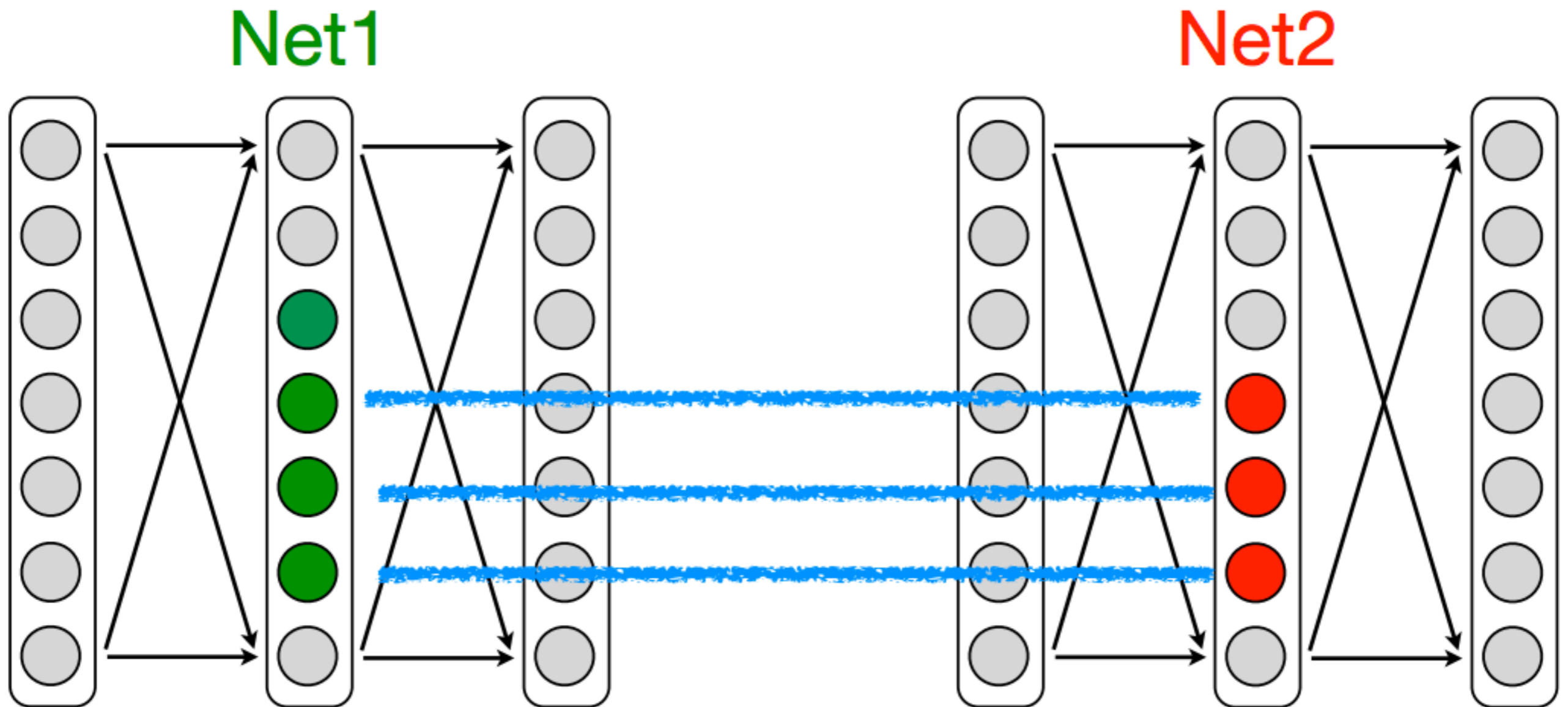


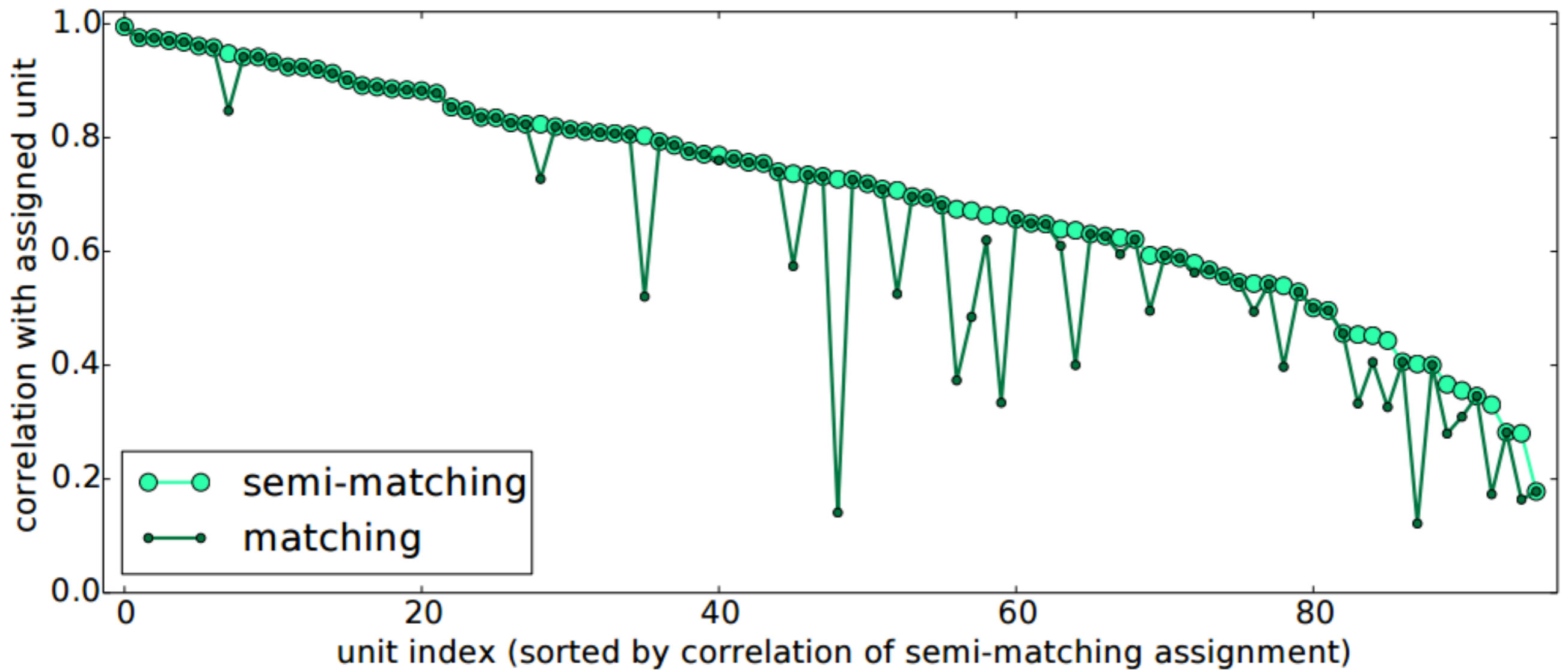
What about non diagonal entries?



