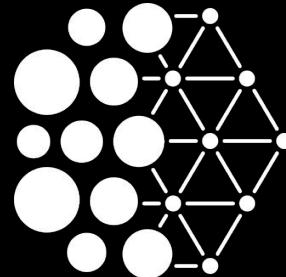


Quebec
Artificial
Intelligence
Institute



Mila

Video Processing, Network Interpretation and Updates in Computer Vision Week 3

Jeremy Pinto

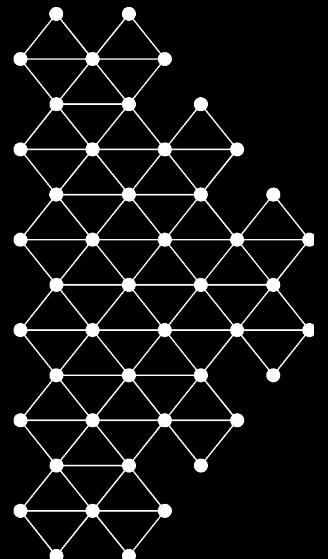
jeremy.pinto@mila.quebec

Content

- CNNs applied to Videos
- Model Interpretation
- Adversarial Attacks
- ImageNet - Updates
- Transformers in Computer Vision



Photo by [Jakob Owens](#) on [Unsplash](#)



From images to videos

What is a Video?

A video can be thought of as a **series** of **images** (*frames*) over **time** (**seconds**)*.

Typical frame rates, i.e. number of images per second, is ~30 *frames per second* (FPS).

The stream of images give our eyes the **illusion** of fluid movement.



[Source](#)

*ignore the sound for now.

Videos

How can we use deep learning on
videos?

Since videos are just **streams** of images, we could just apply image models on **each frame**, independently!

In fact, this is (basically) how cameras work - they capture each frame one at a time.

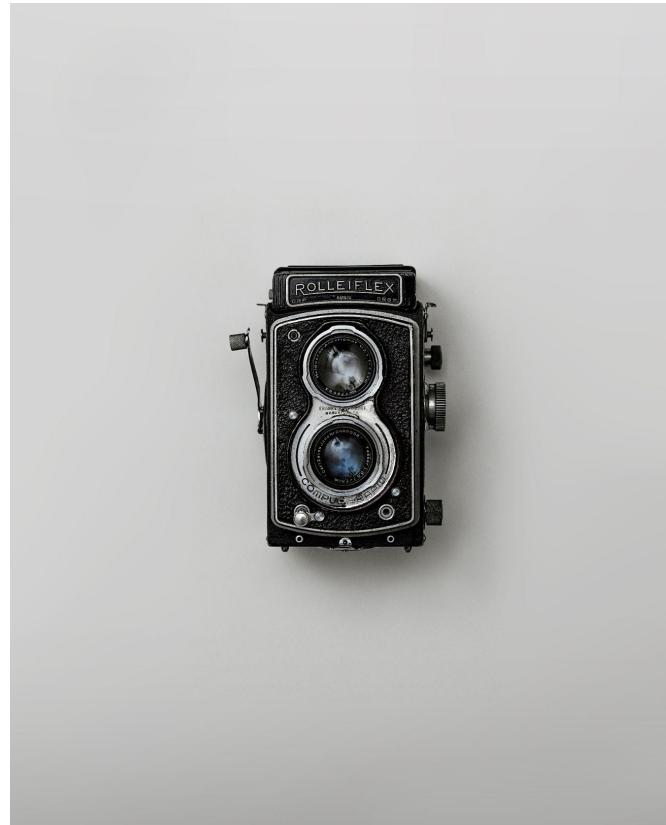
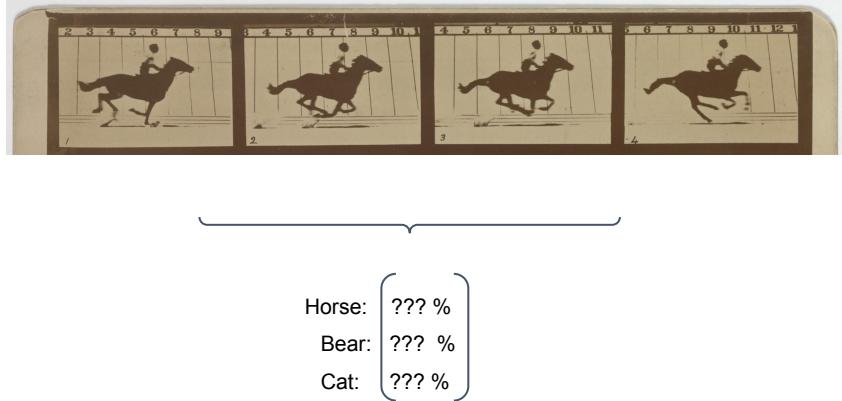


Photo by [Alexander Andrews](#) on Unsplash

Videos

Suppose we record the video of an animal and want to **identify** which animal is in the video.



[Image](#)

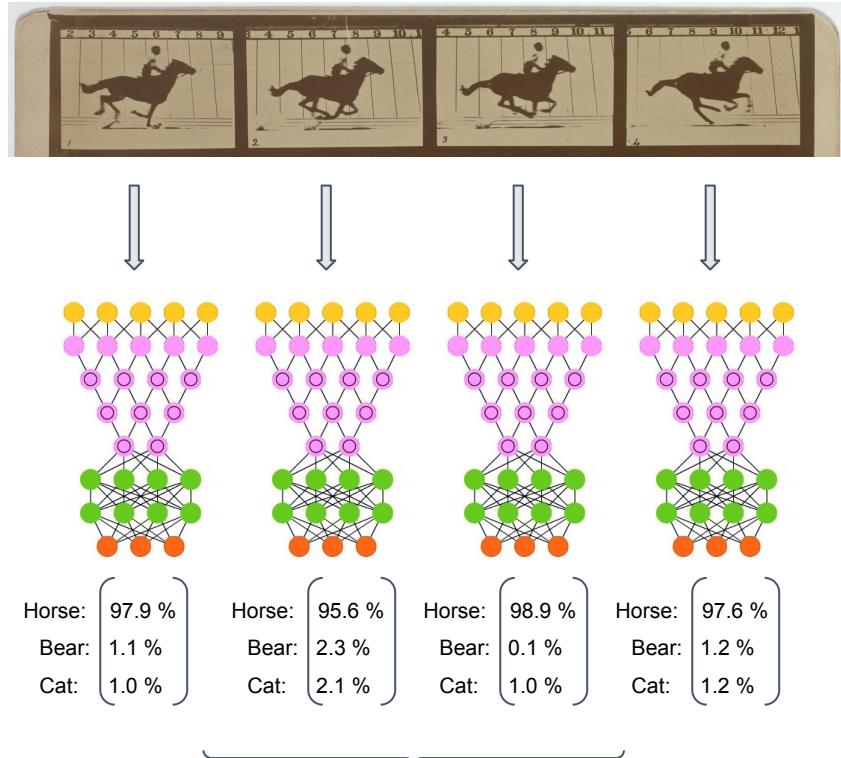
Videos

Suppose we record the video of an animal and want to **identify** which animal is in the video.

We could analyze each frame **independently**.

Every frame is passed through the **same** model, and we can **combine** the results over time to predict what is in our video.

$$p(\text{horse in video}) = (97.9 + 95.6 + 98.9 + 97.6) / 4 = 97.5\%$$



[Image](#)

Videos

This example shows segmentation in **real-time**. Each frame is segmented **independently**.

This approach **can be suitable** for **object detection** since temporal information rarely affects “objectness”.

This particular model was **trained** on single images, not videos.



[source](#)

Videos

Here, the authors compare the speed vs. accuracy **tradeoff** they can get with their models.

They show that they can achieve real-time instance segmentation with comparable precision to state-of-the-art models.

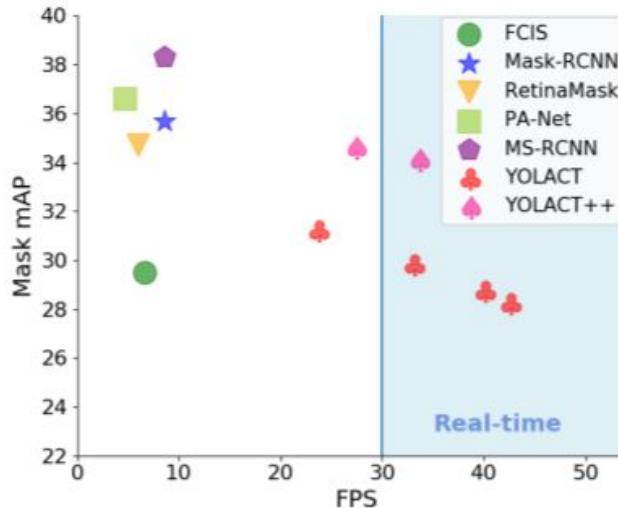


Fig. 1: Speed-performance trade-off for various instance segmentation methods on COCO. To our knowledge, ours is the first *real-time* (above 30 FPS) approach with over 30 mask mAP on COCO test-dev.

<https://arxiv.org/pdf/1912.06218.pdf>

Videos

The previous approaches assume time doesn't impact much the task at hand.

