

Deep Learning

Episode 6

How you *actually do it* :)



Large-scale recognition



CIFAR



- 60k images
 - 10 classes (left)
 - 32x32 RGB
 - subclasses
- | | |
|-----------------|---|
| aquatic mammals | beaver, dolphin, otter, seal, whale |
| fish | aquarium fish, flatfish, ray, shark, trout |
| flowers | orchids, poppies, roses, sunflowers, tulips |
| food containers | bottles, bowls, cans, cups, plates |

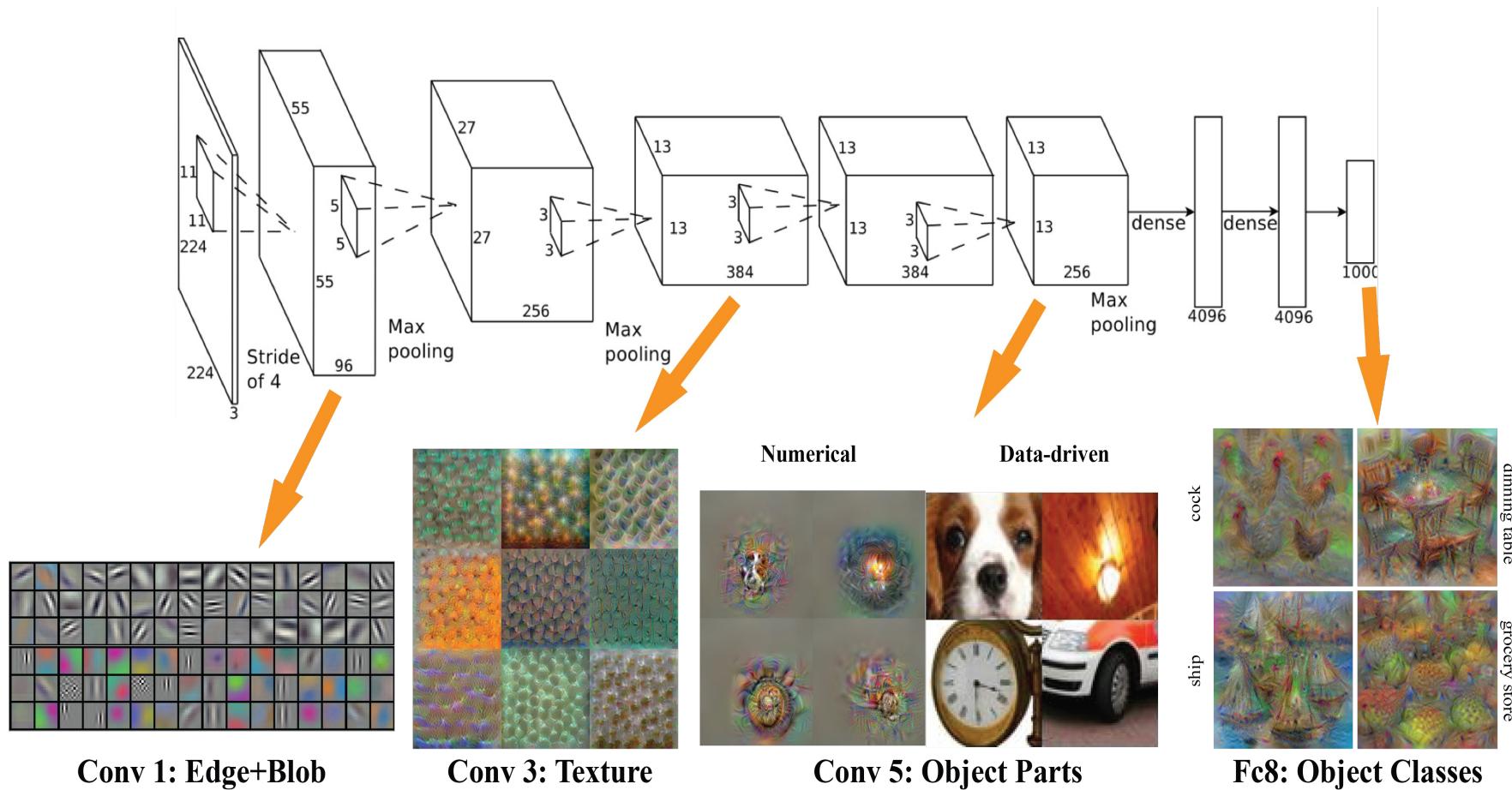
The only problem: you only have 5k images
e.g. transport vs military aircrafts

Ideas?