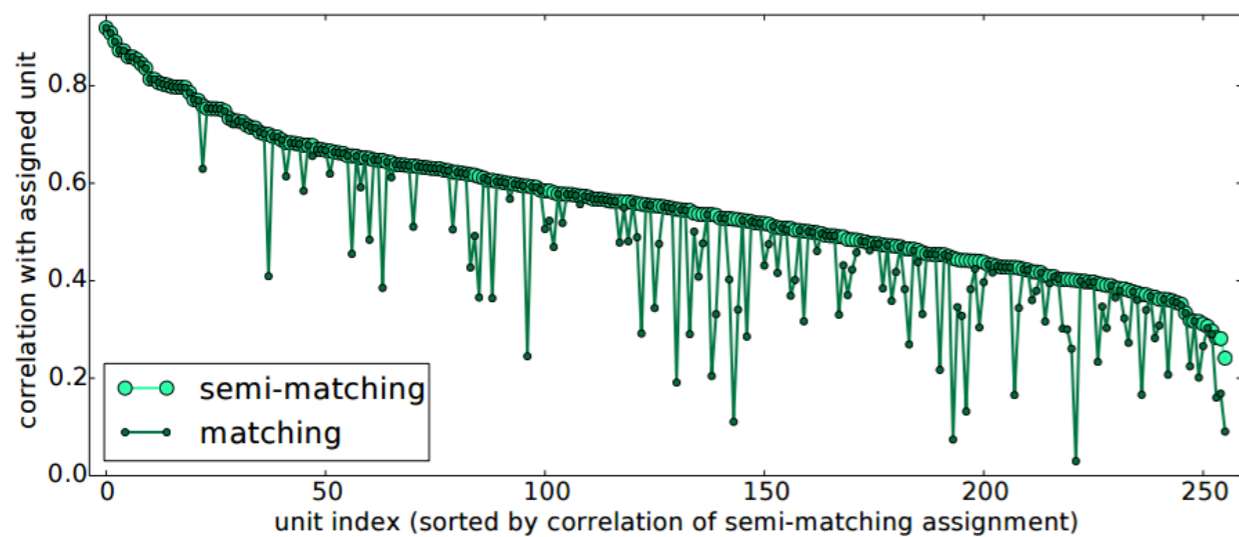


Common filters - chosen by weighted bipartite matching

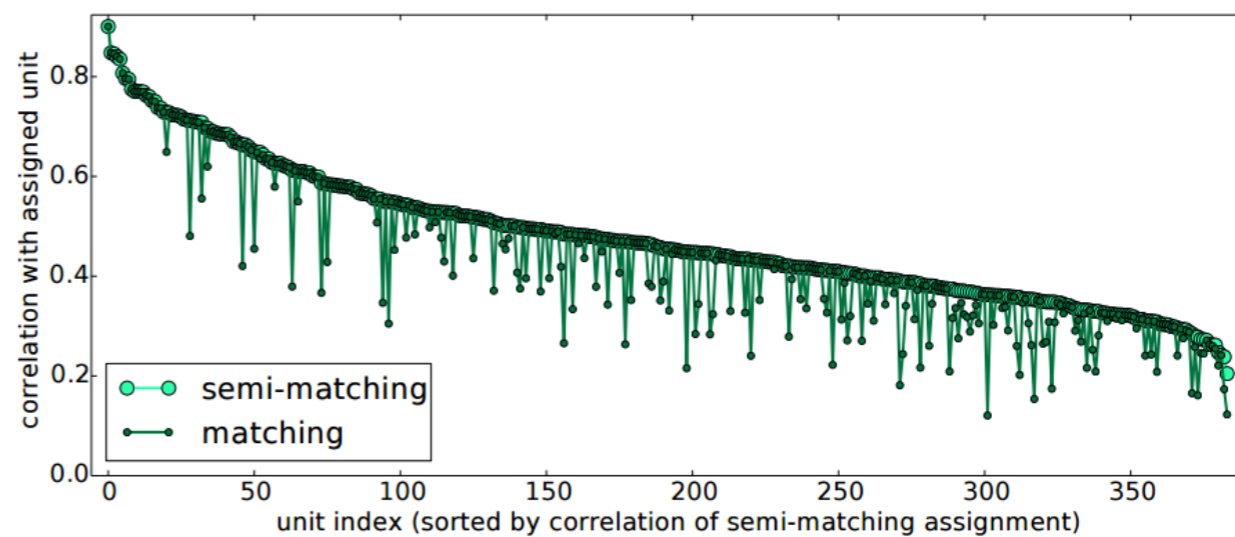




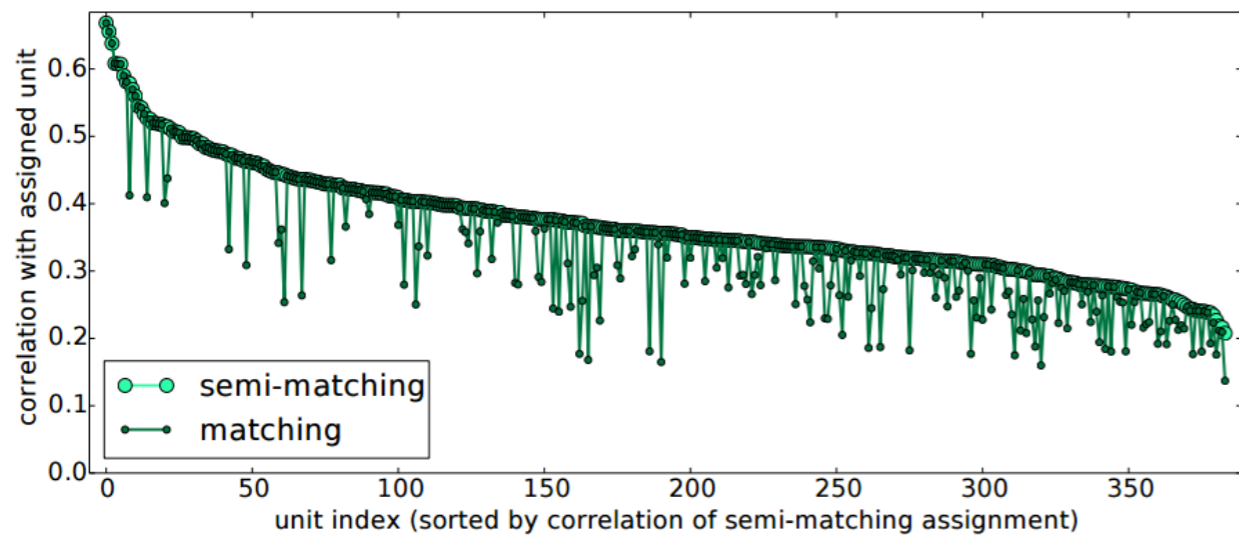
Takeaway: Some units didn't get a match



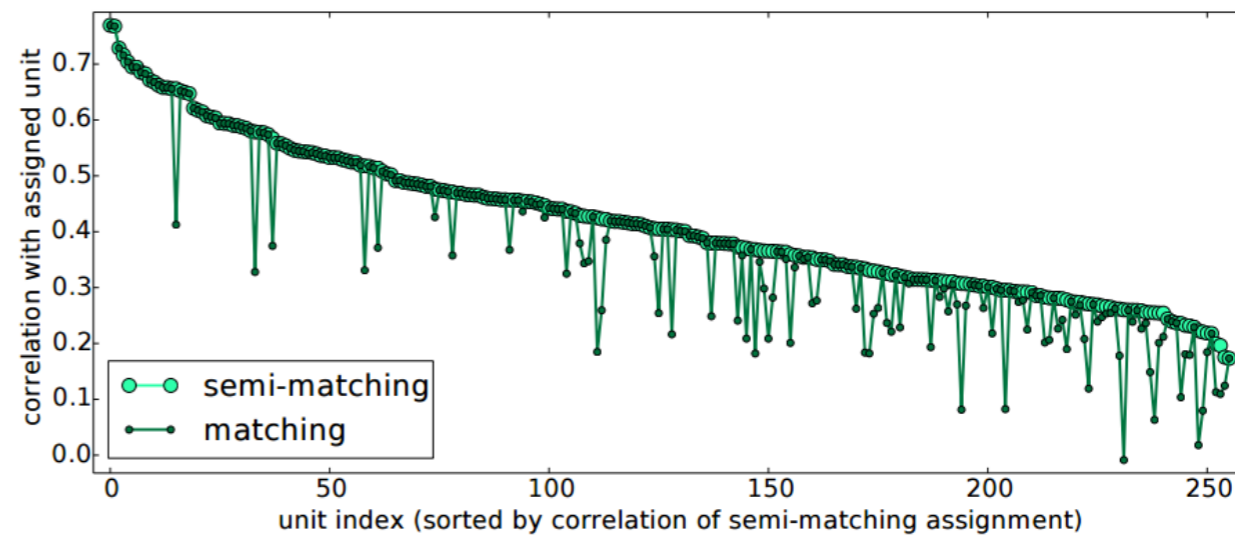
(a) conv2



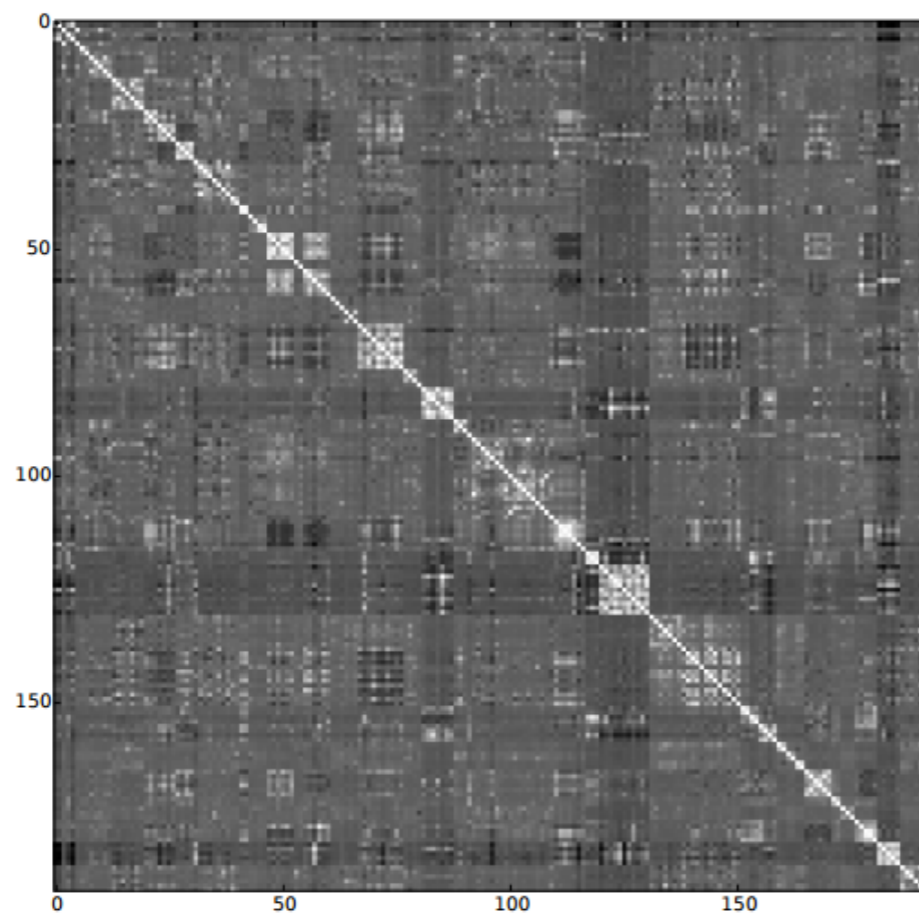
(b) conv3



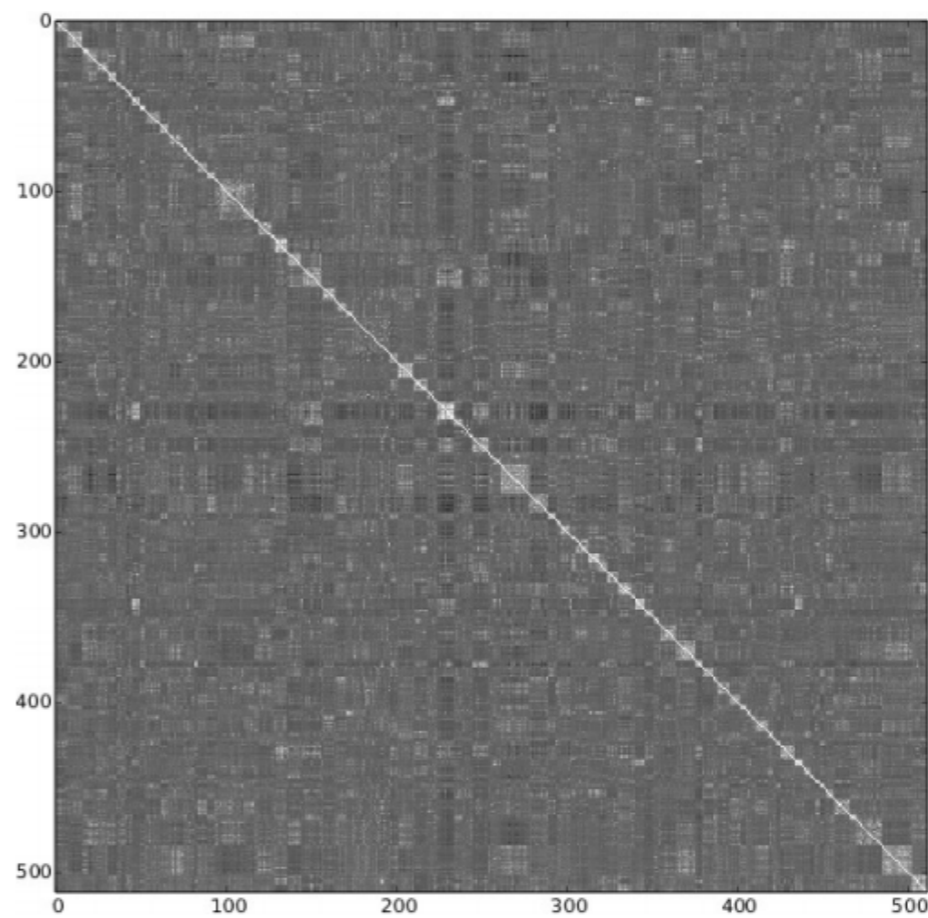