However, some tasks **require** temporal information:

- Action recognition
- Self-driving cars
- Robotics
- Video annotation (scene understanding)
- Etc.



[source](#)

Videos

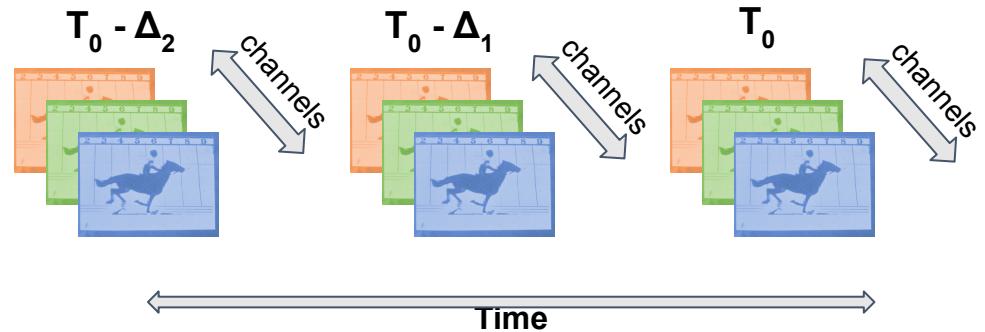
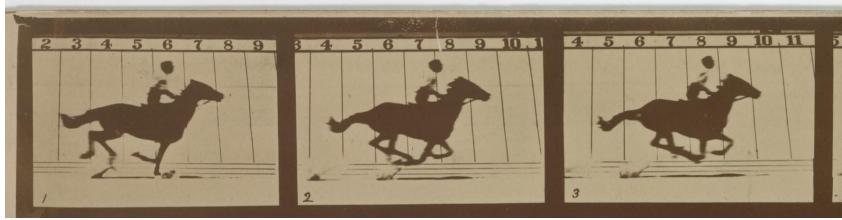
How do we **incorporate time** information to our neural networks?



Photo by [Heather Zabriskie](#) on [Unsplash](#)

Videos

Recall that an image is a **stack** of three channels (RGB) and that a video is a **stack** of multiple images (frames) over time.

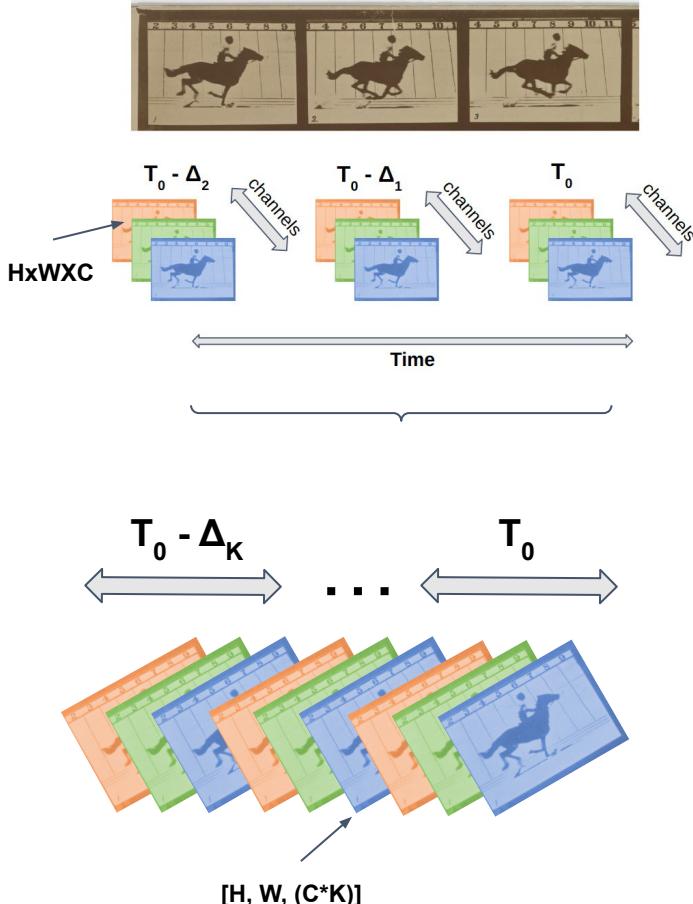


Videos

One approach is to **concatenate** past & present imagery to a single input to a CNN along the channel dimension.

Assume each image is of shape $[H, W, C]$, and we have K frames, where H and W are the height and width of the input crops, C is the number of channels, and K is the number of time frames to use.

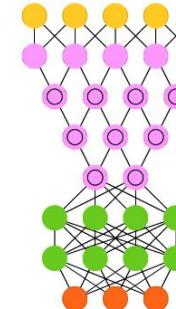
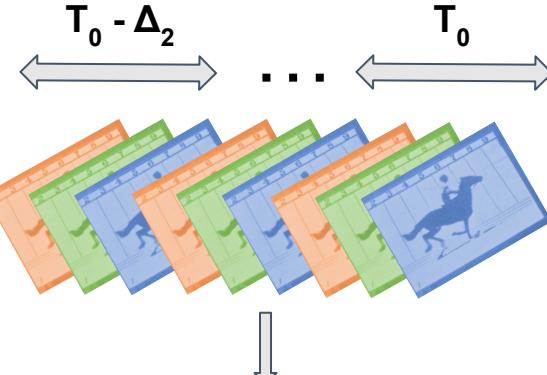
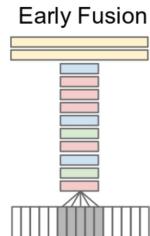
We would then have an input of shape
 $[H, W, (C*K)]$



Videos

Using this approach, you could adapt almost any image CNN classifier architecture (VGG, ResNet, etc.) by changing the number of input channels of the architecture.

This is known as **early fusion**.



VGG,
ResNet,
etc.

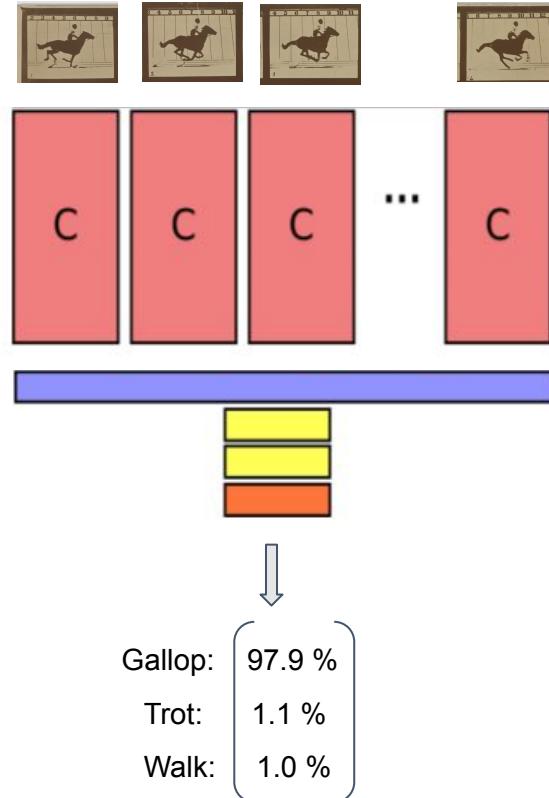
Gallop: 97.9 %
Trot: 1.1 %
Walk: 1.0 %

<http://vision.stanford.edu/pdf/karpathy14.pdf>

Videos

Alternatively, you could use a **shared CNN** backbone to extract features from independent frames and pool them later on in the network. This is known as **late fusion**.

This allows the use of **pre-trained networks** on large image datasets (e.g. ImageNet).

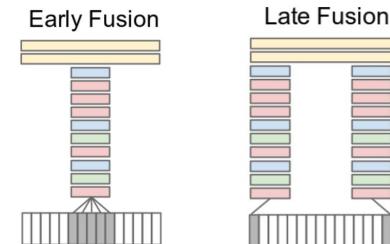
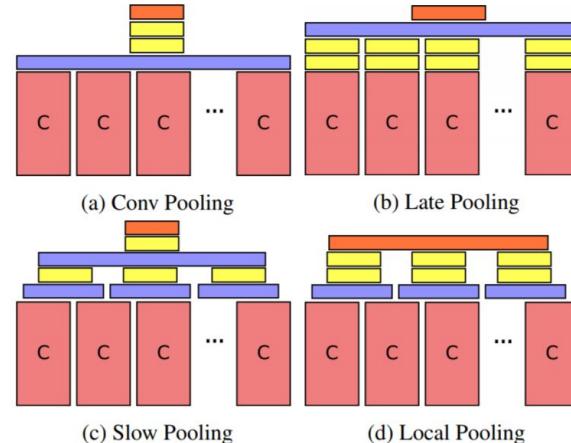


<https://arxiv.org/pdf/1503.08909.pdf>

Videos

There are **many ways** in which you could vary this setup.

You do not have to consider all **successive** frames either, as is shown in the late fusion model.



<https://arxiv.org/pdf/1503.08909.pdf>
<http://vision.stanford.edu/pdf/karpathy14.pdf>

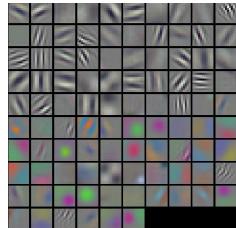
Conv3D

It can be useful to learn patterns of **motion**, just like it can be useful to learn patterns of **objects**. This can be achieved with **3D convolutions**.