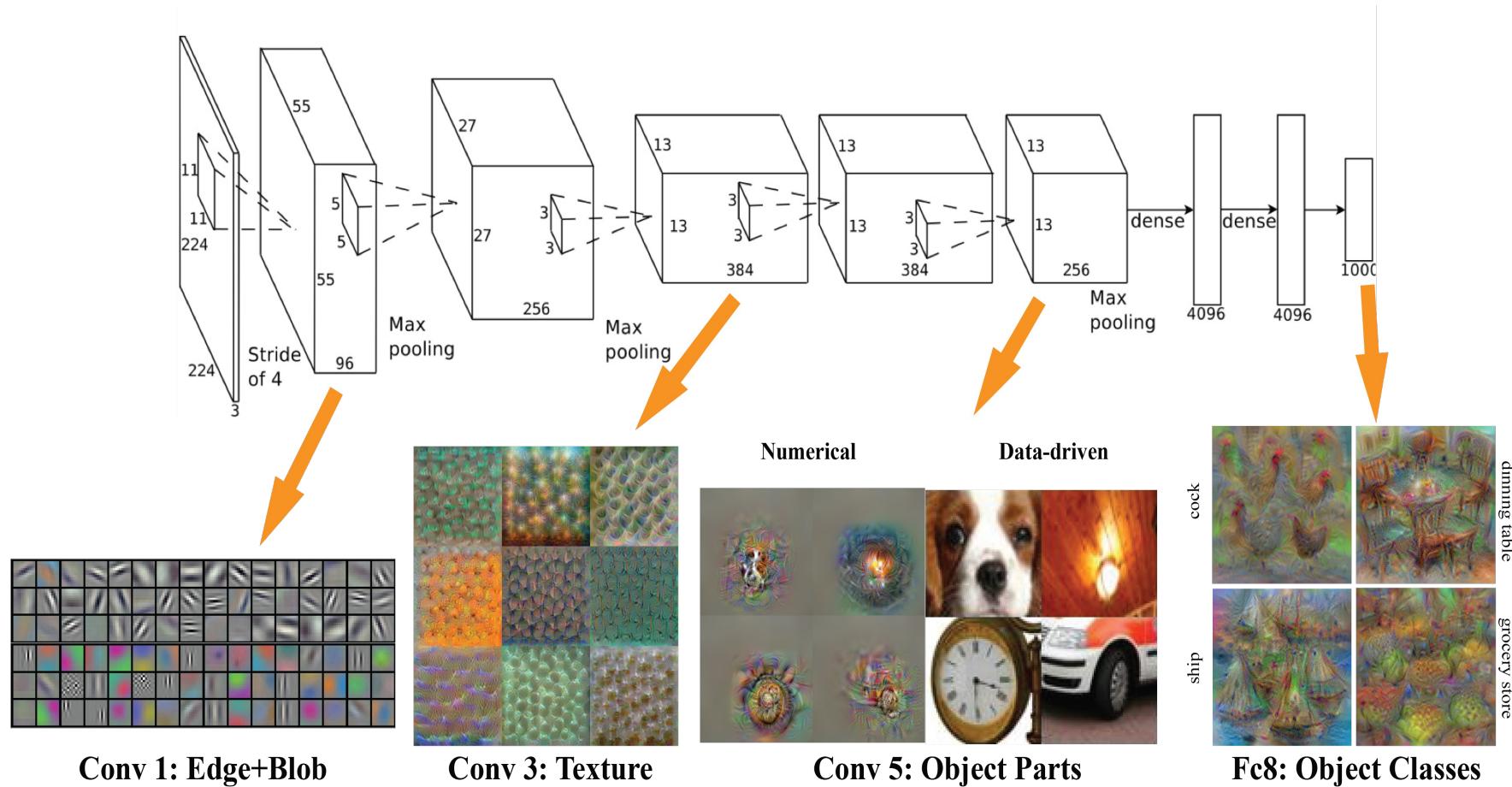
Regular solutions

- Regularize really hard
- Super small network
- Data augmentation
- Whatever

Feature learning



Feature learning



Idea: let's pre-train network on a larger dataset

Pre-training

- 1. Train a network on large dataset

cifar X

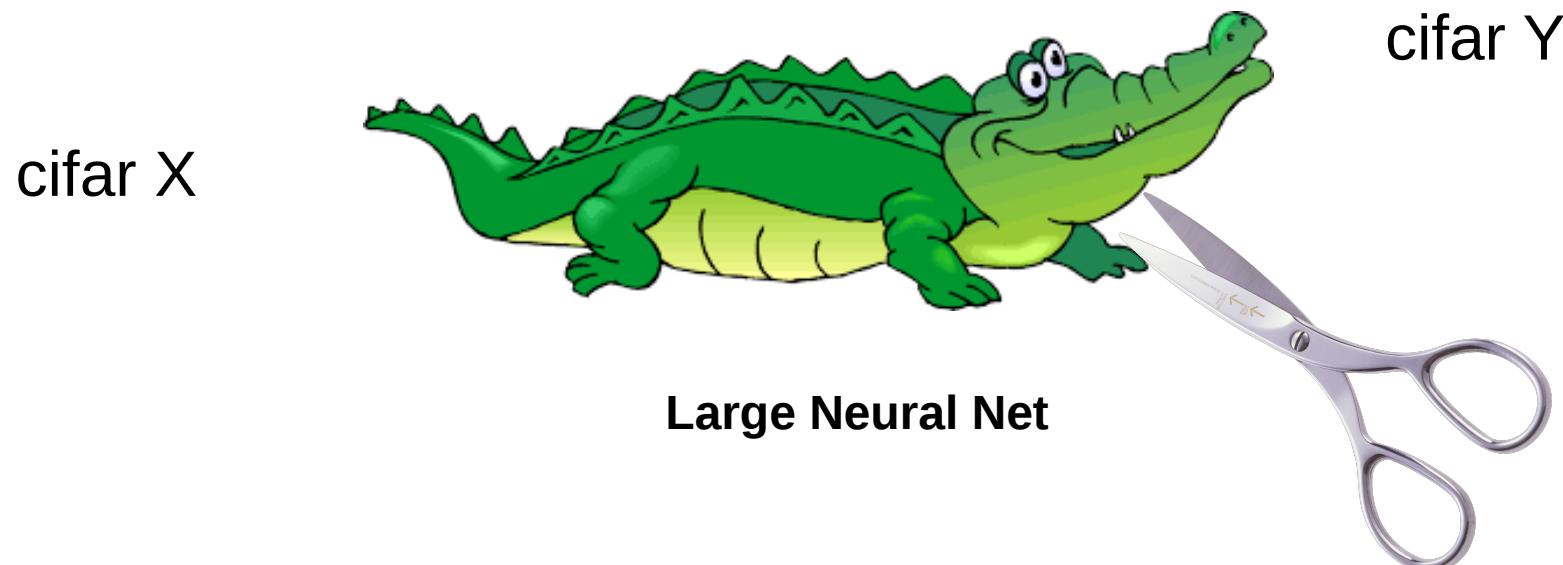


cifar Y

Large Neural Net

Pre-training

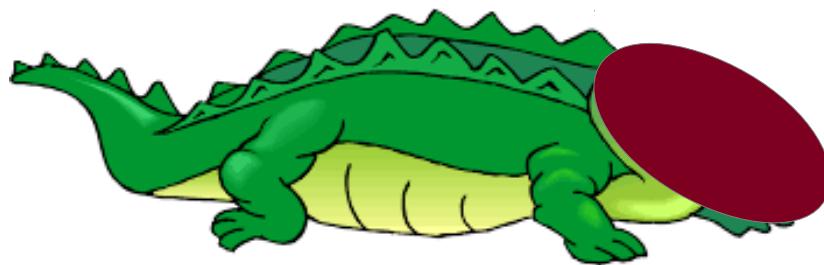
- 1. Train a network on large dataset
- 2. Take some intermediate layer



Pre-training

- 1. Train a network on large dataset
- 2. Take some intermediate layer

cifar X

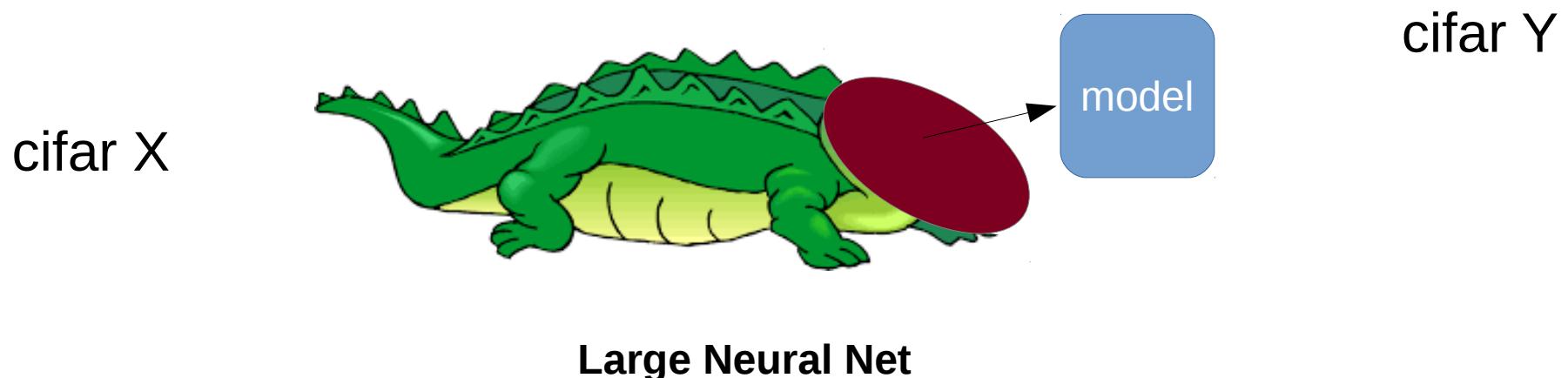


cifar Y

Large Neural Net

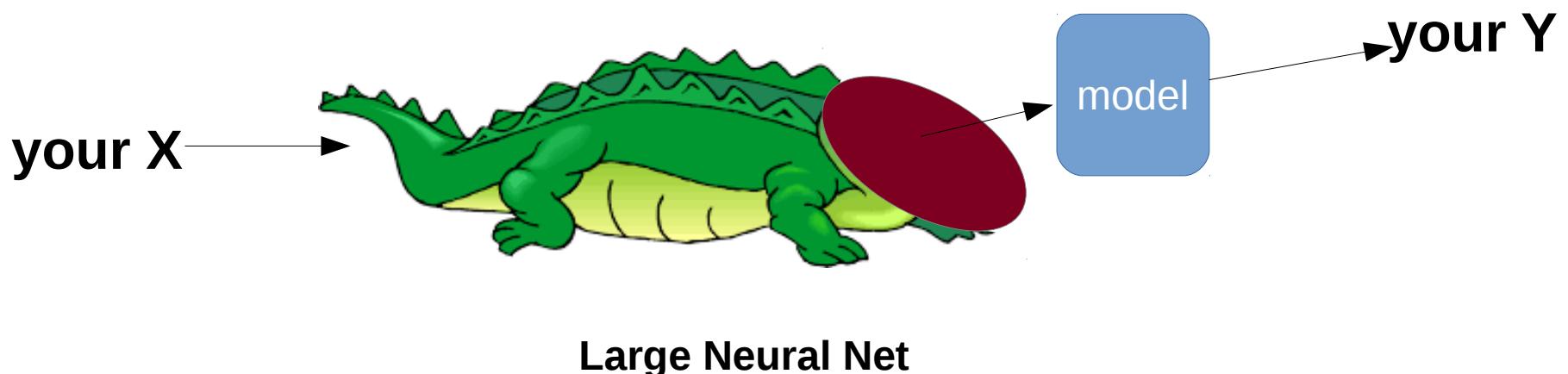
Pre-training

- 1. Train a network on large dataset
- 2. Take some intermediate layer
- 3. Build model on top of it



Pre-training

- 1. Train a network on large dataset
- 2. Take some intermediate layer
- 3. Build model on top of it
- 4. Train model for your objective