(c) conv4



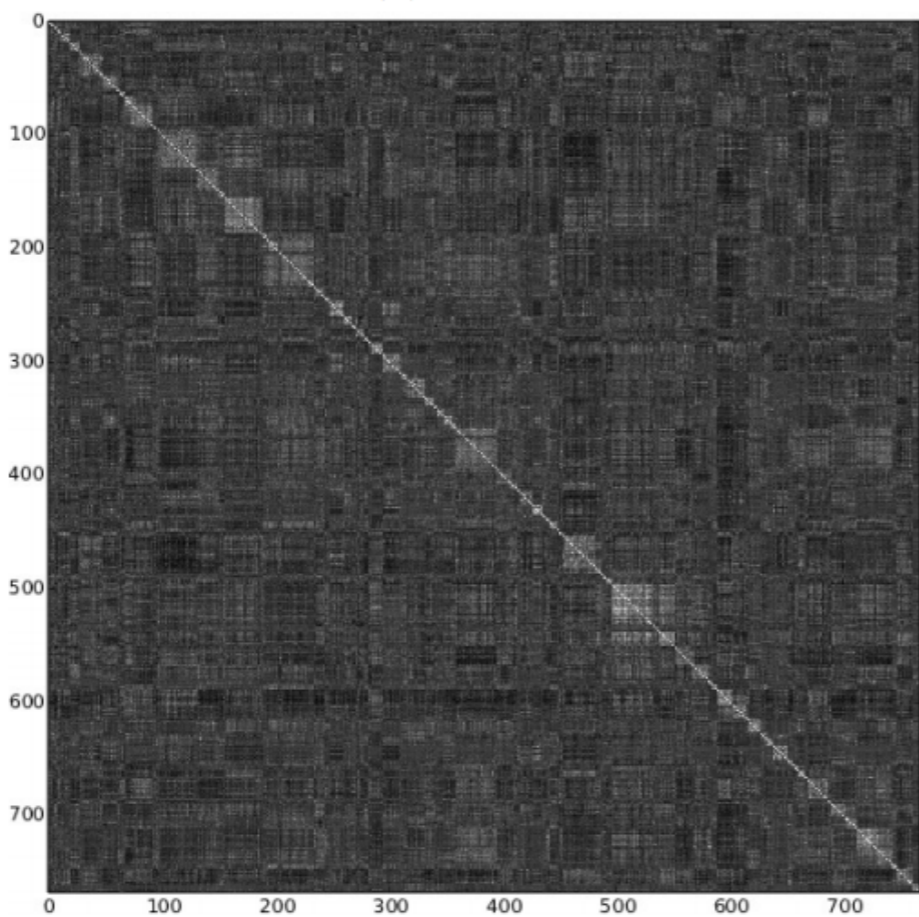
(d) conv5



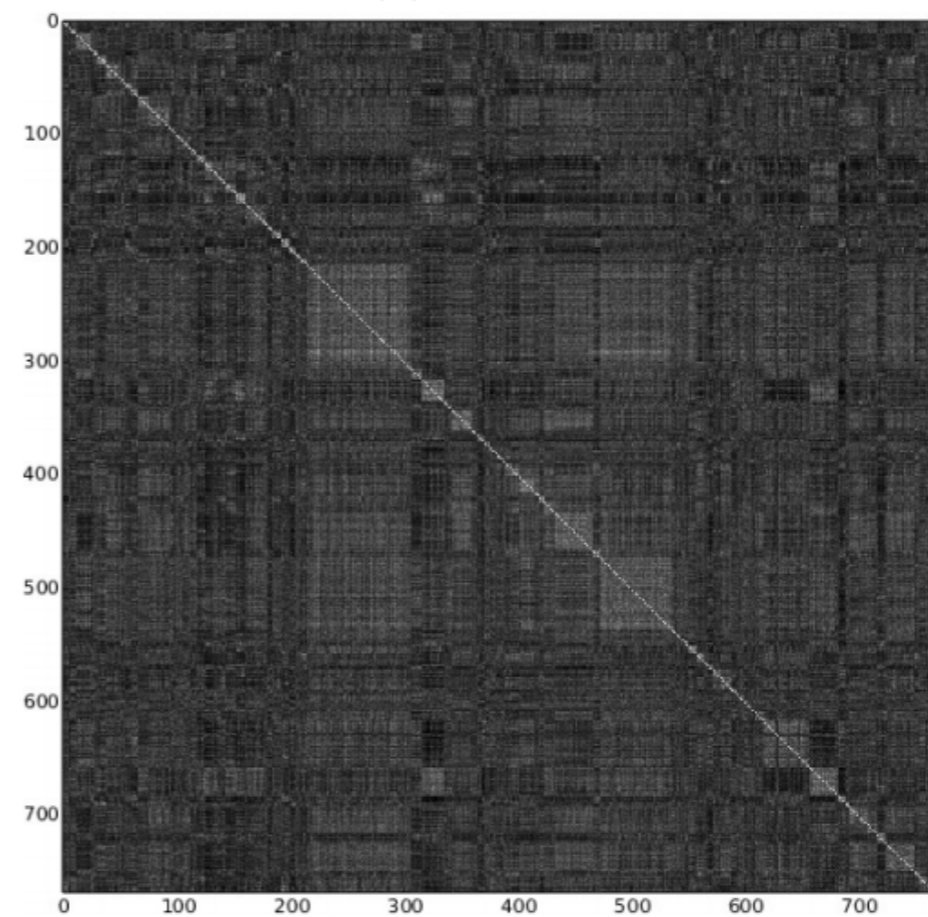
(a) conv1



(b) conv2



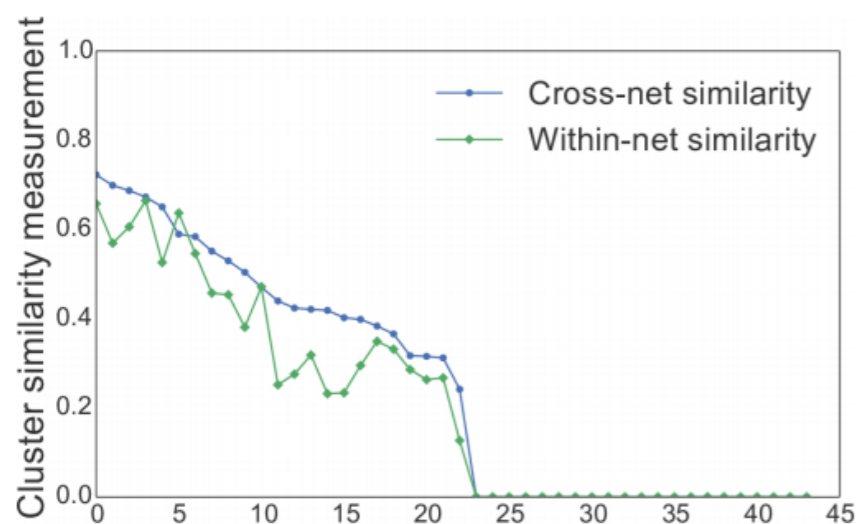
(c) conv3



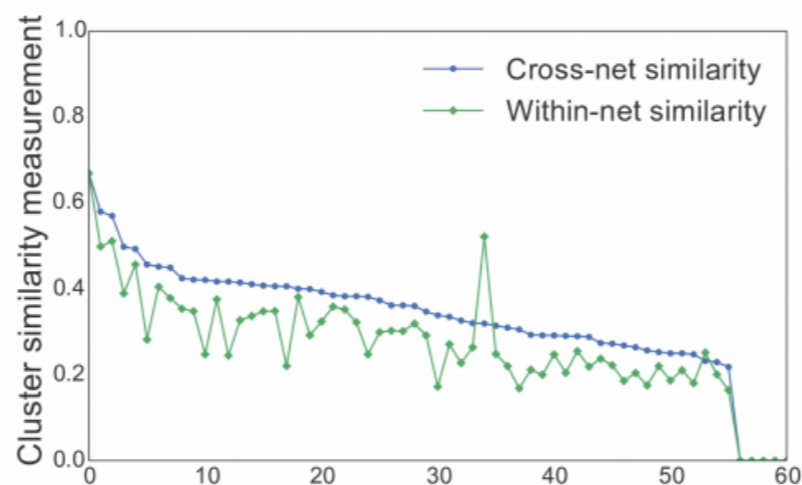
(d) conv4

Between-net similarity: $\text{Sim}_{X_l \rightarrow Y_l} = \frac{\sum_{p=1}^{S_l} \sum_{q=1}^{S_l} \text{corr}(X_l, Y_l)_{pq}}{S_l^2}$