2D Convolutions

Patterns of objects

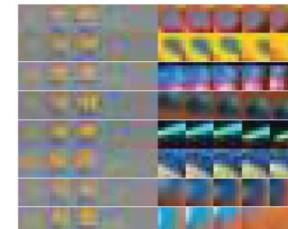
- Edges, contrasts
- ears, eyes noses



3D Convolutions

Patterns of motion

- Speed
- Movement
- Deformation
- Color changes



<https://arxiv.org/pdf/1412.0767.pdf>

Conv3D

In a **2D** convolution, the convolution operation takes place over **space** (x,y).

In a **3D** convolution, we convolve over **space and time** (t).

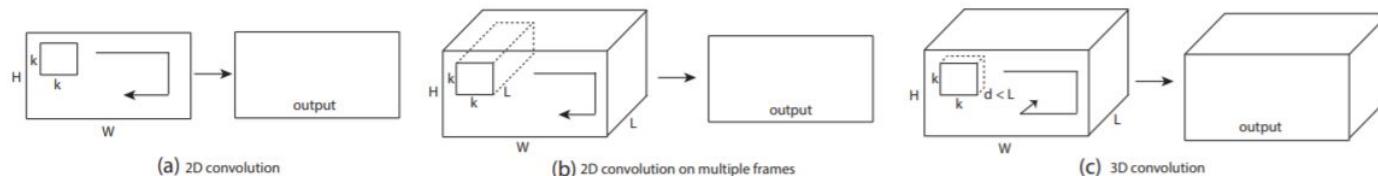


Figure 1. **2D and 3D convolution operations.** a) Applying 2D convolution on an image results in an image. b) Applying 2D convolution on a video volume (multiple frames as multiple channels) also results in an image. c) Applying 3D convolution on a video volume results in another volume, preserving temporal information of the input signal.

Conv3D

Very much like 2D Convolutions, 3D convolutions can be stacked and used sequentially:



Figure 3. **C3D architecture.** C3D net has 8 convolution, 5 max-pooling, and 2 fully connected layers, followed by a softmax output layer. All 3D convolution kernels are $3 \times 3 \times 3$ with stride 1 in both spatial and temporal dimensions. Number of filters are denoted in each box. The 3D pooling layers are denoted from pool1 to pool5. All pooling kernels are $2 \times 2 \times 2$, except for pool1 is $1 \times 2 \times 2$. Each fully connected layer has 4096 output units.

Conv3D

Here is an example of how inception-like modules can be adapted to 3D convolutions.

3D convolutions extend very naturally from 2D convolutions.

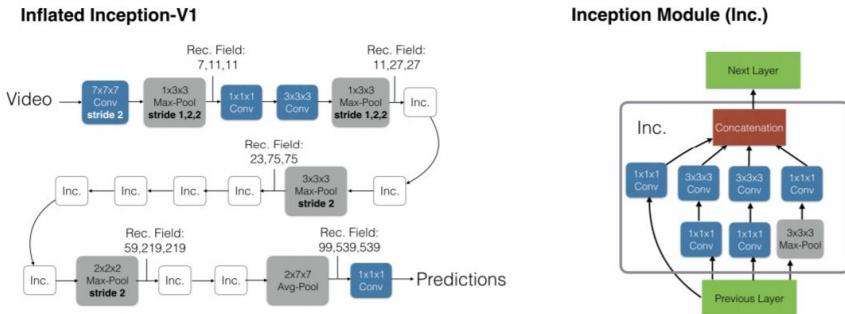


Figure 3. The Inflated Inception-V1 architecture (left) and its detailed inception submodule (right). The strides of convolution and pooling operators are 1 where not specified, and batch normalization layers, ReLu's and the softmax at the end are not shown. The theoretical sizes of receptive field sizes for a few layers in the network are provided in the format "time,x,y" – the units are frames and pixels. The predictions are obtained convolutionally in time and averaged.

<https://arxiv.org/pdf/1705.07750.pdf>

Conv3D

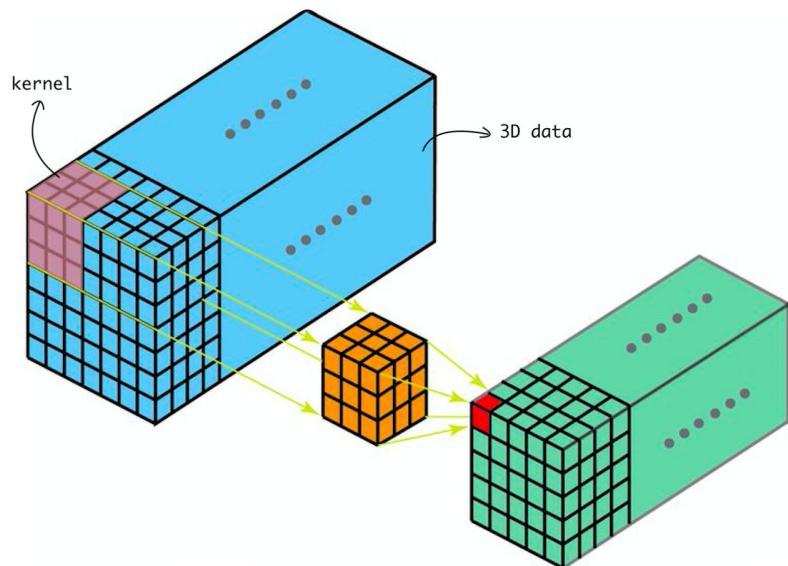
3D Convolutions are **common** in most standard deep learning APIs.

Note that in 2D convolutions, kernels are 3D. In 3D convolutions, kernels are 4D. When we stack many 3D kernels, we stack them along a 5th dimension (the output dimension).

Inputs = (batch, channels, time, height, width)

```
input shapes: torch.Size([1, 16, 30, 128, 128])
weight shapes: torch.Size([32, 16, 3, 3, 3])
output shapes: torch.Size([1, 32, 28, 126, 126])
```

```
import torch
m = torch.nn.Conv3d(16, 32, (3,3,3), stride=1)
input = torch.randn(1, 16, 30, 128, 128)
out = m(input)
```



[Source](#)

Conv3D

However, 3D convolutions are
memory intensive.

One strategy involves using 2D
convolutions on a **frame level** and
stacking features temporally before
using 3D Convolutions.

This also allows for using pretrained
imagenet networks

Model	Layer	Output Size	Pre-trained ResNet-50	VGG-M Like (Scratch)
Base Network	conv1	$112 \times 112 \times c_1$	$7 \times 7, 64, \text{stride } 2$	
			$3 \times 3 \text{ max pool, stride } 2$	
	conv2_x	$56 \times 56 \times c_2$	$\begin{bmatrix} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 1$
	conv3_x	$28 \times 28 \times c_3$	$\begin{bmatrix} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 1$
Embedder Network	conv4_x	$14 \times 14 \times c_4$	$\begin{bmatrix} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{bmatrix} \times 3$	$\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 1$
	Temporal Stacking	$k \times 14 \times 14 \times c_4$	Stack k context frame features in time axis	
	conv5_x	$k \times 14 \times 14 \times 512$	$\begin{bmatrix} 3 \times 3 \times 3, 512 \\ 3 \times 3 \times 3, 512 \end{bmatrix} \times 1$	
	Spatio-temporal Pooling	512	Global 3D Max-Pool	
	fc6_x	512	512	$\begin{bmatrix} 512 \\ 512 \end{bmatrix} \times 1$
Embedding		128	128	

Table 8: Architectures used in our experiments. The network produces an embedding for each frame (and its context window). c_i depends on the choice of the base network. Inside the square brackets, the parameters in the form of: (1) $[n \times n, c]$ refers to 2D Convolution filter size and number of channels respectively (2) $[n \times n \times n, c]$ refers to 3D Convolution filter size and number of channels respectively (3) $[c]$ refers to channels in a fully-connected layer. Downsampling in ResNet-50 is done using convolutions with stride 2, while in VGG-M models we use MaxPool with stride 2 for downsampling.

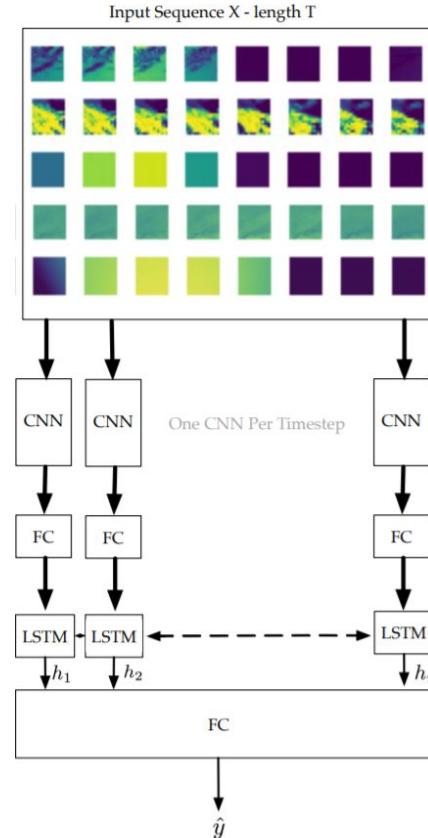
<https://arxiv.org/pdf/1904.07846.pdf>

CNN + LSTM

Another strategy is to use **2D convolution** features as an input to recurrent models (like LSTMs).