CIFAR



- 60k images
- 10 classes (left)
- 32x32 RGB
- subclasses

aquatic mammals

beaver, dolphin, otter, seal, whale

fish

aquarium fish, flatfish, ray, shark, trout

flowers

orchids, poppies, roses, sunflowers, tulips

food containers

bottles, bowls, cans, cups, plates

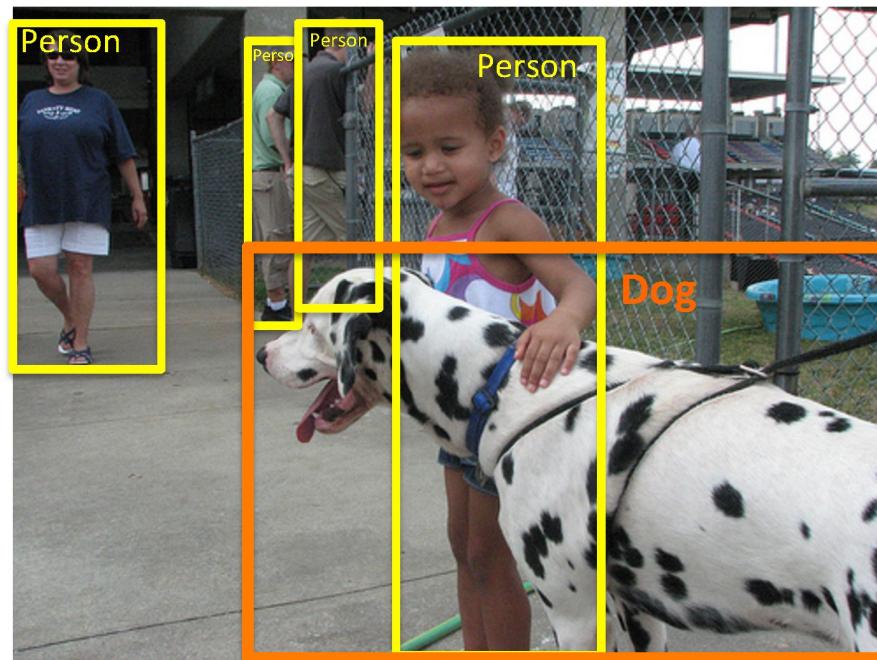
Pfft... weakling

IMAGENET Large Scale Visual Recognition Challenge (ILSVRC) 2010-2014

200 object classes
1000 object classes

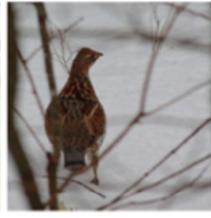
456,567 images
1,431,167 images

DET
CLS-LOC

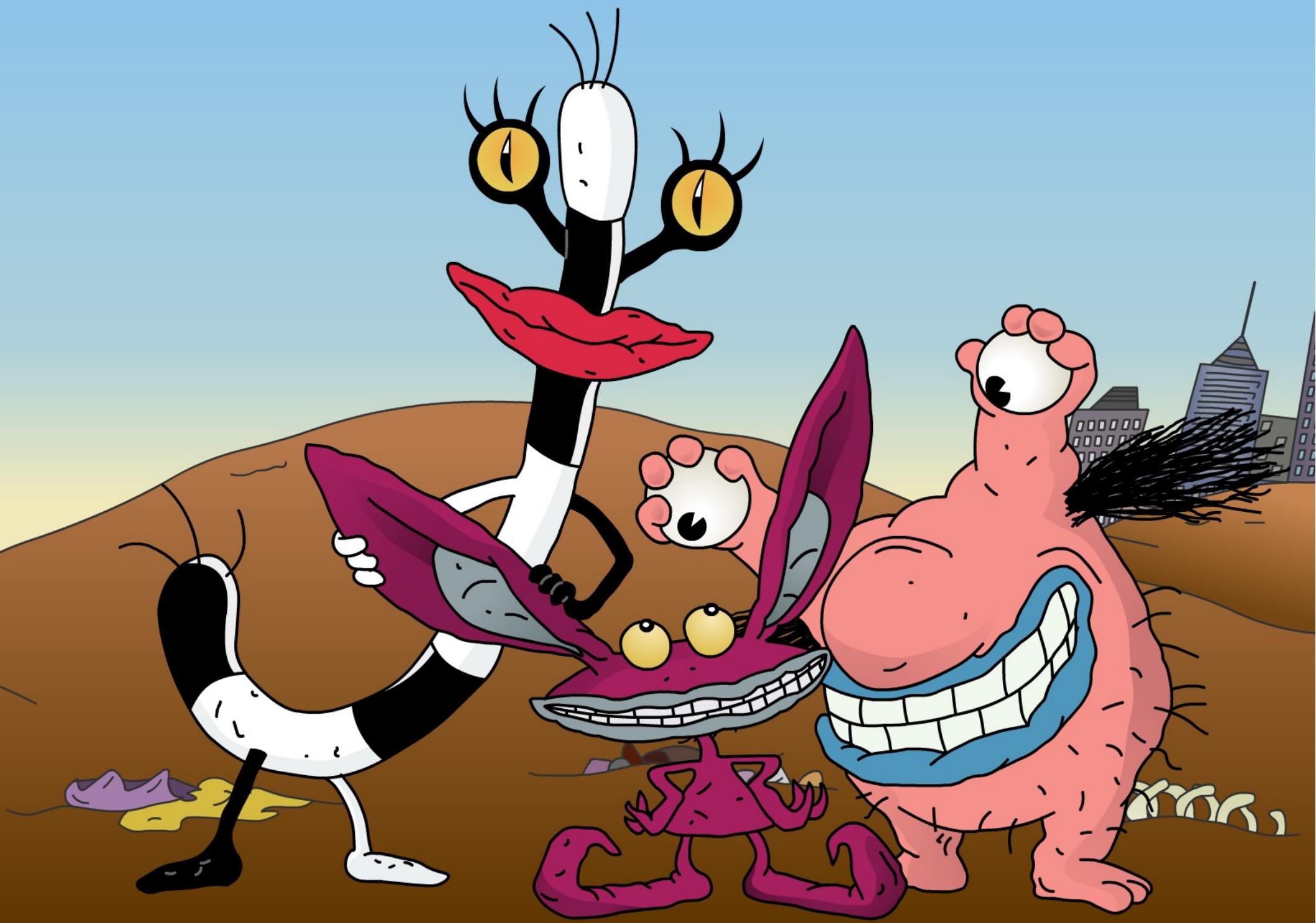


<http://image-net.org/challenges/LSVRC/>

Variety of object classes in ILSVRC

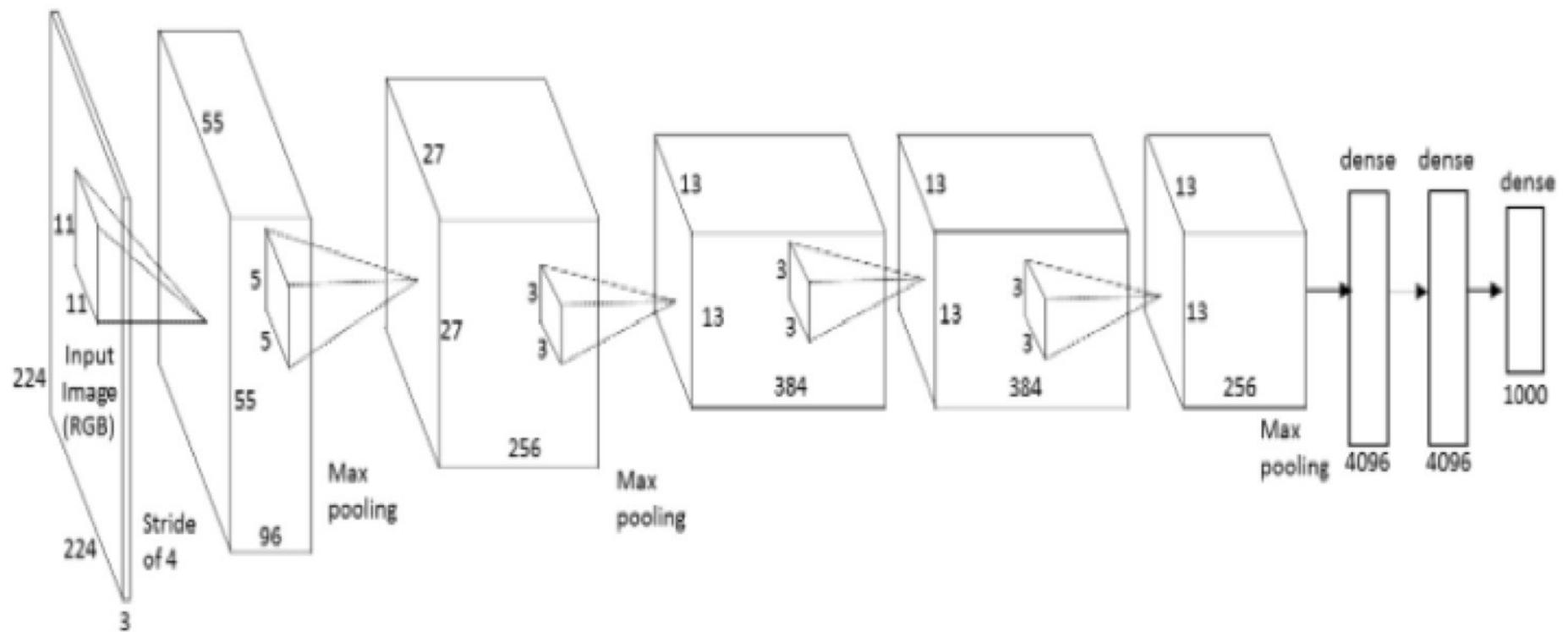
	DET	CLS-LOC				
birds						
bottles						
cars						

What kind of network will handle that?