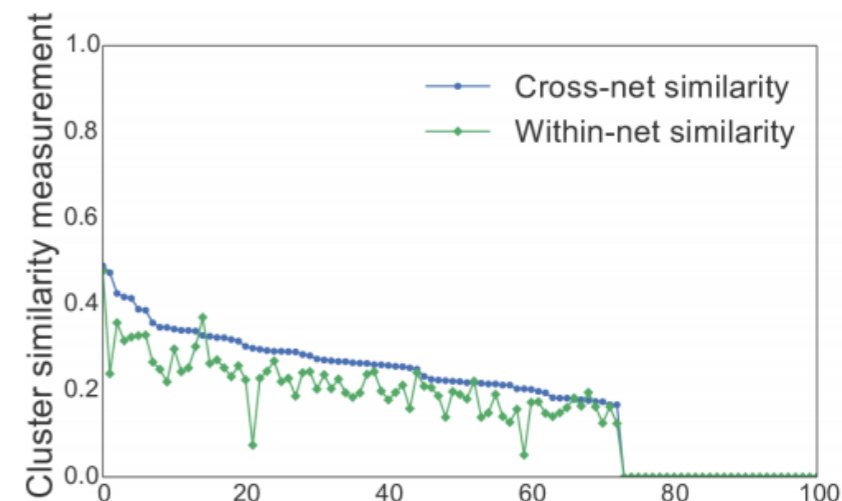
Within-net similarity: $\text{Sim}_{X_l, Y_l} = (\text{Sim}_{X_l \rightarrow X_l} + \text{Sim}_{Y_l \rightarrow Y_l})/2$



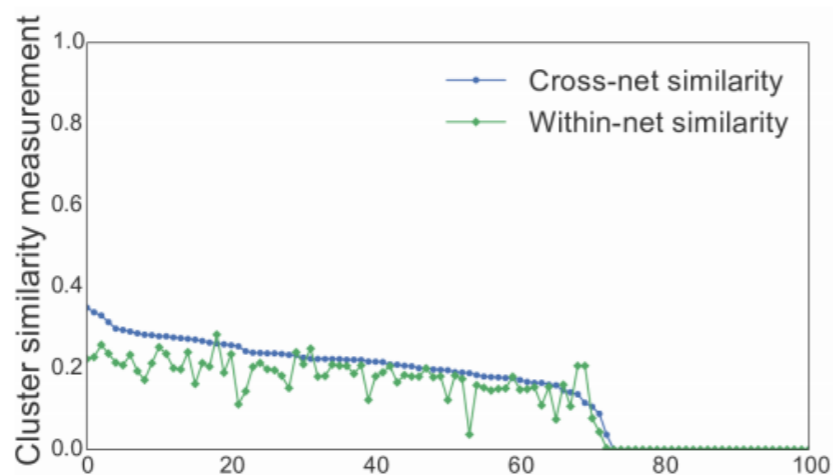
(a) conv1



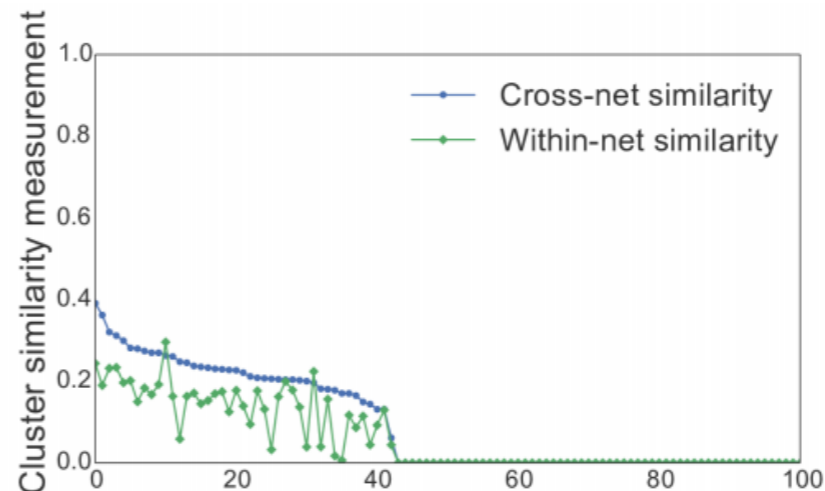
(b) conv2



(c) conv3



(d) conv4



(e) conv5

Figure S9: The distribution of between-net and within-net similarity measurement after clustering neurons (conv1 – conv5). The x-axis represents obtained clusters, which is reshuffled according to the sorted between-net similarity value.

Transferring features by
'pre-training'

*side knowledge: layer maximization

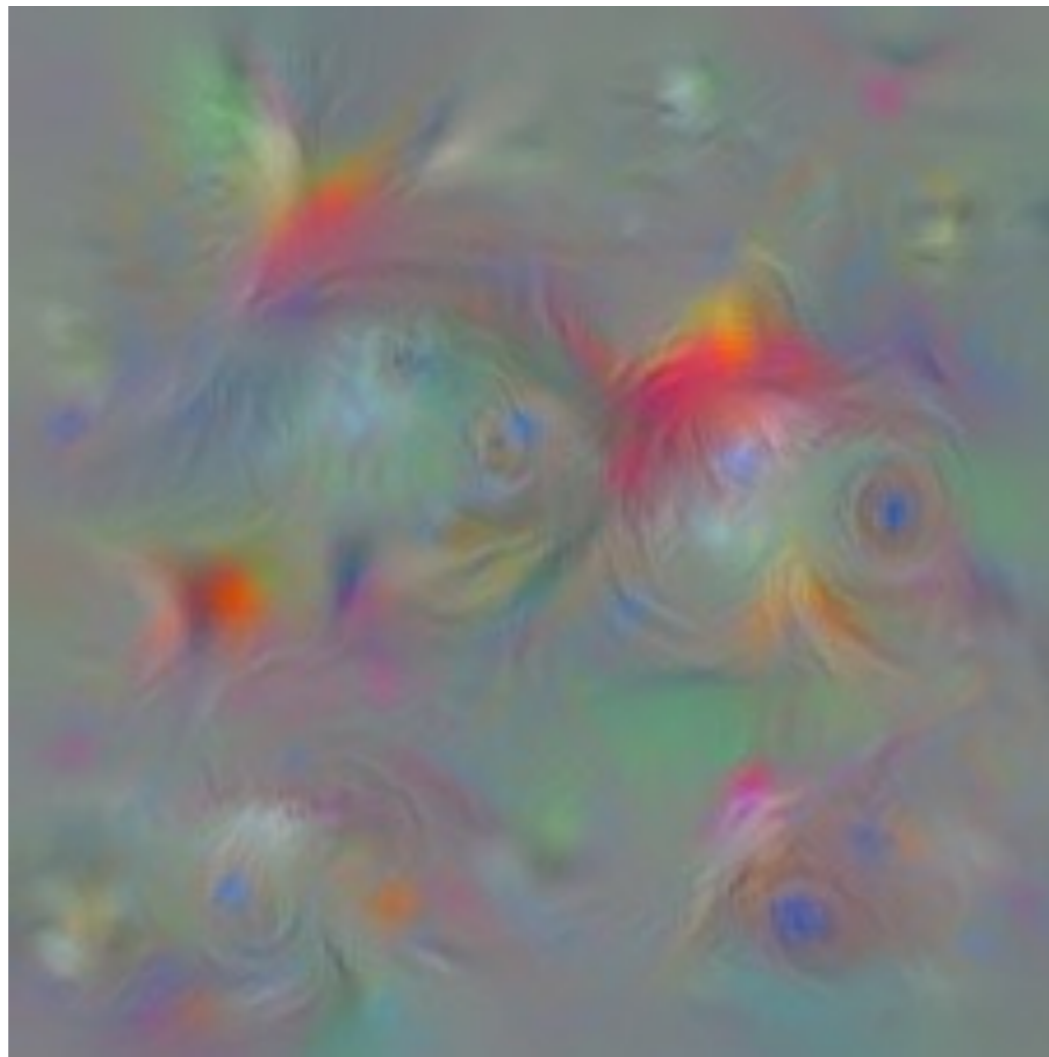
1) Input random noise

2) Maximize (steepest ascent) a chosen layer with respect to input

Result: input image that will maximize a chosen layer (aka we see what the layer is “expecting” to see)



Examples



goldfish

Why useful?