This also allows for using **pretrained** ImageNet networks.

You will see **recurrent models** in more detail on week 4 with Mirko.



<https://arxiv.org/abs/1902.01453>

ConvLSTM

2D convolutions can also be passed through a modified **LSTM** which accepts matrices as input (instead of vectors). The model was proposed to predict rainfall in the near future.

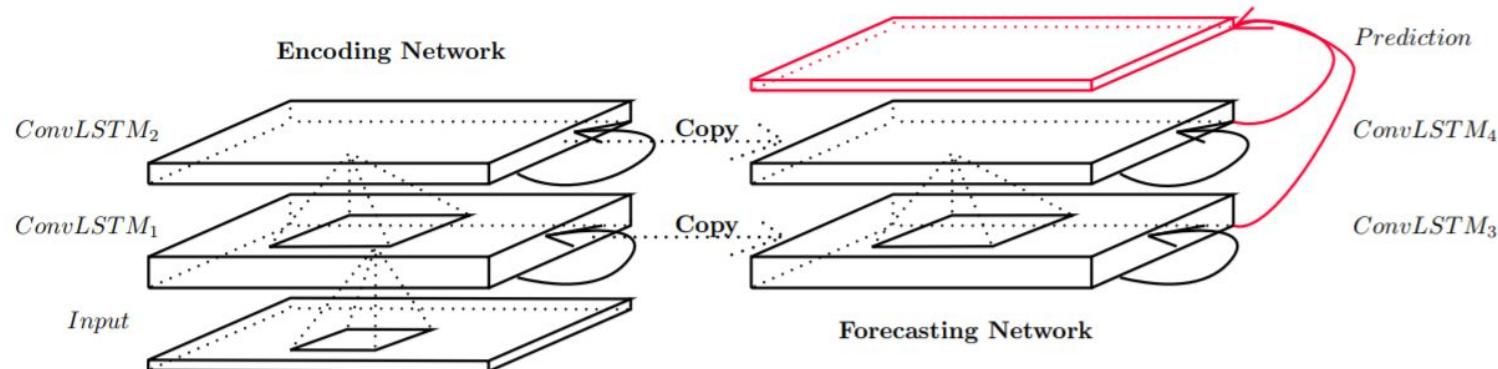
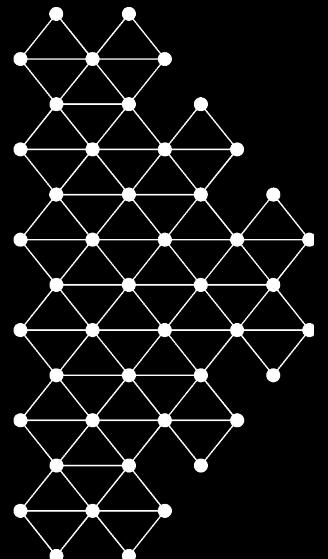


Figure 3: Encoding-forecasting ConvLSTM network for precipitation nowcasting

<https://papers.nips.cc/paper/5955-convolutional-lstm-network-a-machine-learning-approach-for-precipitation-nowcasting.pdf>

Model interpretation



Model interpretation

Due to the inherent structure of CNNs, it can be very difficult to **interpret how** they make decisions.

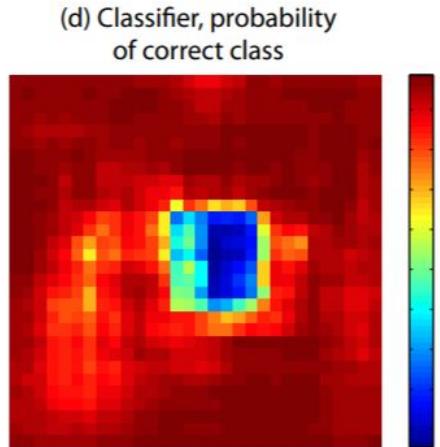
There are nevertheless certain techniques which help us **understand** what is happening **under the hood**.



Photo by [Zach Vessels](#) on [Unsplash](#)

Model interpretation

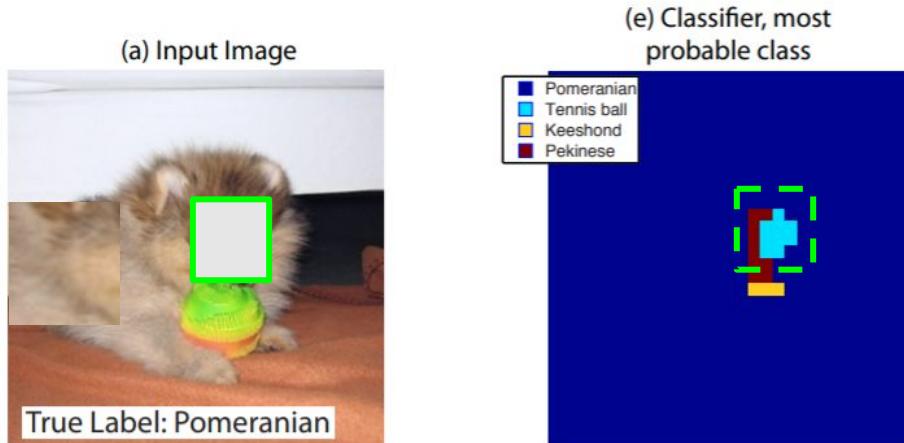
One strategy involves **blocking** out inputs of an image and looking at how the predictions vary over the network. These are known as **perturbation methods**.



<https://arxiv.org/pdf/1311.2901.pdf>

Model interpretation

When we block out the face of the dog, the **tennis ball** get's predicted with high probability.

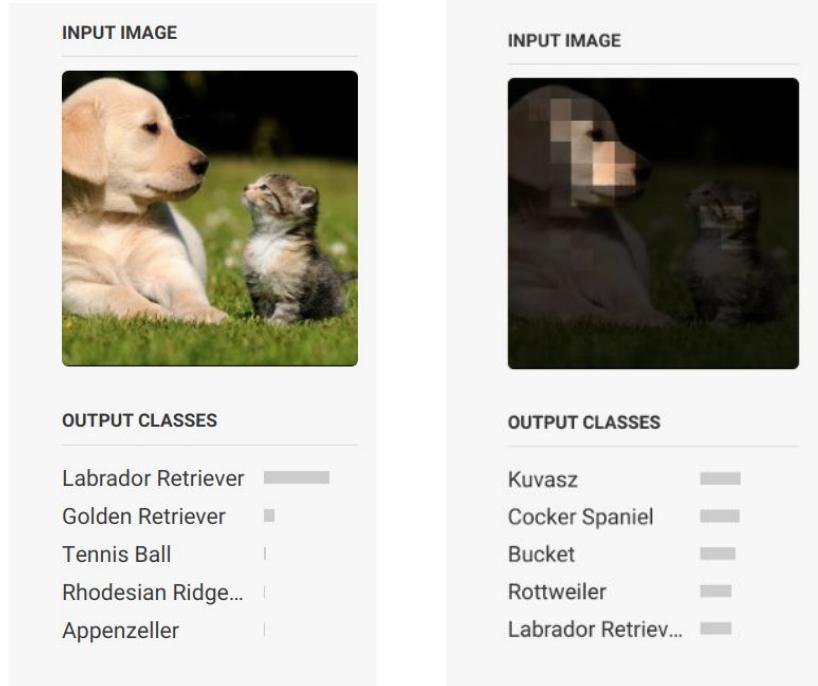


<https://arxiv.org/pdf/1311.2901.pdf>

Model Interpretation

We can also use the activation from a given layer to correlate the parts of an image with its predicted classes.

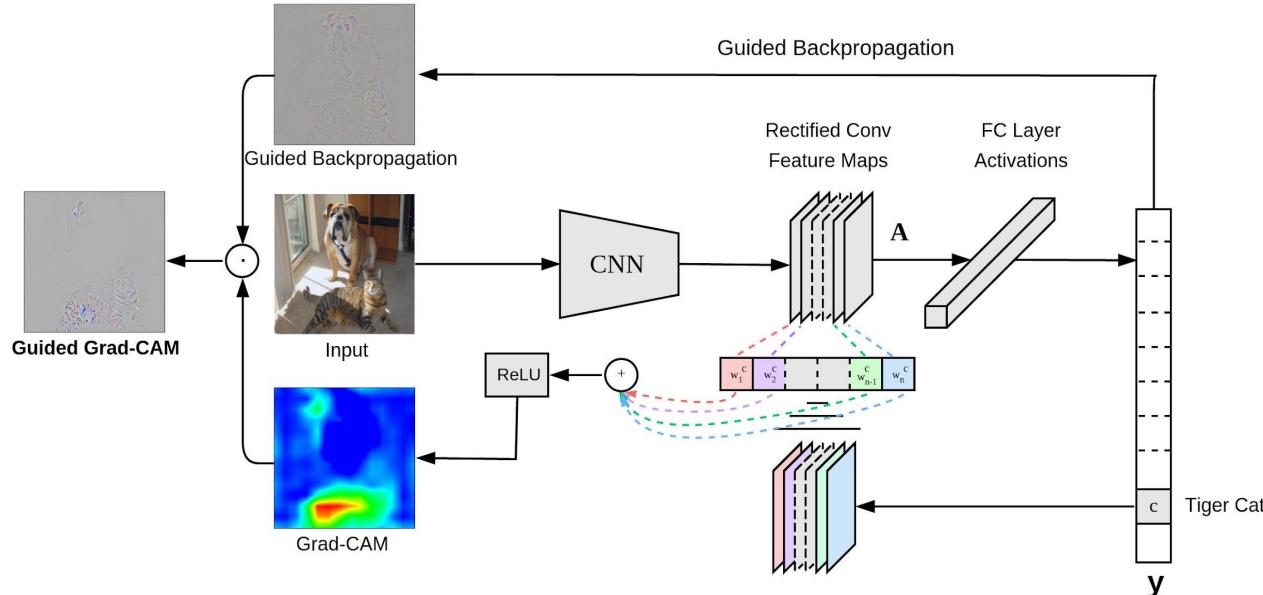
We see here that when focusing on a dog's face, the network hesitates between various dog classes.