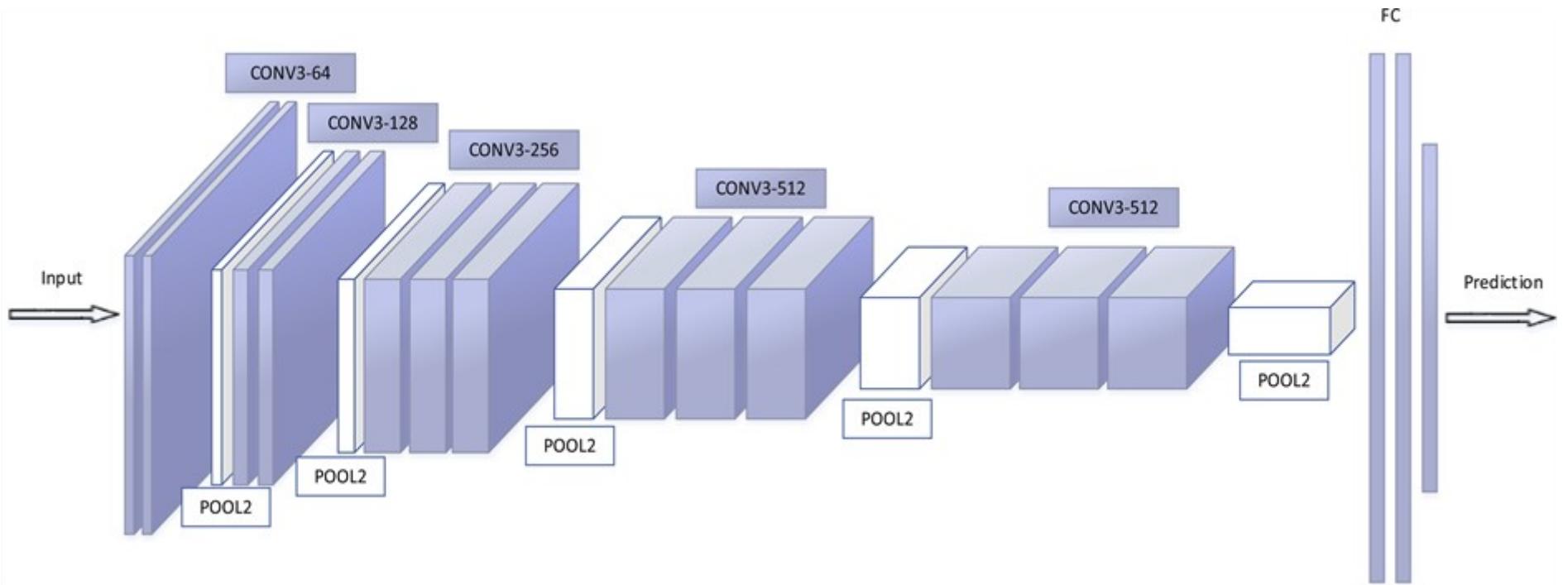
BY JAVIER 2006

Alexnet



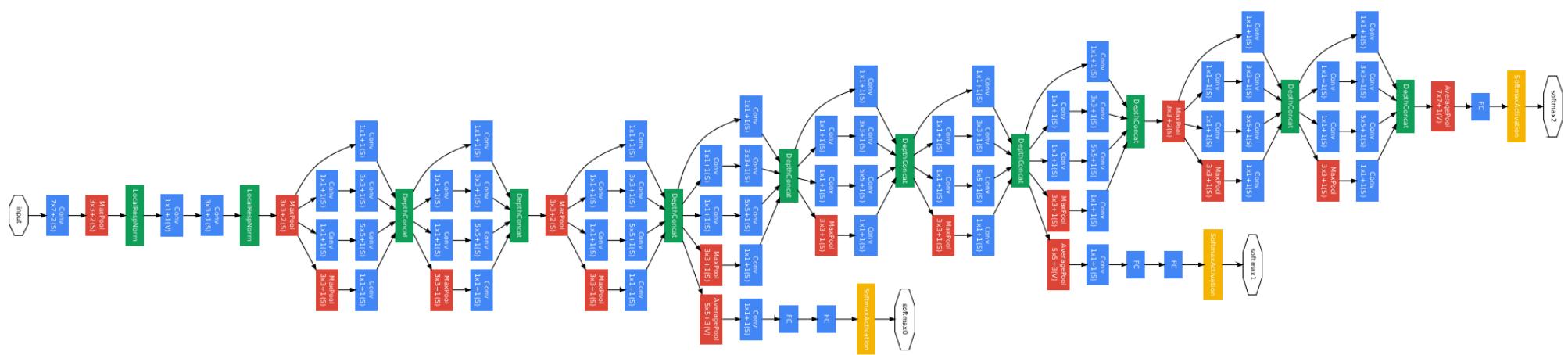
(Krizhevsky et al.)

VGG



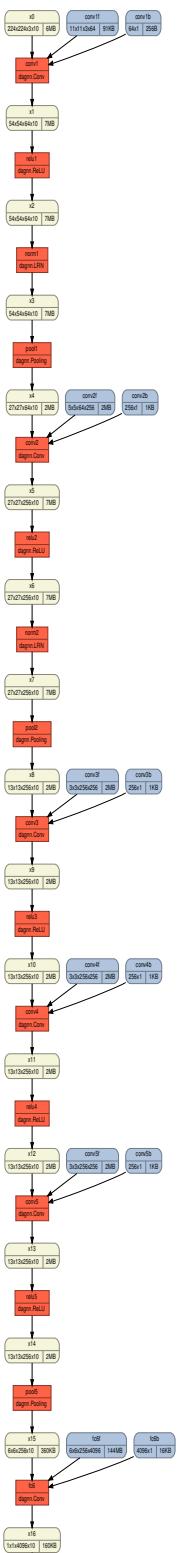
Karen Simonyan and Andrew Zisserman

Googlenet (inception)