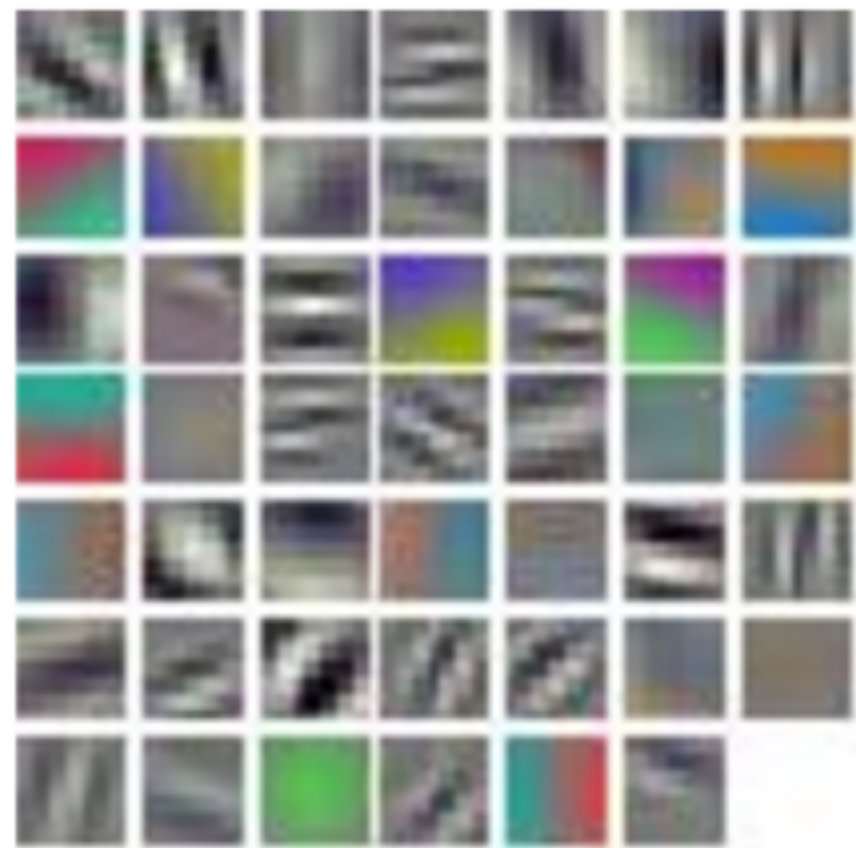


Traffic Light DNN classifier

Conv 1 filters

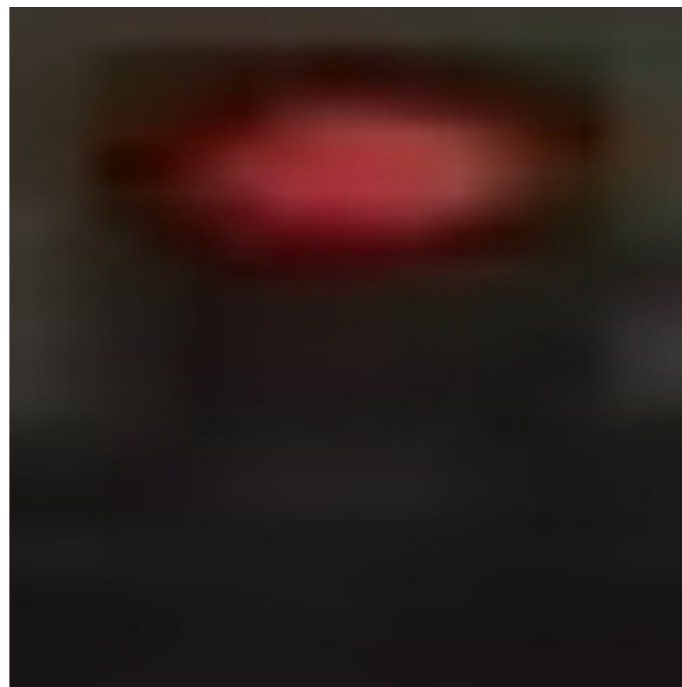


from scratch

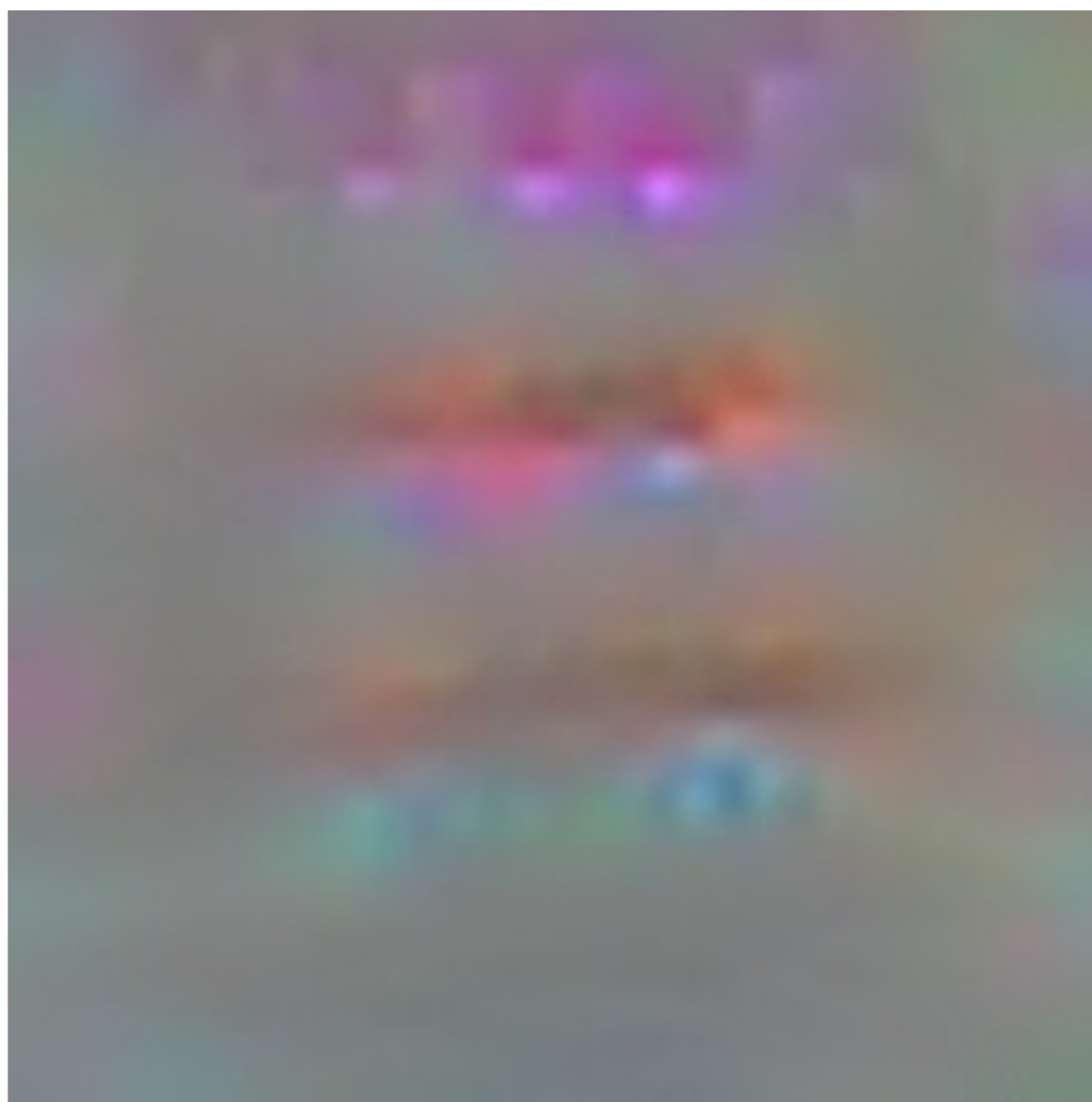


pre-trained

Data Example



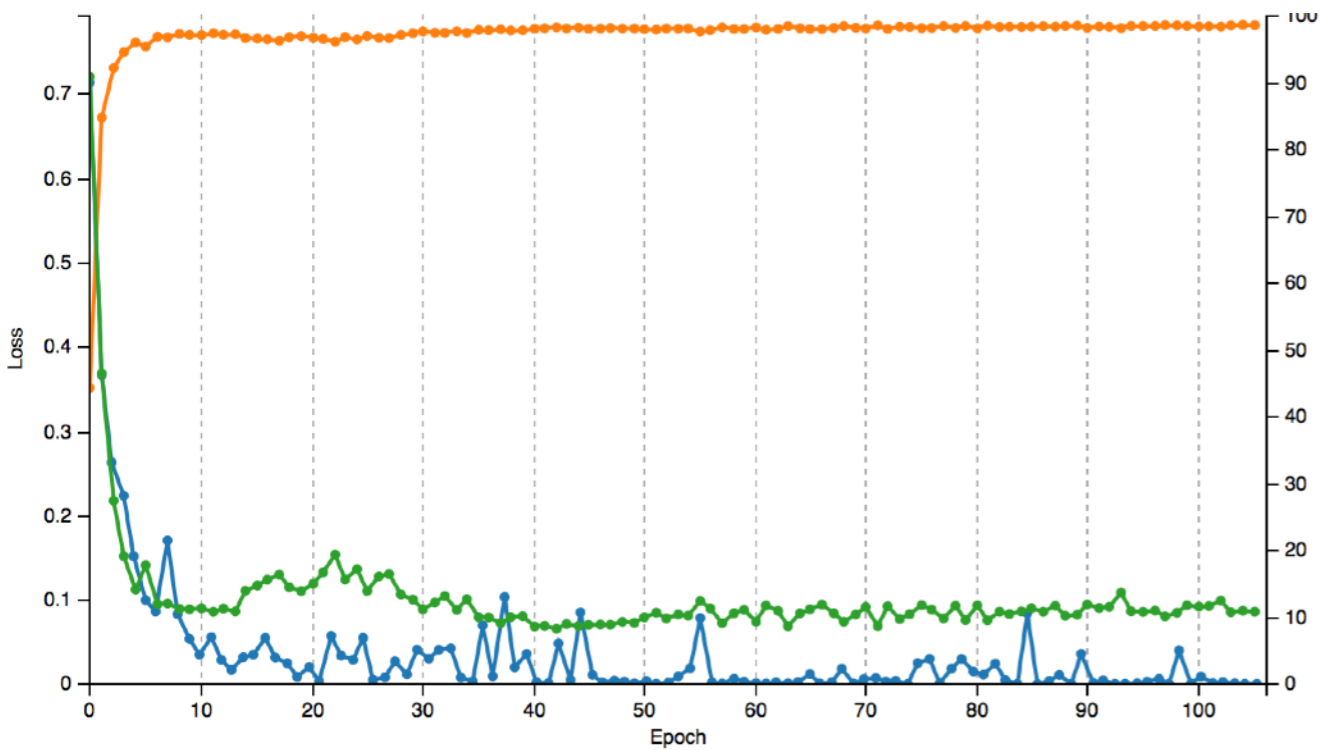
Traffic Light?