<https://distill.pub/2018/building-blocks/>

Model interpretation

Other methods, such as **GradCam**, use the information of the **gradient** of a single-class prediction to interpret what the model is looking at.



<https://arxiv.org/abs/1610.02391>

Model interpretation

These methods can be very useful to **identify bias** attributed to certain classes in a network.

Here, we see that in a biased model, a network focuses on the face of a female doctor when incorrectly attributing the label “nurse”.

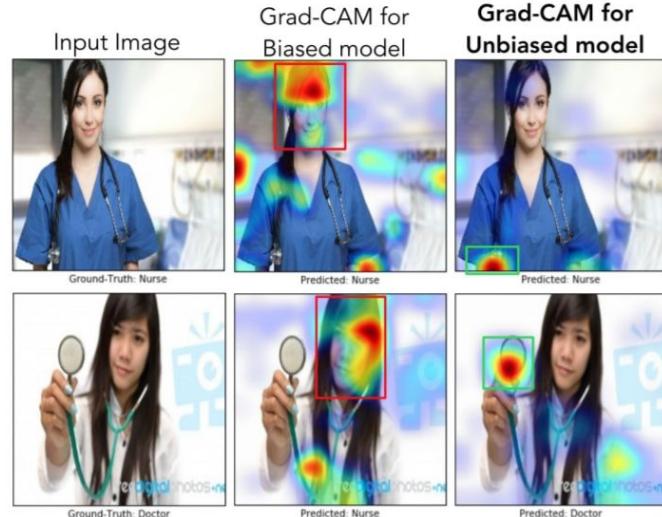
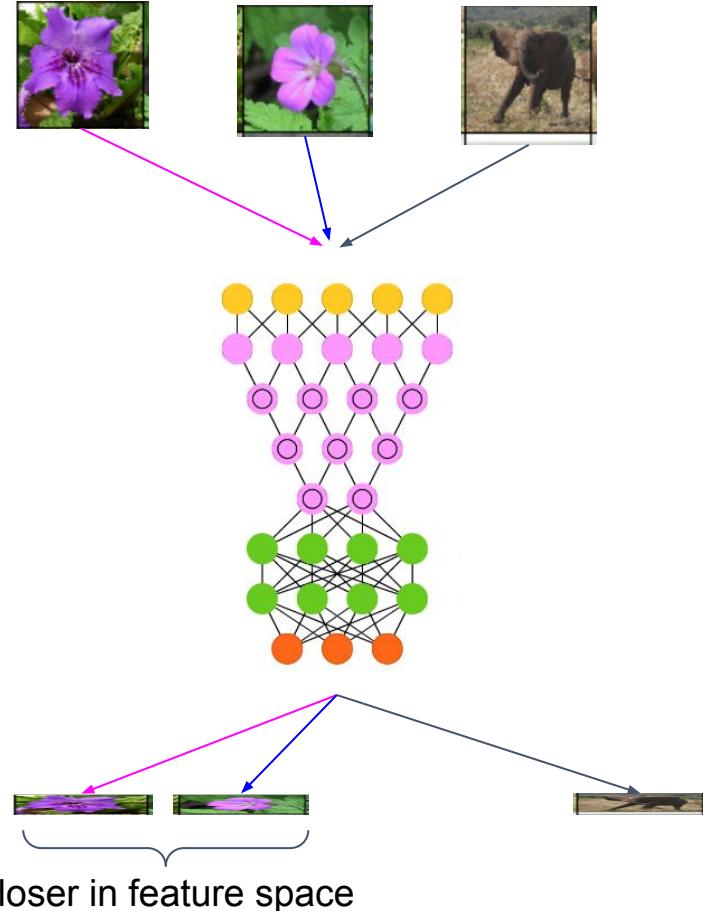


Fig. 8: In the first row, we can see that even though both models made the right decision, the biased model (model1) was looking at the face of the person to decide if the person was a nurse, whereas the unbiased model was looking at the short sleeves to make the decision. For the example image in the second row, the biased model made the wrong prediction (misclassifying a doctor as a nurse) by looking at the face and the hairstyle, whereas the unbiased model made the right prediction looking at the white coat, and the stethoscope.

Model interpretation

Another useful trick to visualize what your model might be picking up on is by analyzing the **feature space**.

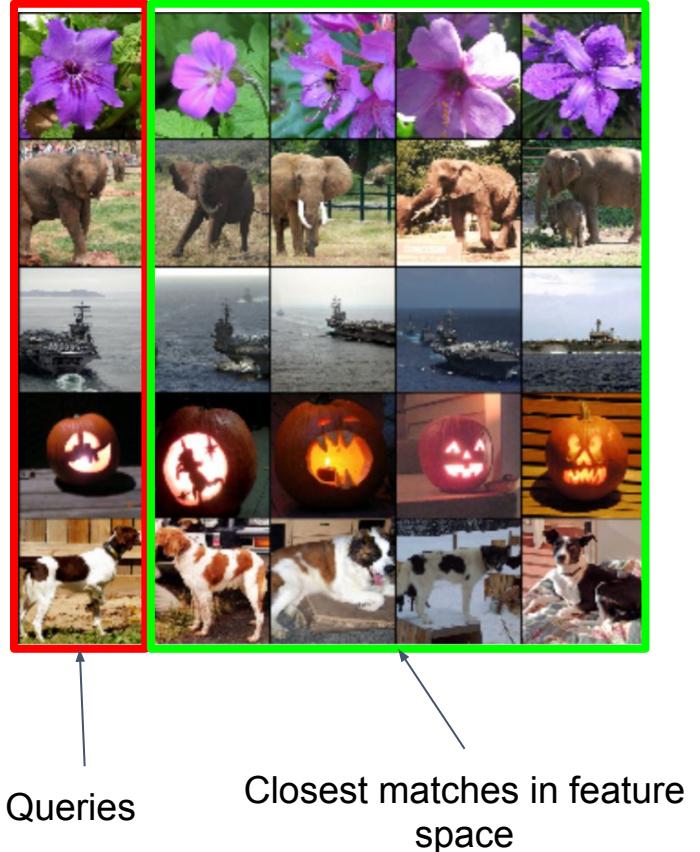
For example, in AlexNet, the authors reported the images that were closest in the **feature space** (last hidden layer) to the query sample.



Model interpretation

Analyzing the correlations in feature space can help **debug** if your network is learning “meaningful” representations or not for your given task.

Here, a dot product is used to find the most **similar samples** to the query in the feature space.

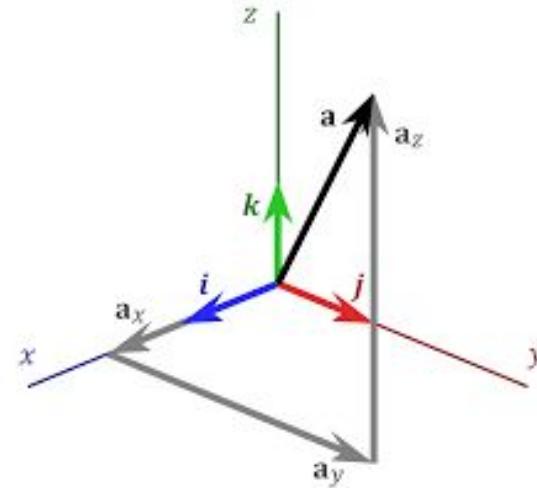


Dimensionality Reduction

However, these **feature spaces** can be very large. In AlexNet, the last dense layer is a vector with 2048 dimensions.

These numbers are **abstract**, non-linear relationships that the network learns. They are very difficult to interpret.

Dimensionality reduction can help to untangle these high dimensional mappings.



Dimensionality Reduction

Dimensionality reduction can really be a useful tool to see if **meaningful clusters** emerge.

So **how** does it work?