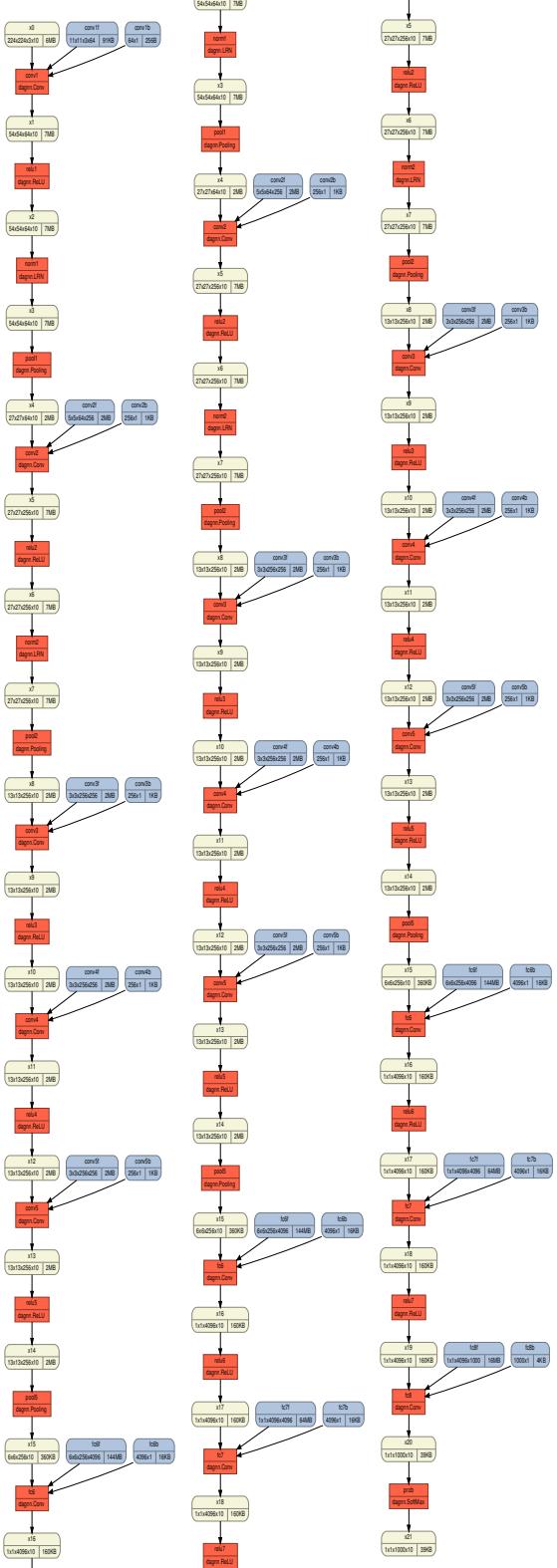


Szegedy et al

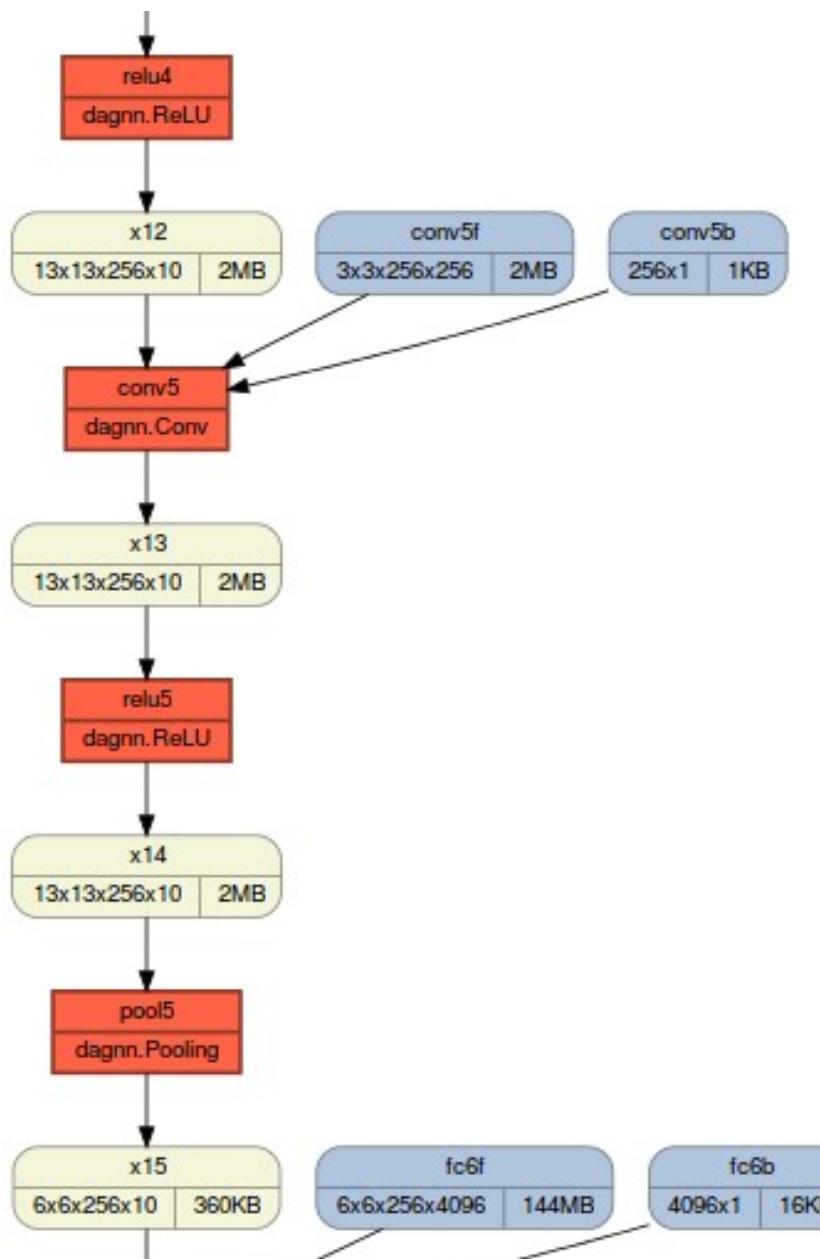
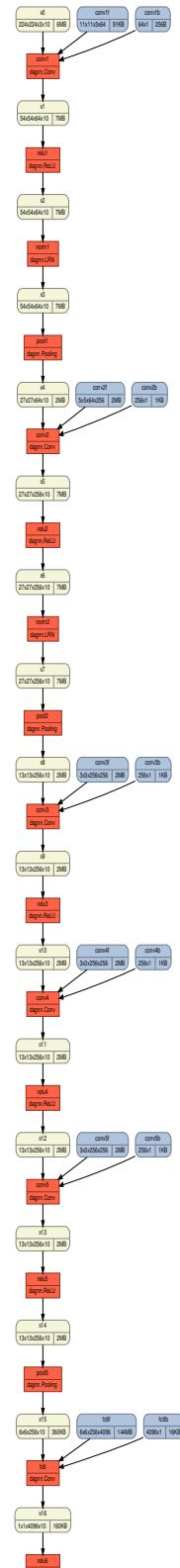
ResNet-152