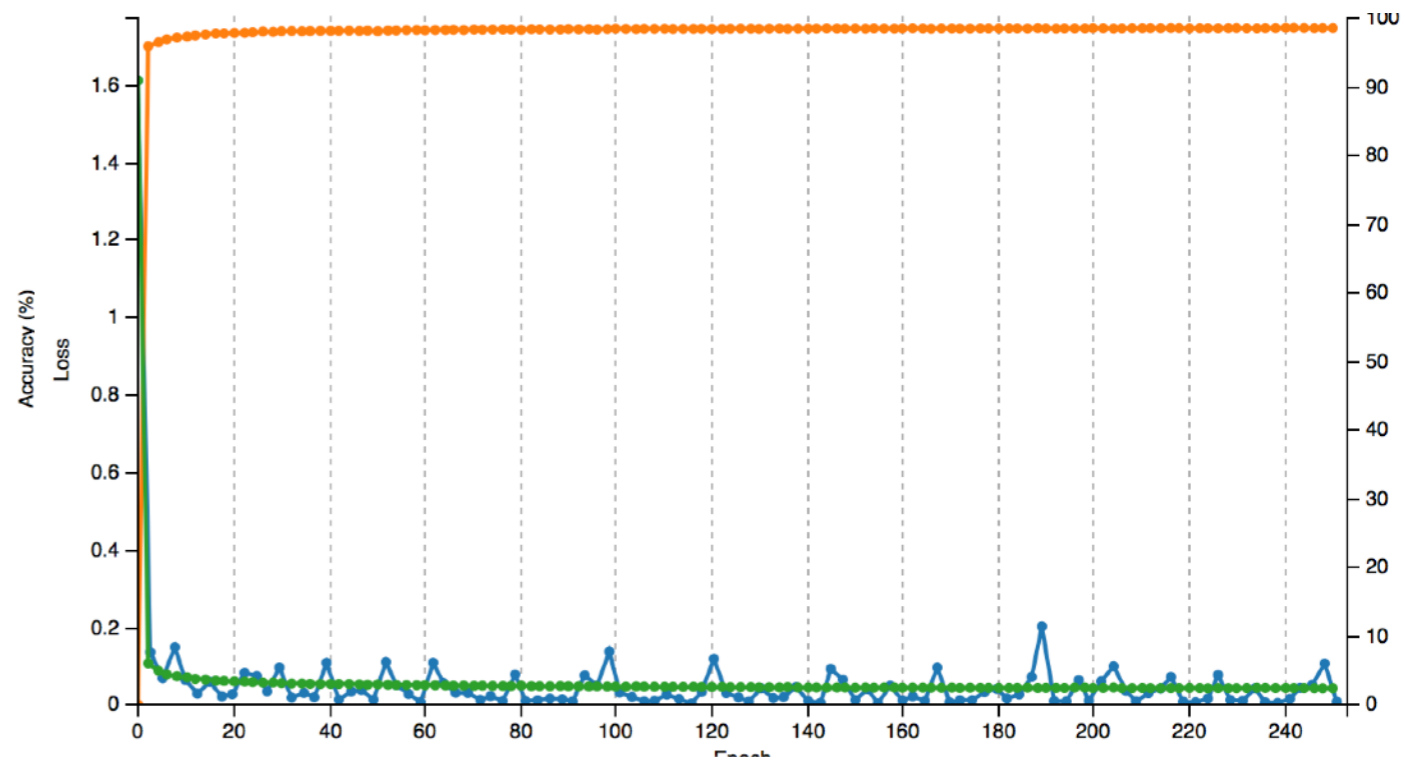
Essentially, we found a hack



Loss functions

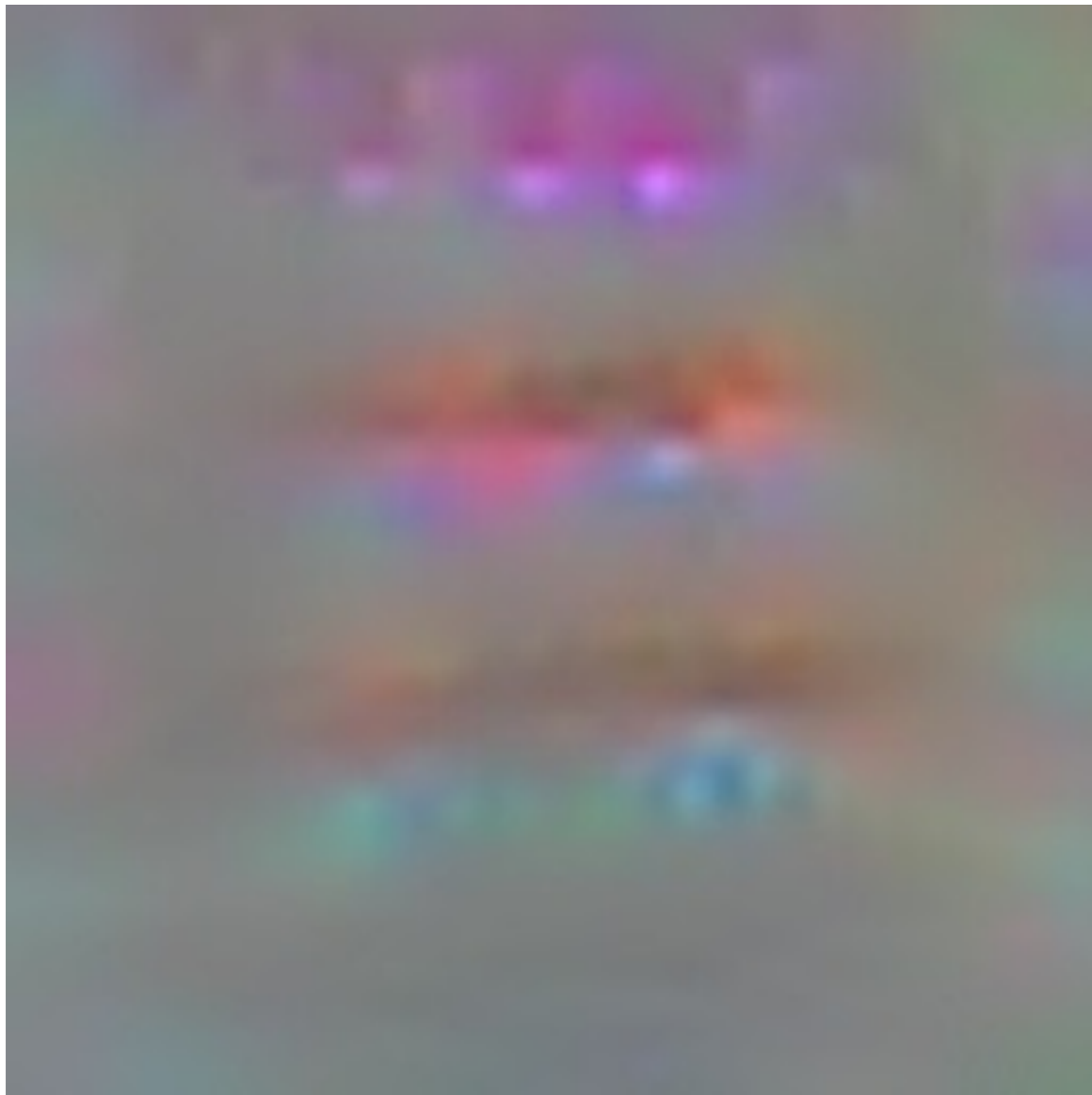


from scratch

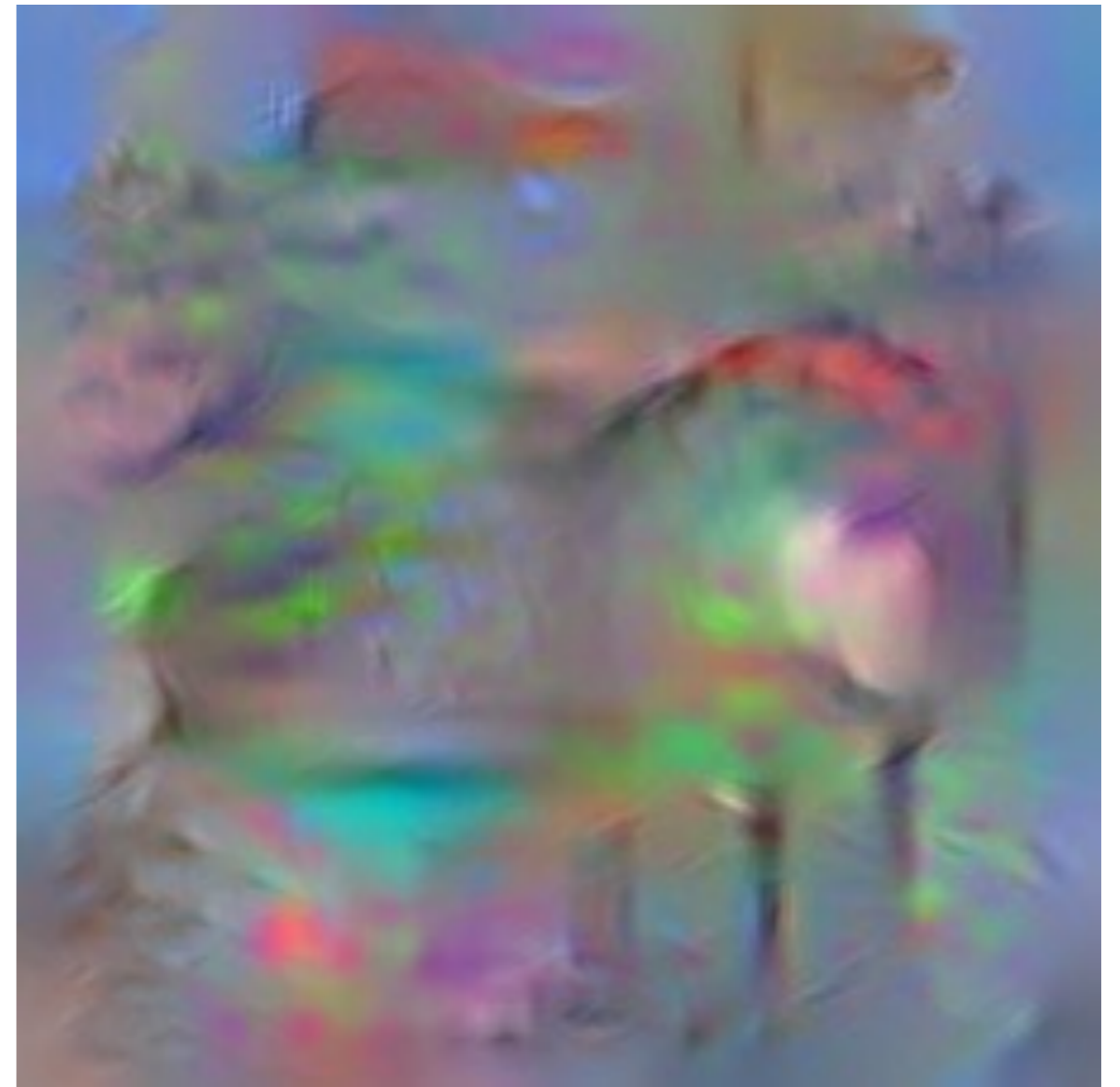


pre-trained

'Prob' layer (Traffic Light)



from scratch