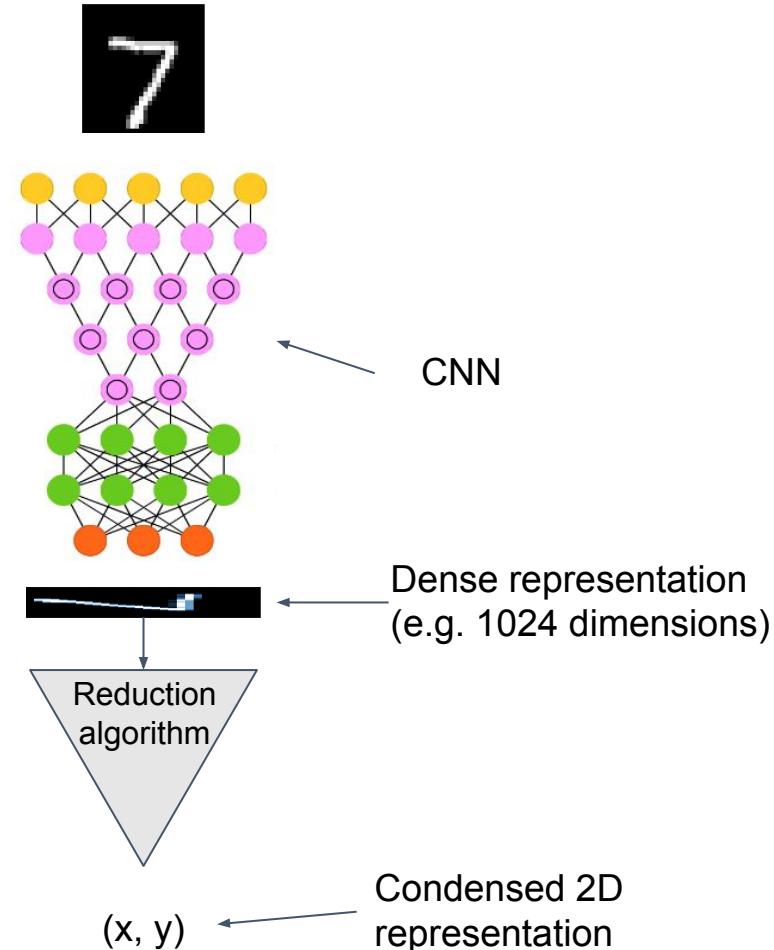


[Image](#)

Dimensionality Reduction

Dimensionality reduction allow us to preserve “**useful**” relationships while **significantly reducing** the number of dimensions of the vectors.

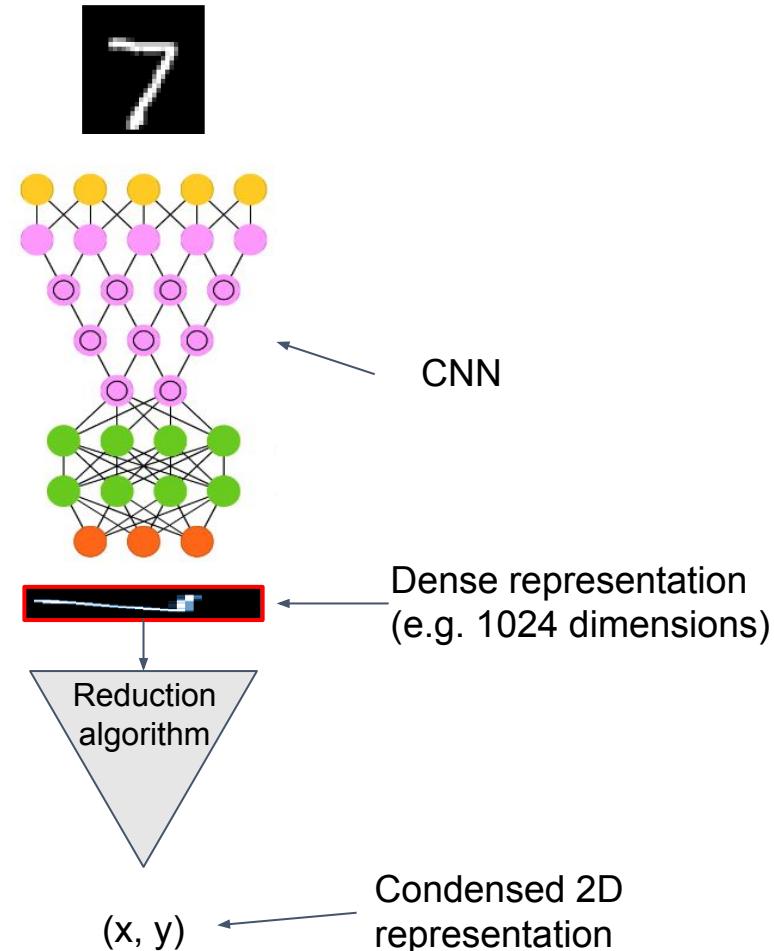
Dimensionality reduction can be used on the **dense representation** of the input - typically the last layers of a CNN.



Dimensionality Reduction

Some commonly used algorithms are:

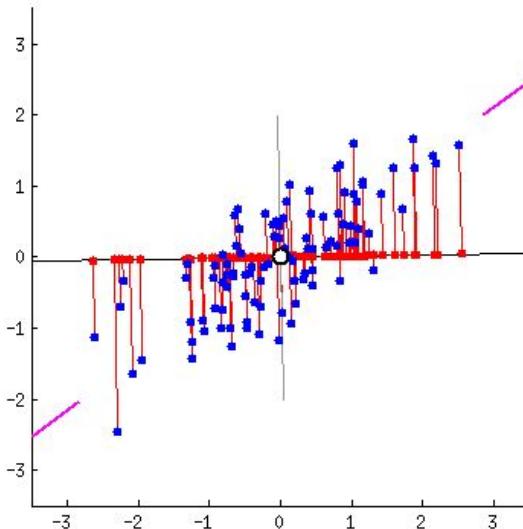
- Statistical methods
 - PCA (linear)
 - T-sne (non-linear)
- Learned-based methods
 - Autoencoders (non-linear)



PCA

Principal component analysis (PCA) attempts to **reduce** the number of features representing while **maximizing the variance** of the dataset.

In this example, the “purple” direction represents the **direction** along which **maximal information** is preserved.



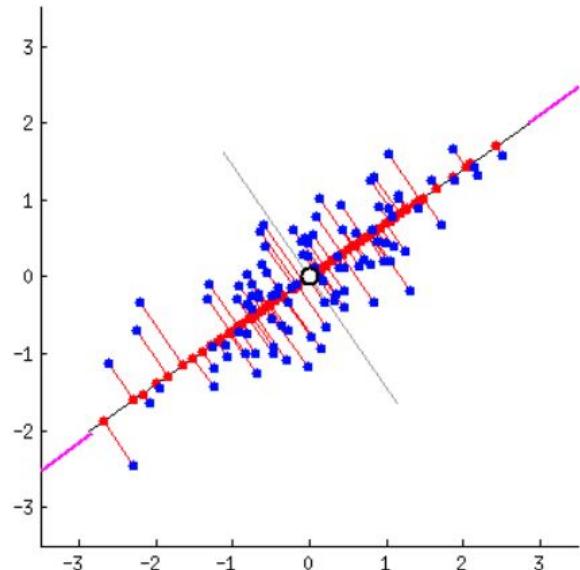
[source](#)

PCA

Principal component analysis (PCA) attempts to **reduce** the number of features representing while **maximizing the variance** of the dataset.

In this example, the “purple” direction represents the **direction** along which **maximal information** is preserved.

We can think of reducing our current scatter (x,y) to the form $y = ax + b$, going from 2 variables, to 1 variable.



[source](#)

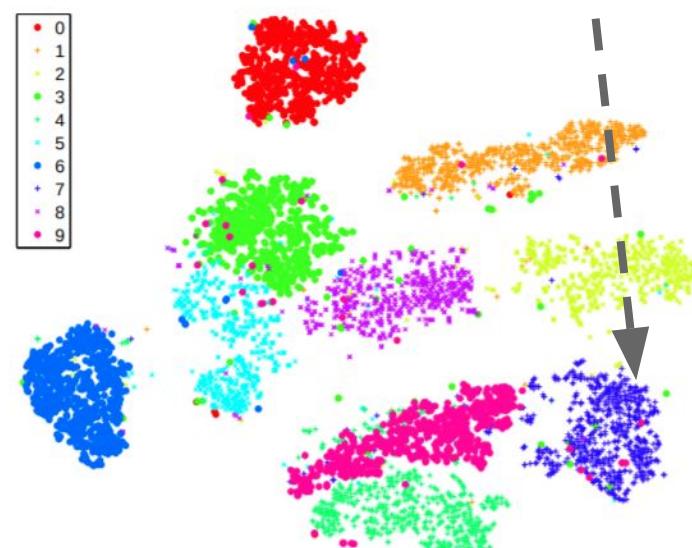
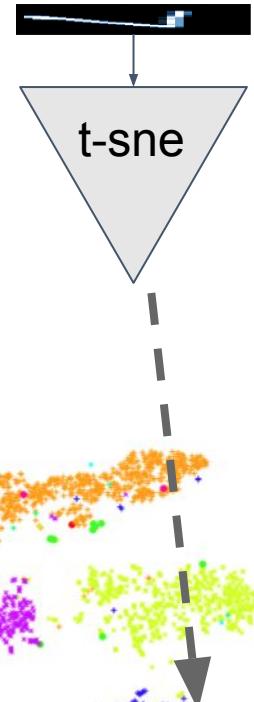
t-sne

T-SNE looks for mappings that **preserve** the relative distances between samples across feature spaces.