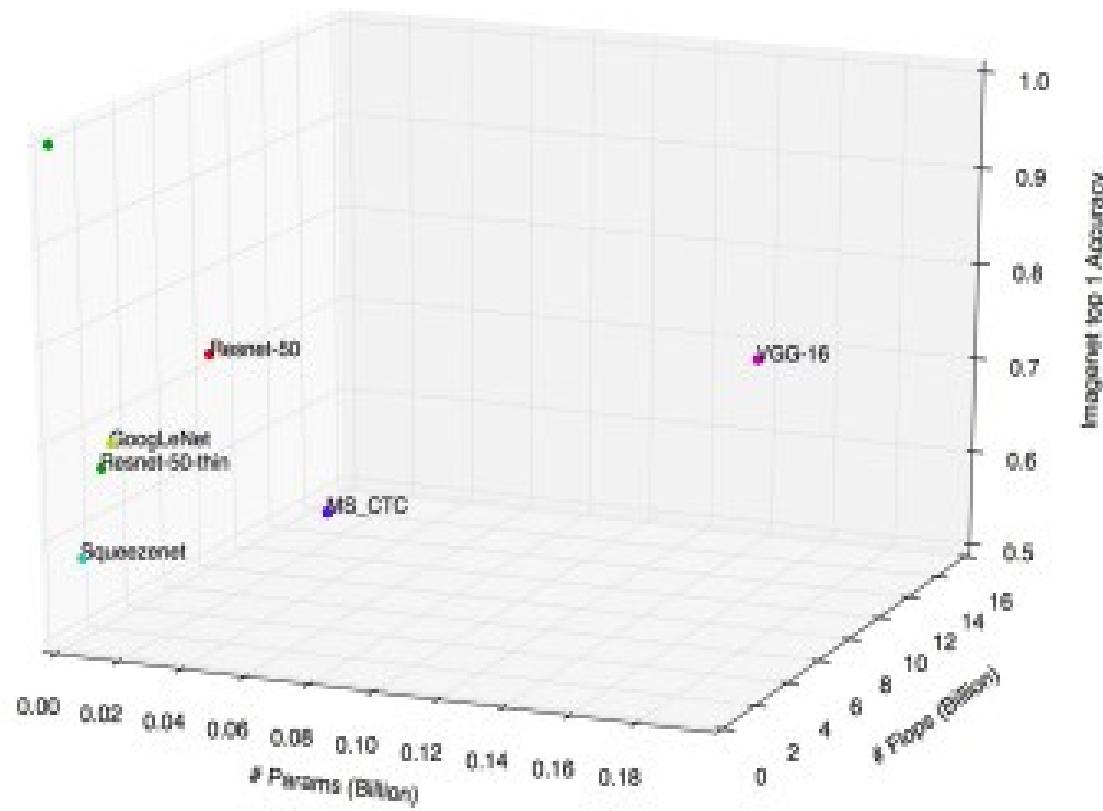
ResNet-152



ResNet-152



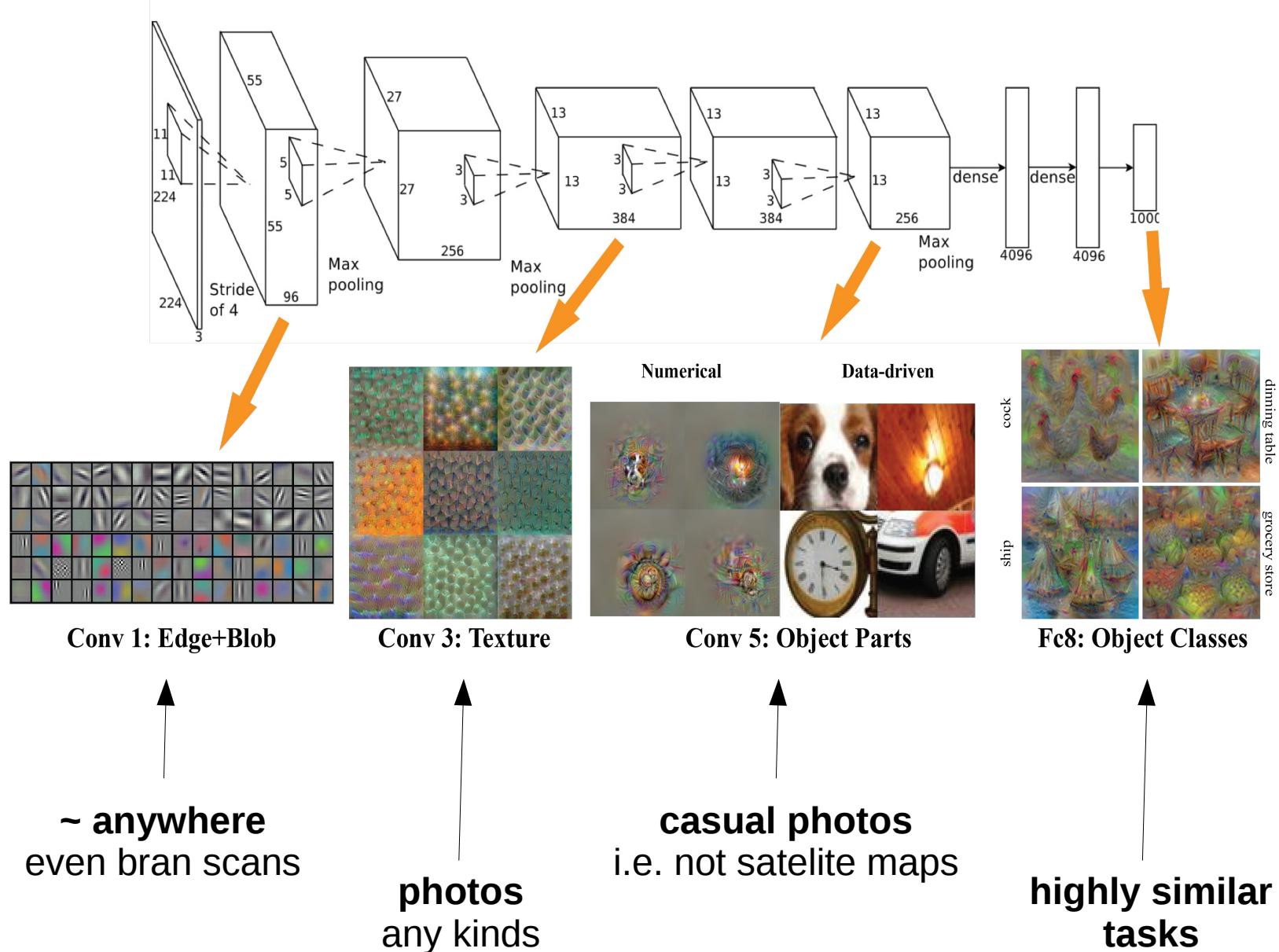
Monster stats



Can we use them?*

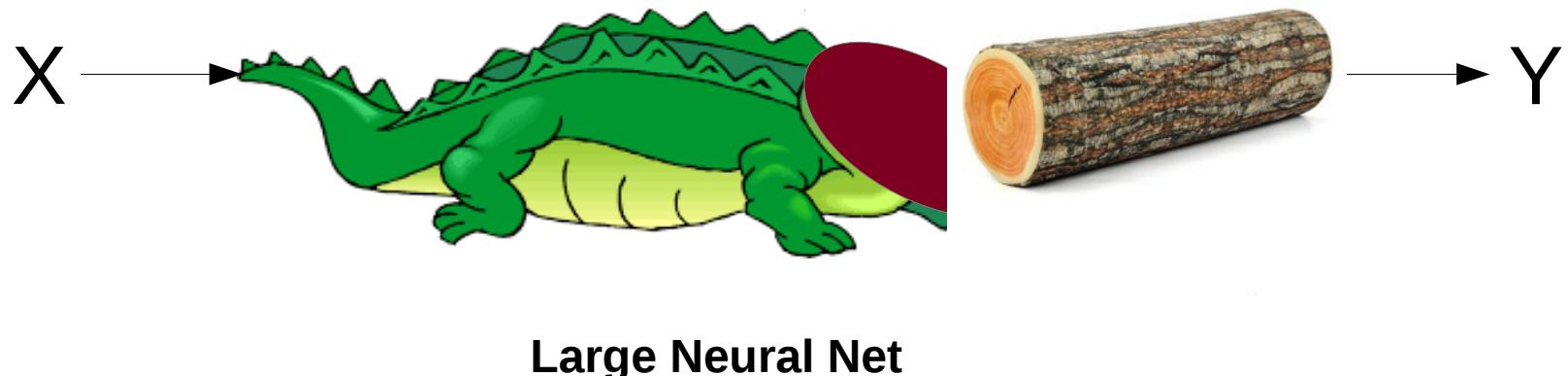
*for aircraft problem

Can we use them?



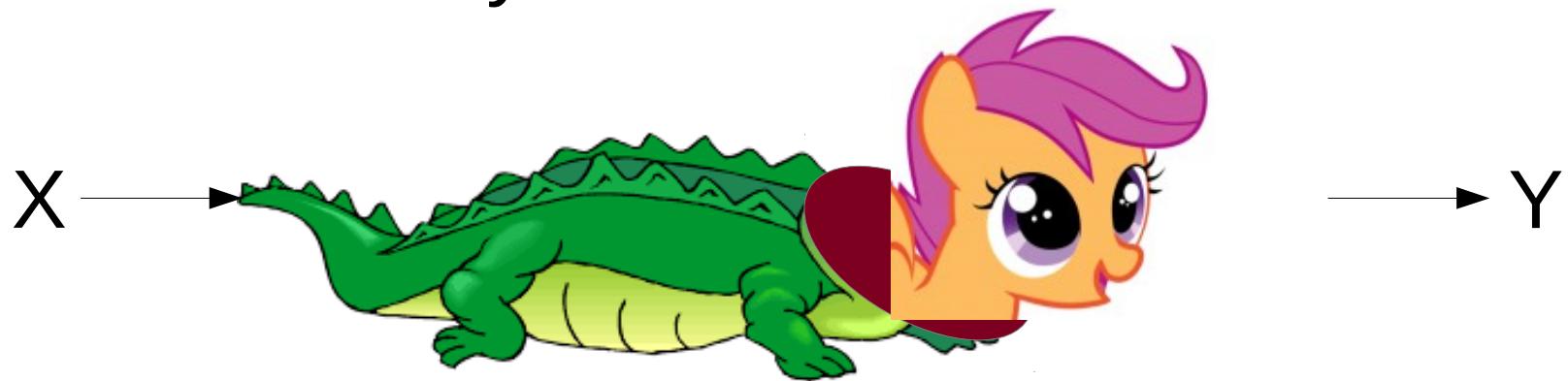
Reusing features

- Chop off “head”
- use “neck” as feature extractor
- Train ANY classifier
 - even random forest will do



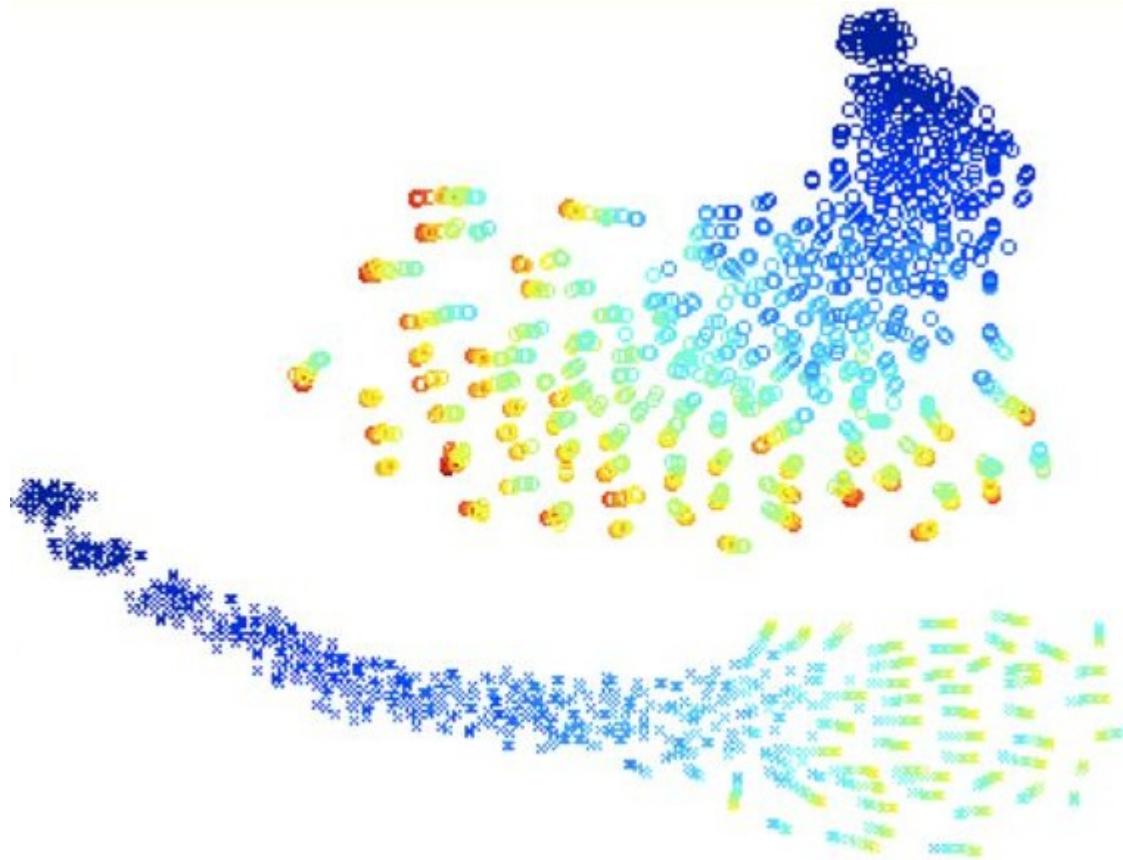
Fine-tuning

- Chop off “head”
- “freeze” body (consider constant)
- Build new neural network in it's place
- Train “head” only for several iterations
- Un-freeze body and train **full network**



Large Neural Net

Comparison



- TSNE of network weights
- Crosses = pre-training, Rounds = no pre-training
- Color = epoch (hotter means higher)
- No overlap!