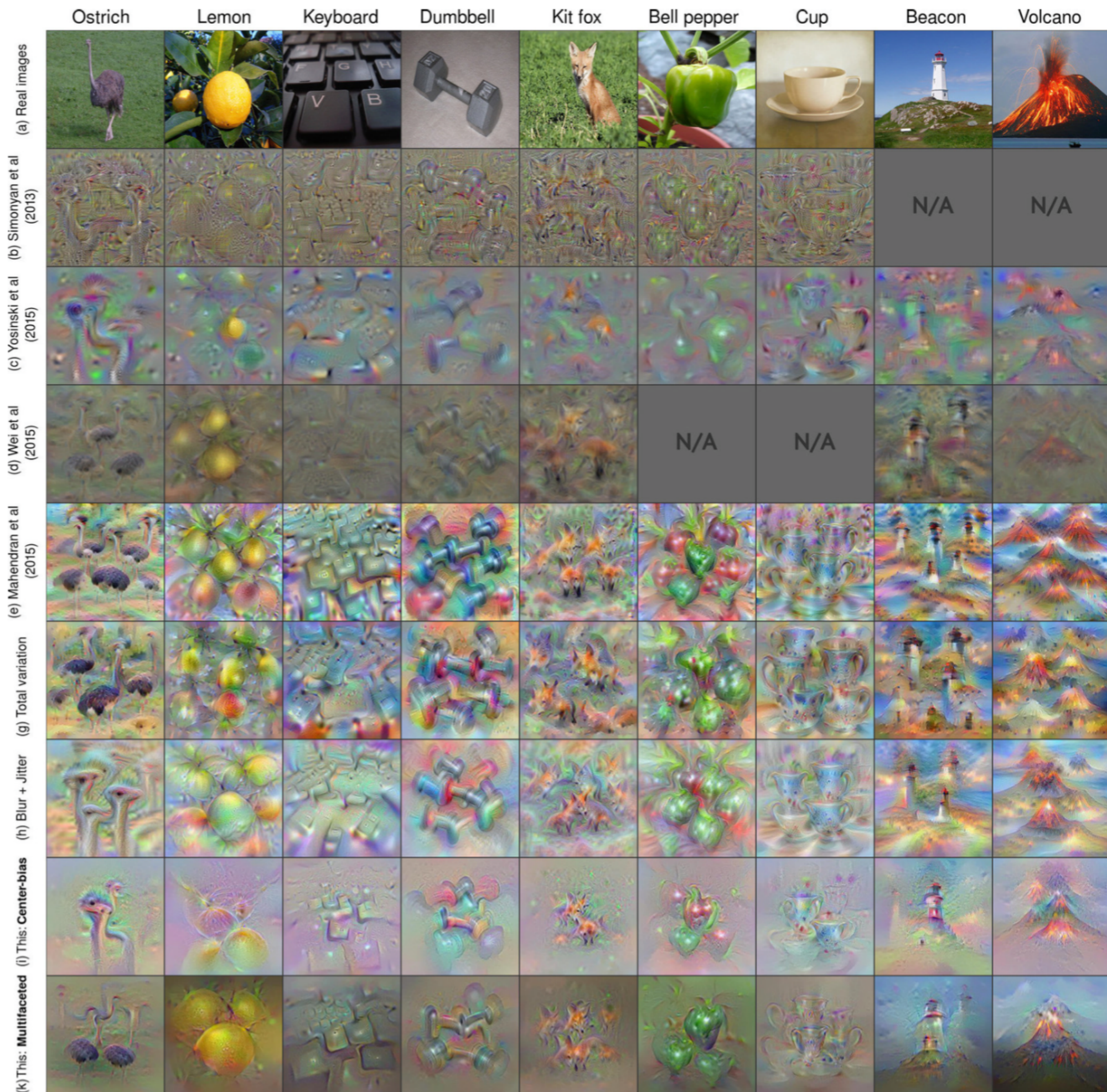


pre-trained

In Practice



Methods



How transferable are features in deep neural networks?

Jason Yosinski,¹ Jeff Clune,² Yoshua Bengio,³ and Hod Lipson⁴

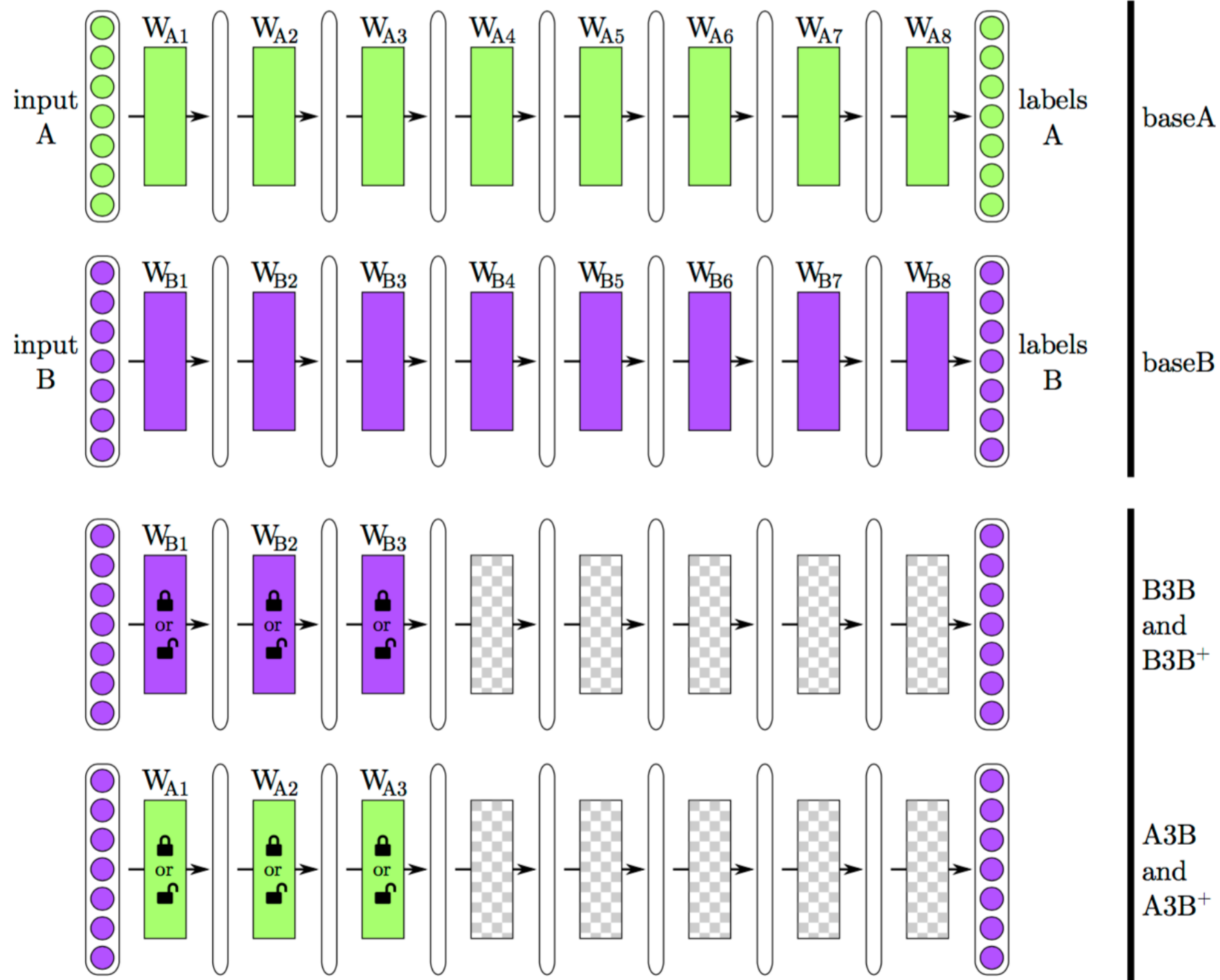
¹ Dept. Computer Science, Cornell University