Here, we **reduce** the feature vectors in MNIST (e.g. 1568 in LeNet) to 2 dimensions using t-sne.

By **plotting** the 2 features and color-coding by category, we can **observe clusters** forming.

Dense representation



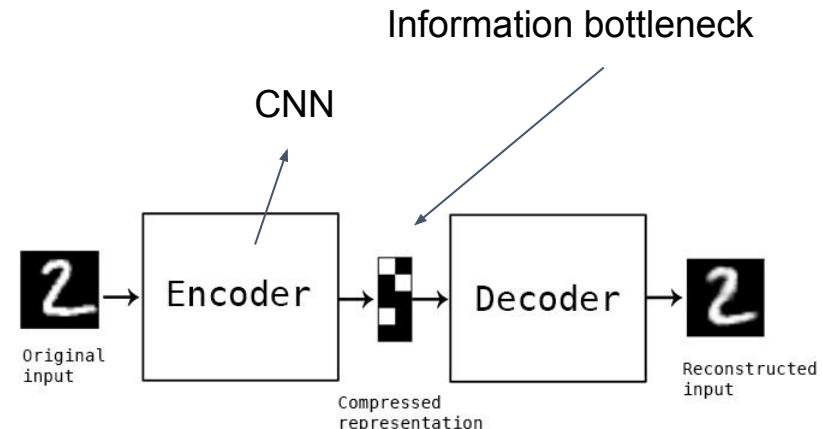
https://lvdmaaten.github.io/publications/papers/JMLR_2008.pdf

Autoencoders

Autoencoders learn directly from the data to **reduce** the feature representations.

Given an input, we learn to **reconstruct** the exact same output but pass the input through an **information bottleneck**.

Once the network is trained, the **encoder** representation can be used in downstream tasks.

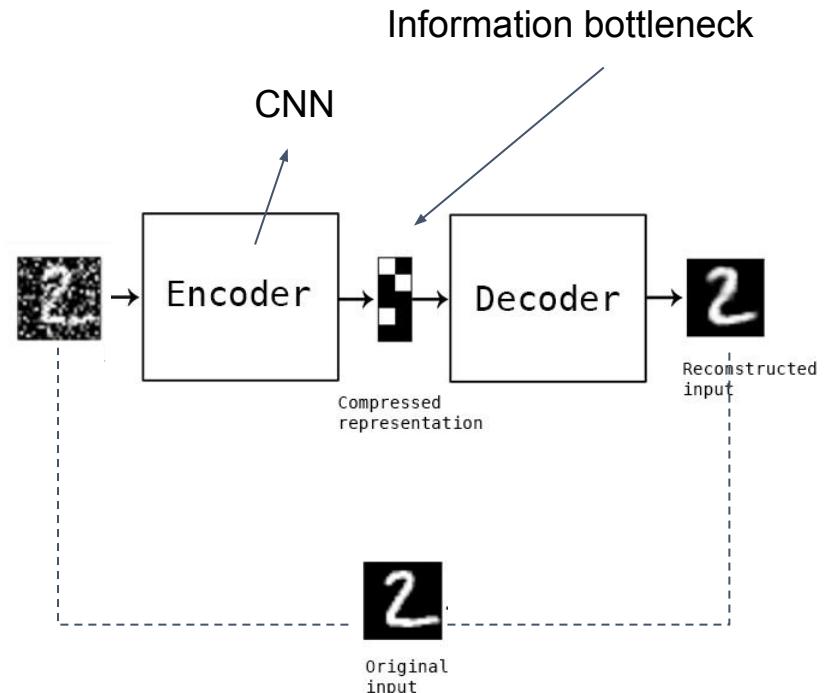


[source](#)

Denoising Autoencoder

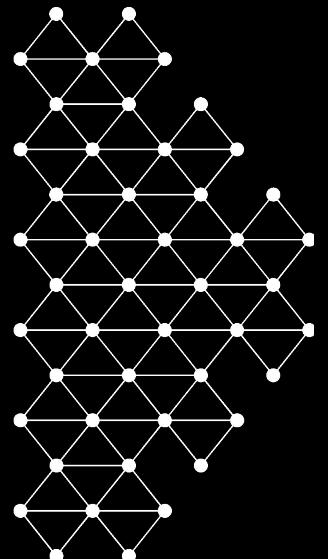
In a denoising autoencoder,
we **add noise** to our input and learn to
reconstruct the clean original input.

This gives us a more robust compressed representation - the network needs to learn to focus on **important details** while filtering out the noise.



[source](#)

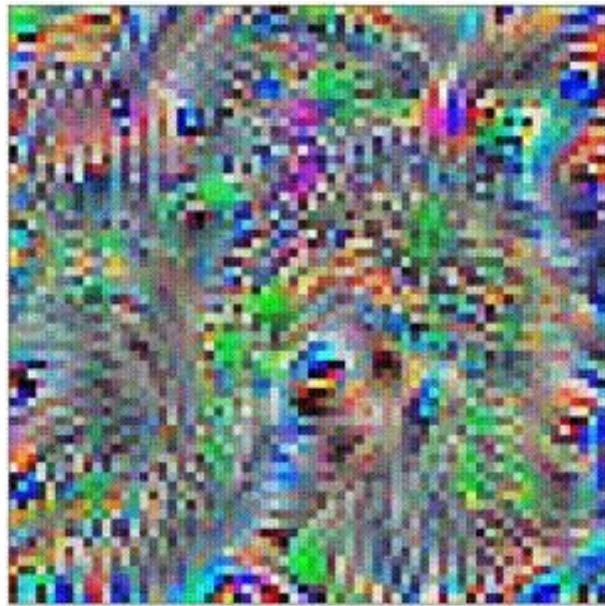
Feature Visualization



Feature Visualization

We've seen through dimensionality reduction how we can observe the clustering of our data.

However, there are also techniques to **peek** at what our neural network has learned.



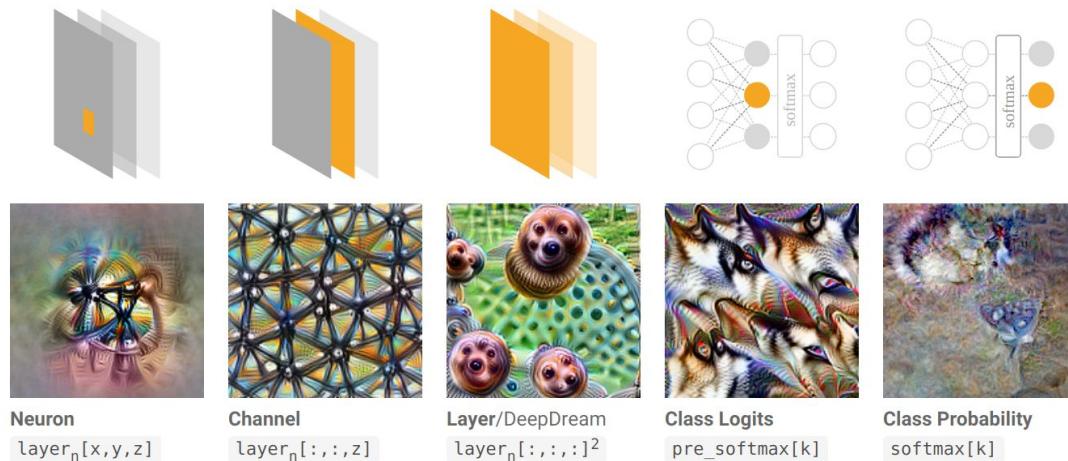
[source](#)

Feature Visualization

We can visualize what activates **components** in various areas of a network using **optimization**. It is helpful to “see” what the network might be thinking.

Different optimization objectives show what different parts of a network are looking for.

- layer index
- x, y spatial position
- z channel index
- k class index

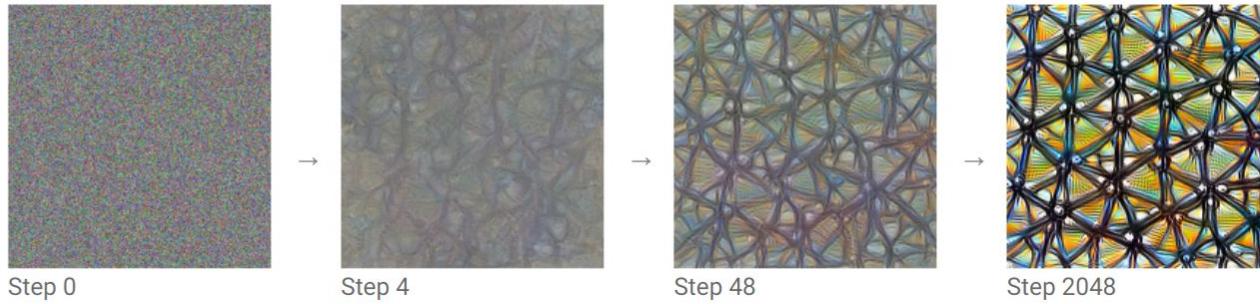


<https://distill.pub/2017/feature-visualization/>

Feature Visualization

Starting from a trained network, we use **noise** as an input and **optimize** the input to maximally activate the neuron we are seeking to interpret.

Starting from random noise, we optimize an image to activate a particular neuron (layer mixed4a, unit 11).