Other kinds of NNs

Fine-tuning not restricted to computer vision

- Using word2vec as pre-trained embedding
- Using skip-thought as pre-trained encoder-decoder
- Pre-training simple shallow networks on similar domains

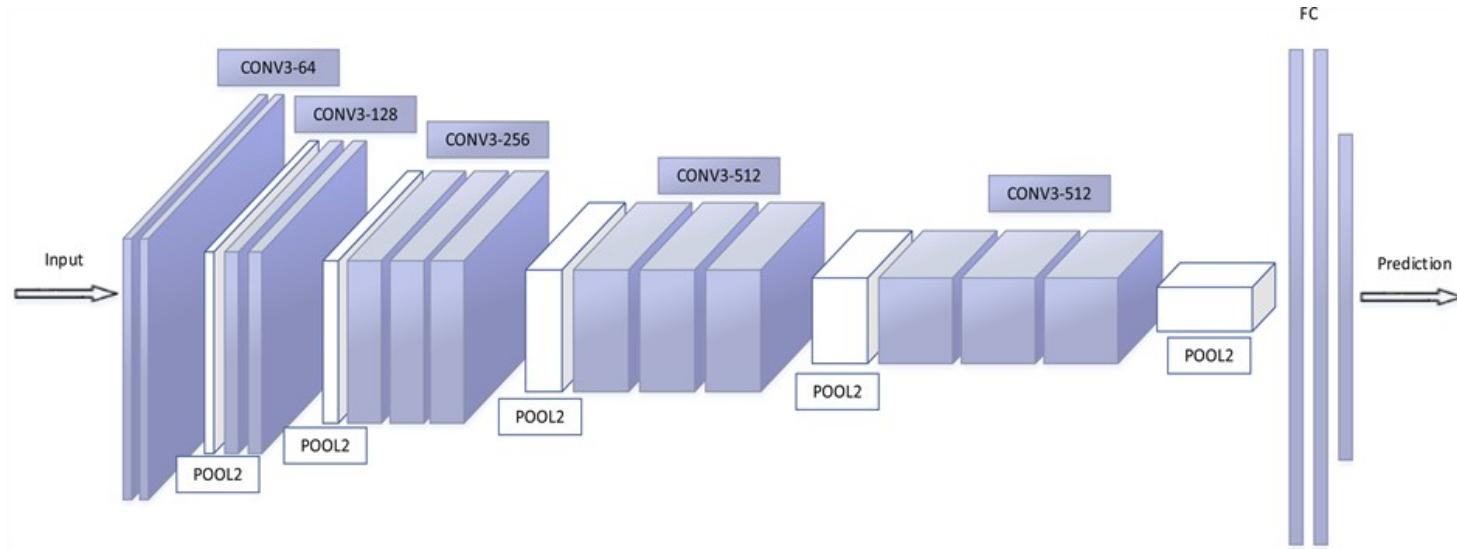
Alright, we did it!

But how do we fit large NN into an Iphone?

- That's simple!
 - We just train a smaller network
 - Or use a remote server with GPUs
 - Maybe use float16 instead of float32

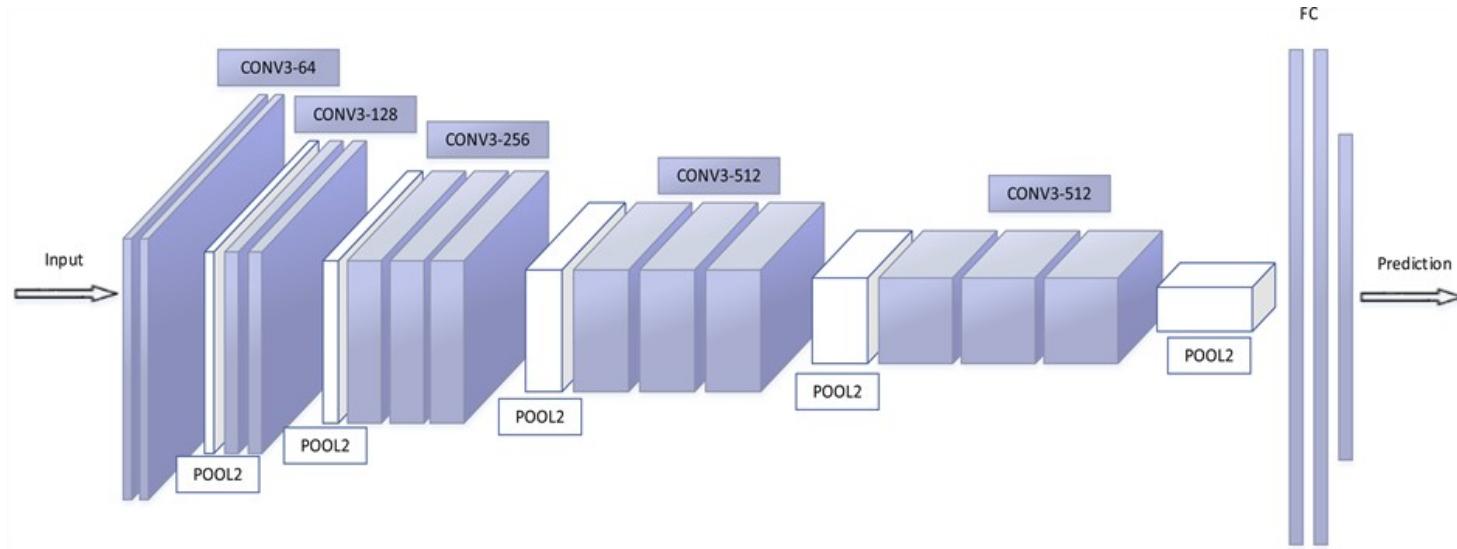
Can we do better?

Decomposition trick



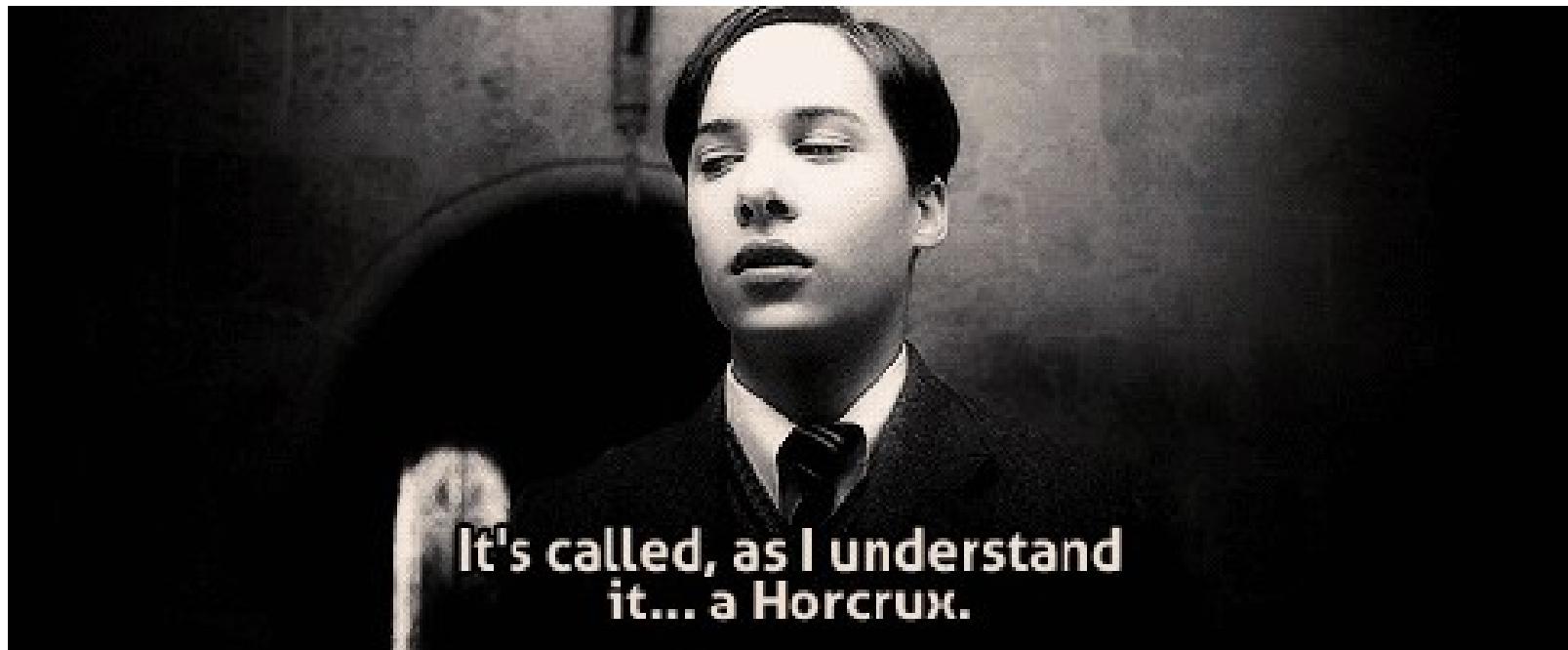
- Pre-final dense layer 4096x4096 matrix
- Filters 512x256x3x3
- Can we reduce the size/complexity?