² Dept. Computer Science, University of Wyoming

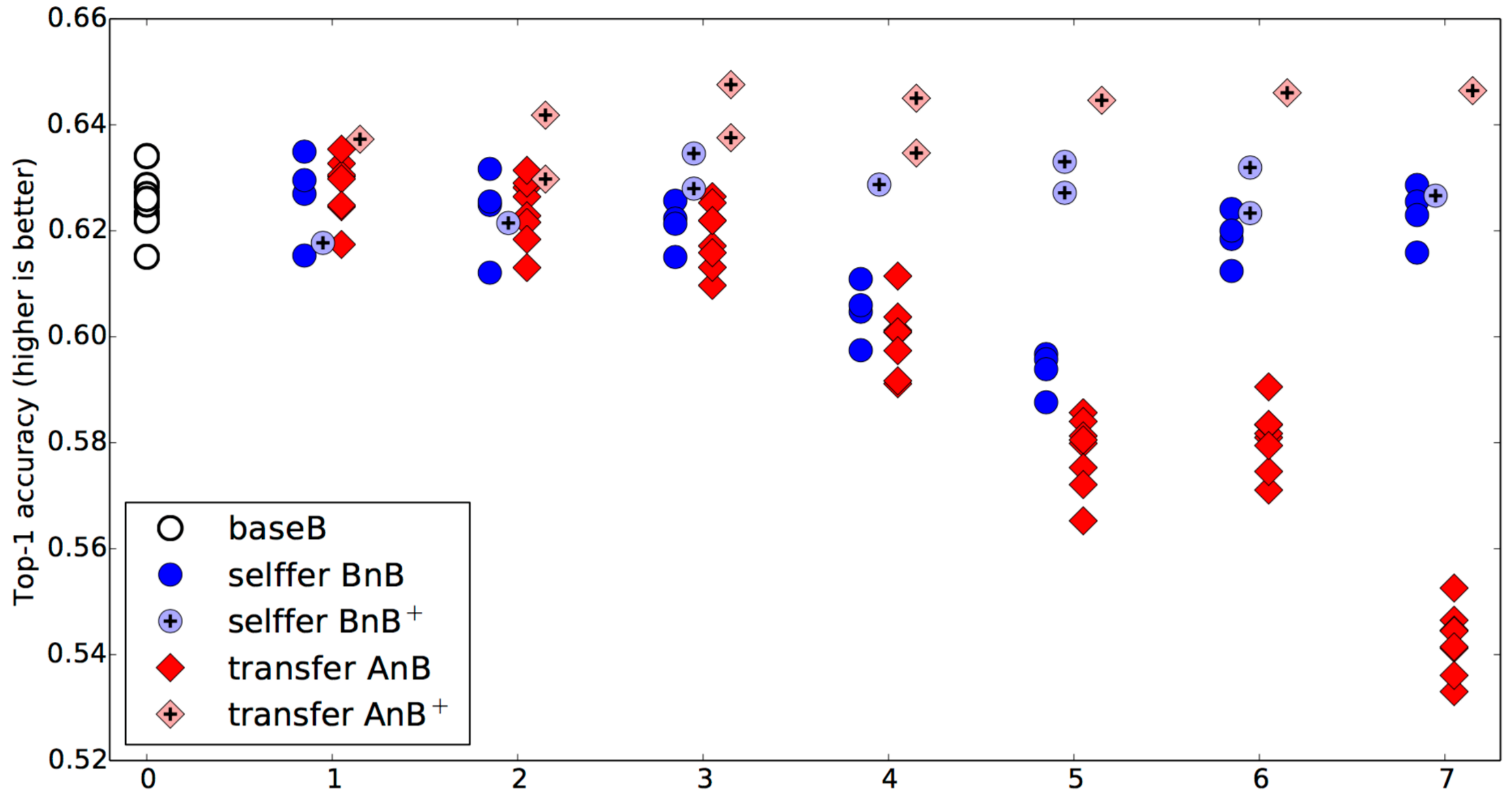
³ Dept. Computer Science & Operations Research, University of Montreal

⁴ Dept. Mechanical & Aerospace Engineering, Cornell University

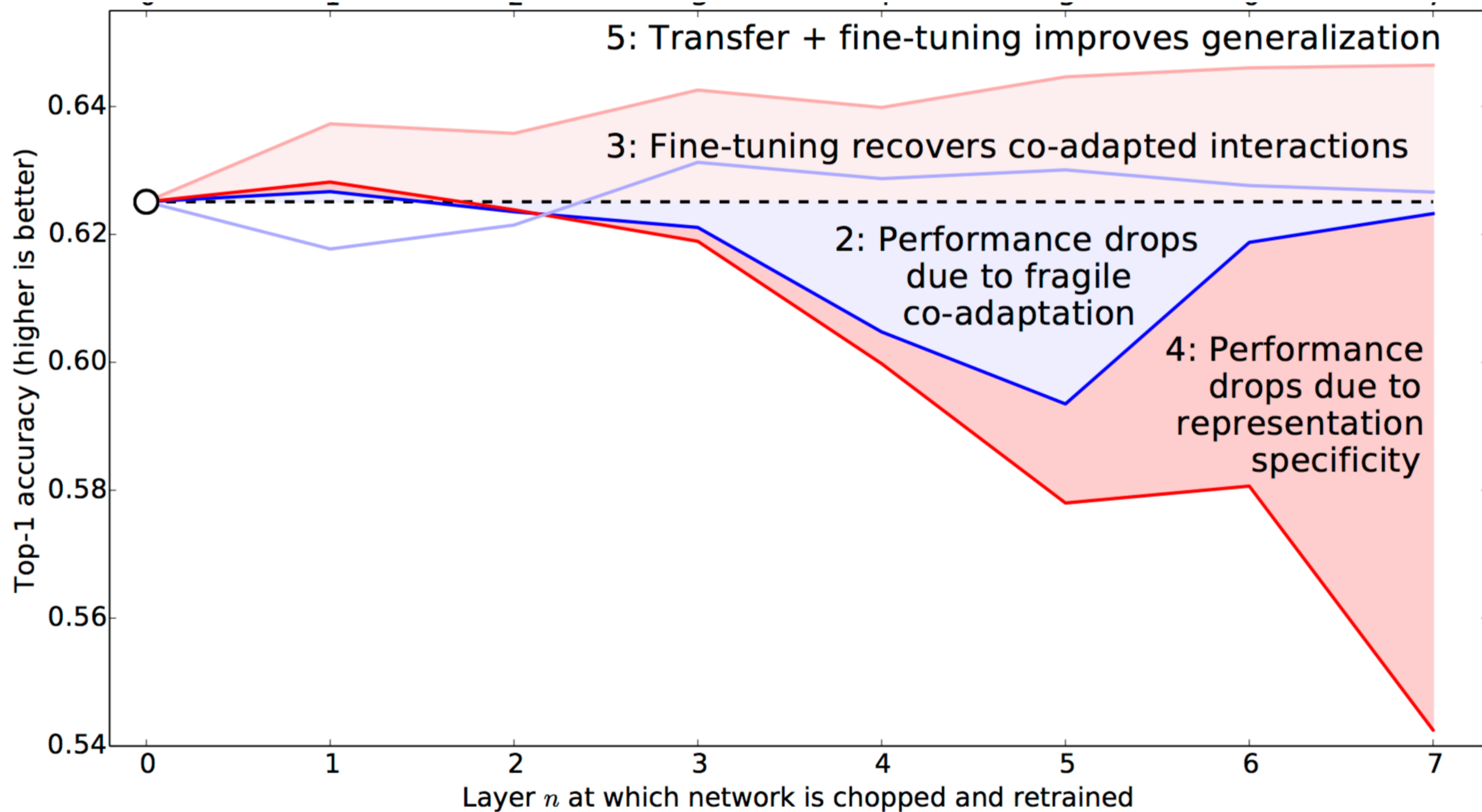
Experiment setup



Where to chop?



Interpretation



Conclusion

- 1) features are difficult to learn from small datasets
- 2) features are transferable
- 3) pre-training is a valid solution to lack of data if task is similar
- 4) layer visualization is a good way to assess DNN quality

Fin