Decomposition trick



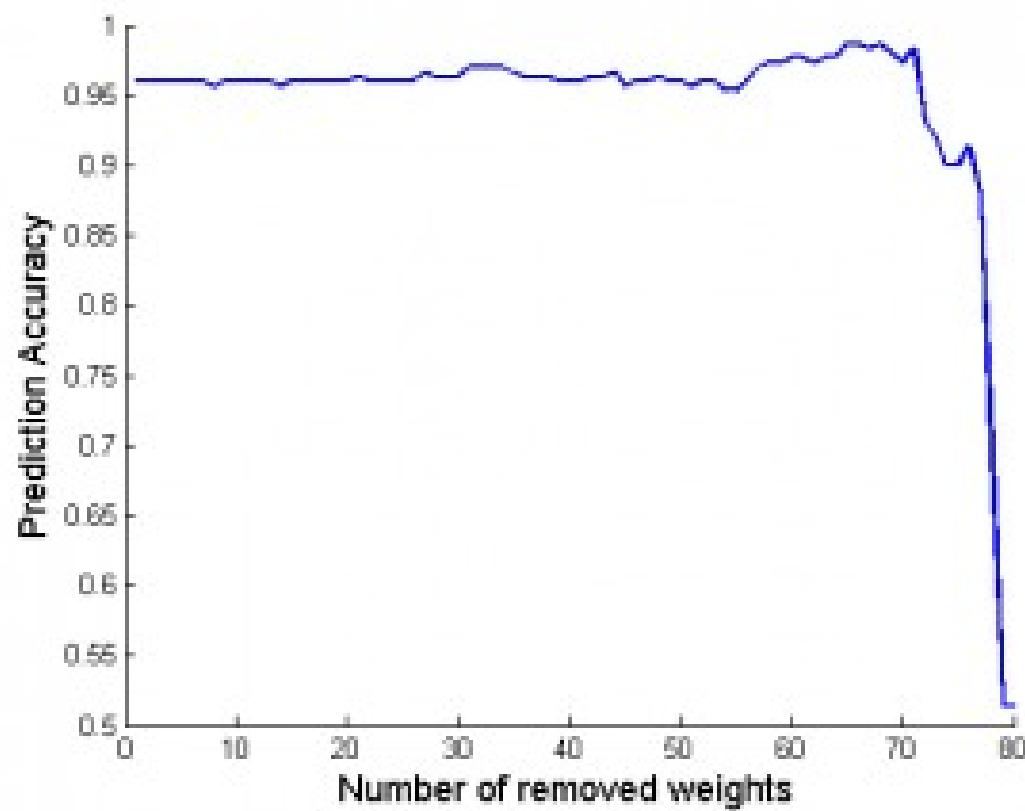
- Pre-final dense layer 4096x4096 matrix
- Filters 512x256x3x3
- Can we reduce the size/complexity?
 - Yes, we can! (c)
 - SVD (flatten to matrix), TensorTrain, ...

Reducing network size



**It's called, as I understand
it... a Horcrux.**

Optimal Brain Damage



Optimal Brain Damage

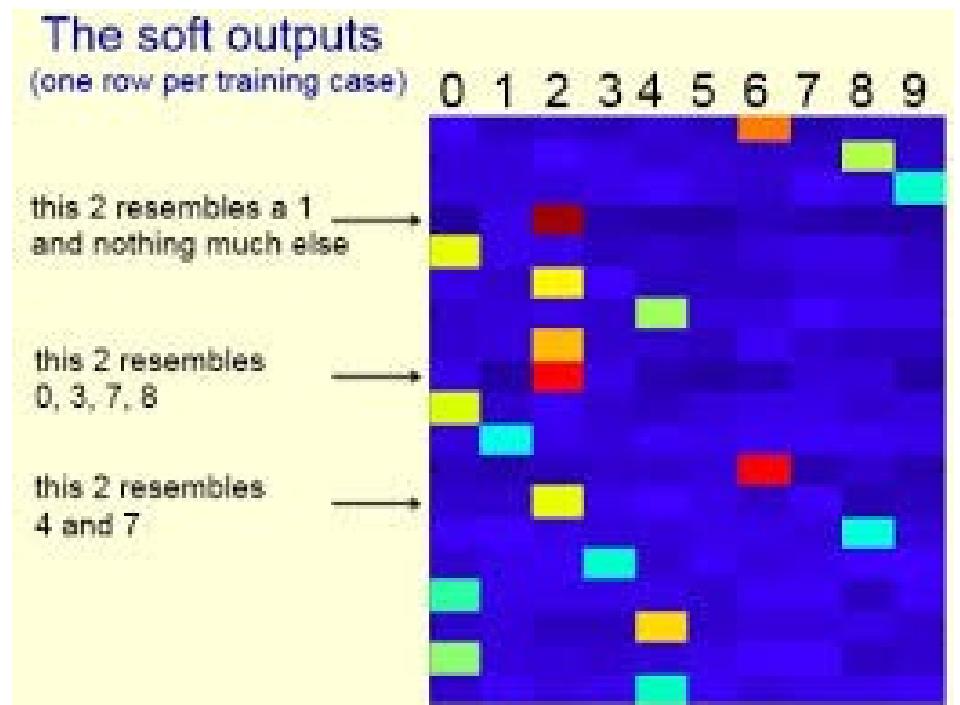
- Define “neuron saliency” ~ how bad it is to remove it

$$s_k = \frac{H_{kk} \cdot a_k^2}{2} = \frac{\delta^2 L}{\delta a_k} \cdot a_k^2 / 2$$

- Remove neurons with lowest saliency
- ???
- PROFIT!!!
- Similar idea: Regularize activations with L1 and prune zeros

Soft Targets

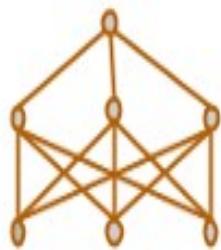
- Take the larger network
(or any model / ensemble)
- Predict probabilities for
training samples
- Use some function of proba-
bilities as objective to train
a new smaller network
- Some intuition: if an image
is truly ambiguous, this way
you will be able to handle that



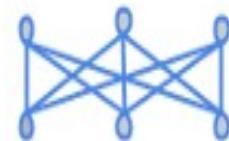
Net2net

- Same, but expanding
- Start small, extend, fine-tune

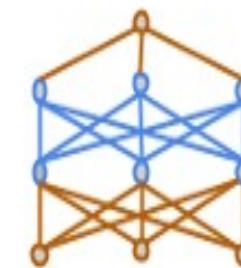
Original Model

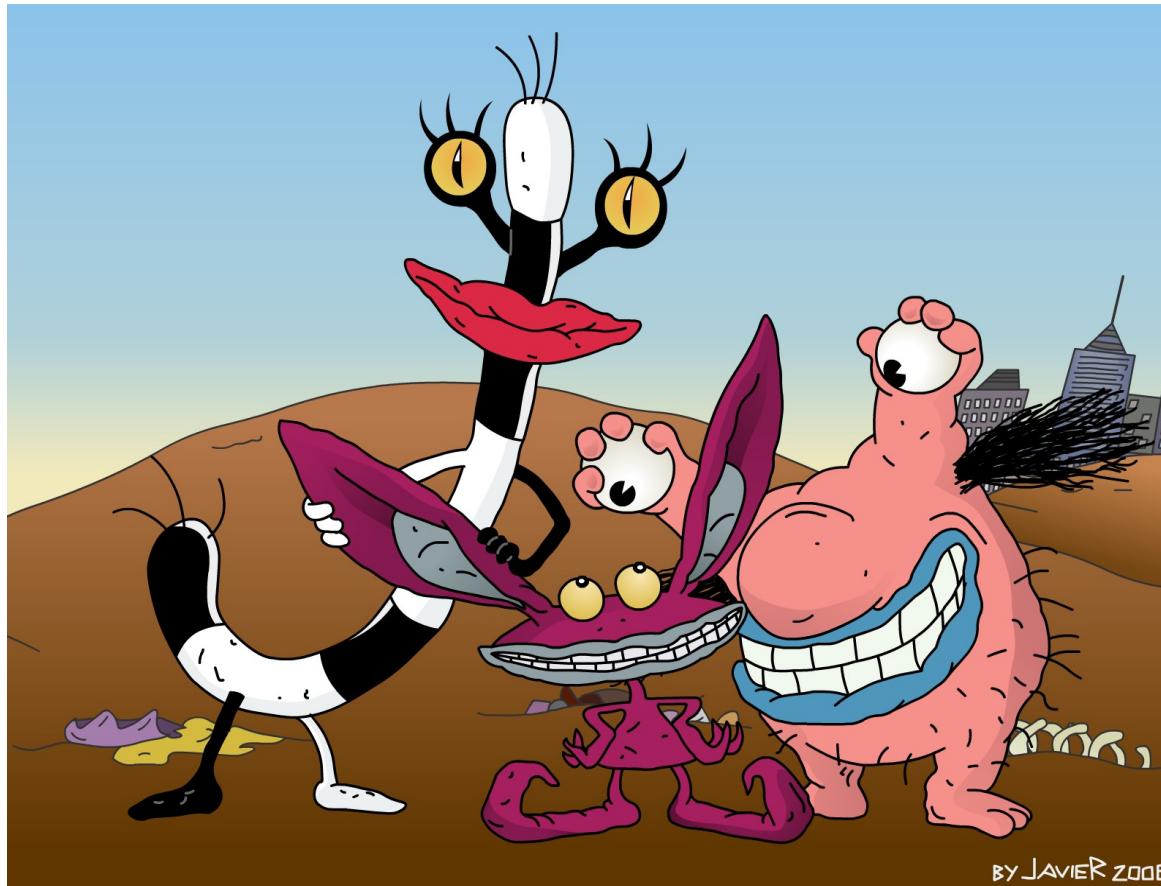


Layers that Initialized as Identity Mapping



A Deeper Model Contains Identity Mapping Initialized Layers





Come at me bro!