<https://distill.pub/2017/feature-visualization/>

Feature Visualization

We can use this **optimization** as a tool to probe what causes certain neurons to fire in our network.

Dataset Examples show us what neurons respond to in practice

Optimization isolates the causes of behavior from mere correlations. A neuron may not be detecting what you initially thought.



Baseball—or stripes?
mixed4a, Unit 6

Animal faces—or snouts?
mixed4a, Unit 240

Feature Visualization

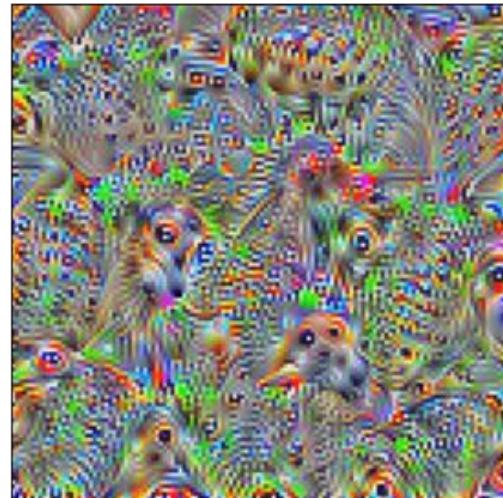
This can be implemented in just a few lines of pytorch code - however there are a lot of “tricks” that need to be implemented to make it look pretty.

```
class Vgg16(torch.nn.Module):
    def __init__(self, stop_layer, device, pretrained=True):
        super(Vgg16, self).__init__()
        vgg16 = models.vgg16(pretrained=pretrained)
        vgg16 = vgg16.to(device)
        self.features = torch.nn.Sequential(
            *list(vgg16.features.children())[:stop_layer]
        )
    def forward(self, x):
        return self.features(x)

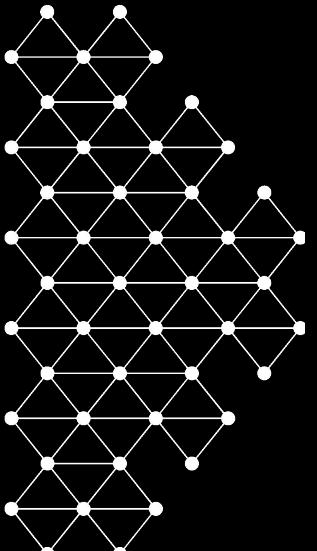
device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")

lr = 0.05
wd = 1e-6
n_iters = 2000
batch_size = 1
stop_layer = 24
target = init_target(width=128, height=128, device=device) # normalized random noise
model = Vgg16(stop_layer=stop_layer, pretrained=True, device=device)
optimizer = torch.optim.Adam([target], lr=lr, weight_decay=wd)
for n_iter in tqdm(range(1, n_iters + 1)):
    optimizer.zero_grad()

    out = model.forward(target)
    loss = -torch.mean(out[:, :, :, :, :])
    loss.backward()
    optimizer.step()
```



[source](#)

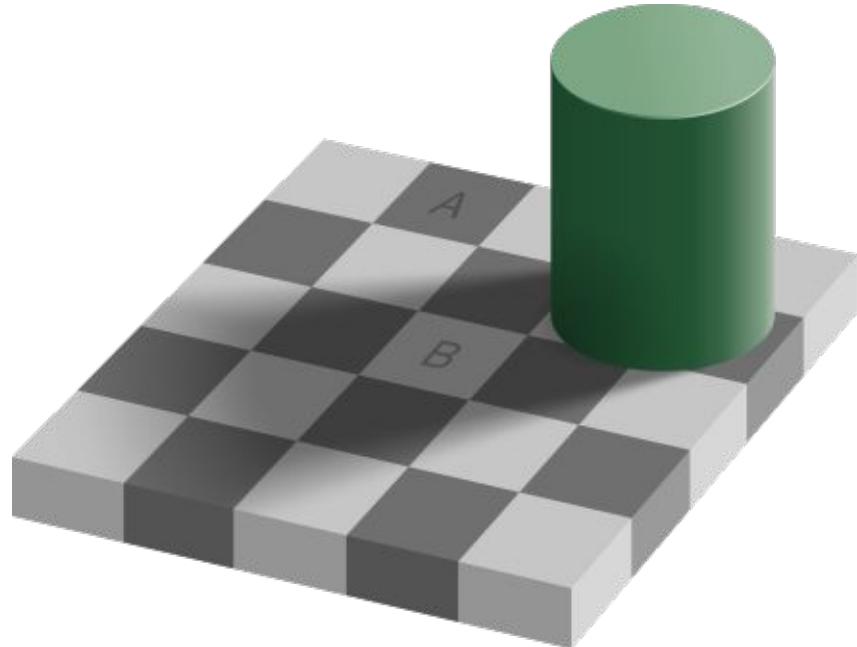


Adversarial Attacks

Adversarial Attacks

Our brain can easily be fooled using trickery. For example, look at this grid:

Are the colours of square A and B identical?

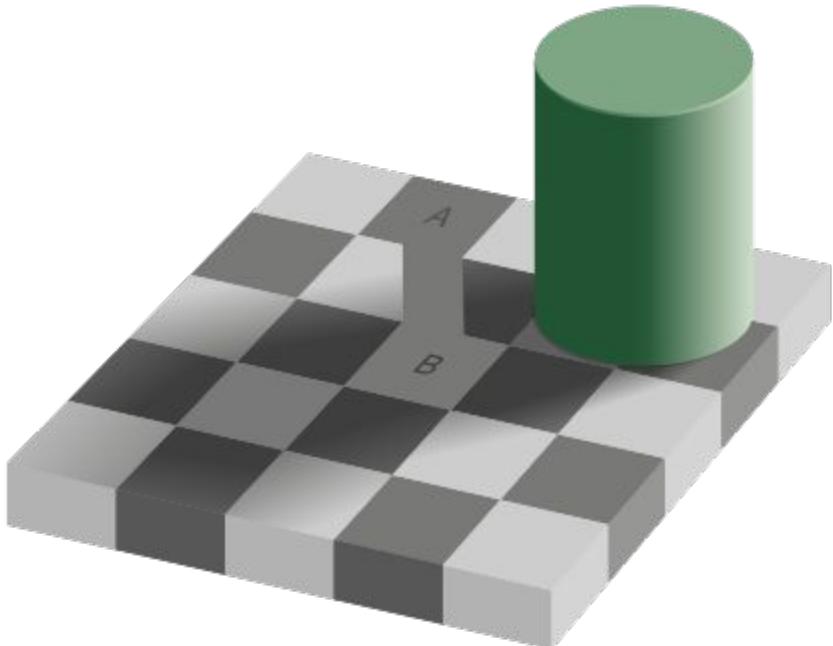


[source](#)

Adversarial Attacks

It turns out that they are. No **trickery** was used.

By simply playing with the **shadows**, our brains get fooled due to our preconceived expectations of what a regular grid should look like.



[source](#)

Adversarial Attacks

It is equally possible to **fool** networks with adversarial attacks. Here, we use a noisy perturbation based on the gradient of the network. We can easily overturn the classification of a class in a network:



\mathbf{x}
“panda”
57.7% confidence

$+ .007 \times$



$\text{sign}(\nabla_{\mathbf{x}} J(\boldsymbol{\theta}, \mathbf{x}, y))$
“nematode”
8.2% confidence

=

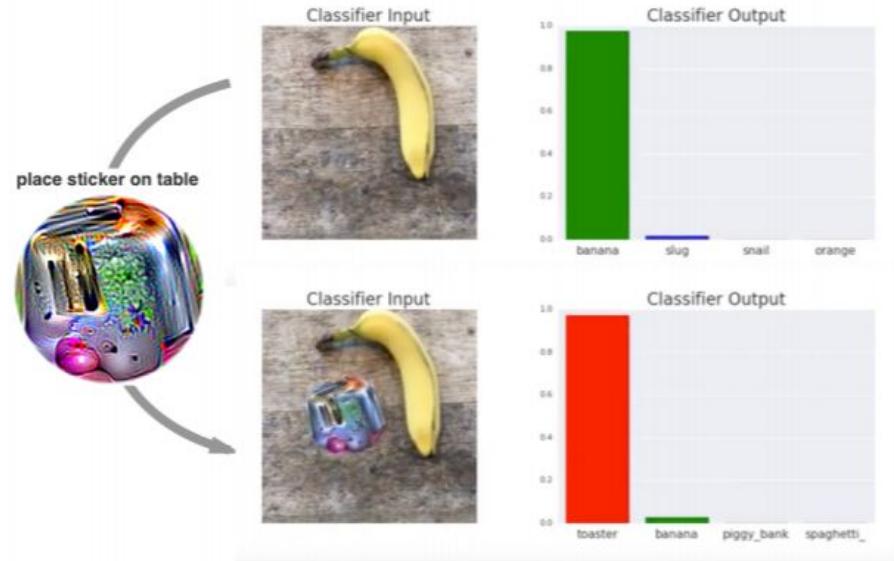


$\mathbf{x} + \epsilon \text{sign}(\nabla_{\mathbf{x}} J(\boldsymbol{\theta}, \mathbf{x}, y))$
“gibbon”
99.3 % confidence

<https://arxiv.org/pdf/1412.6572.pdf>

Adversarial Attacks

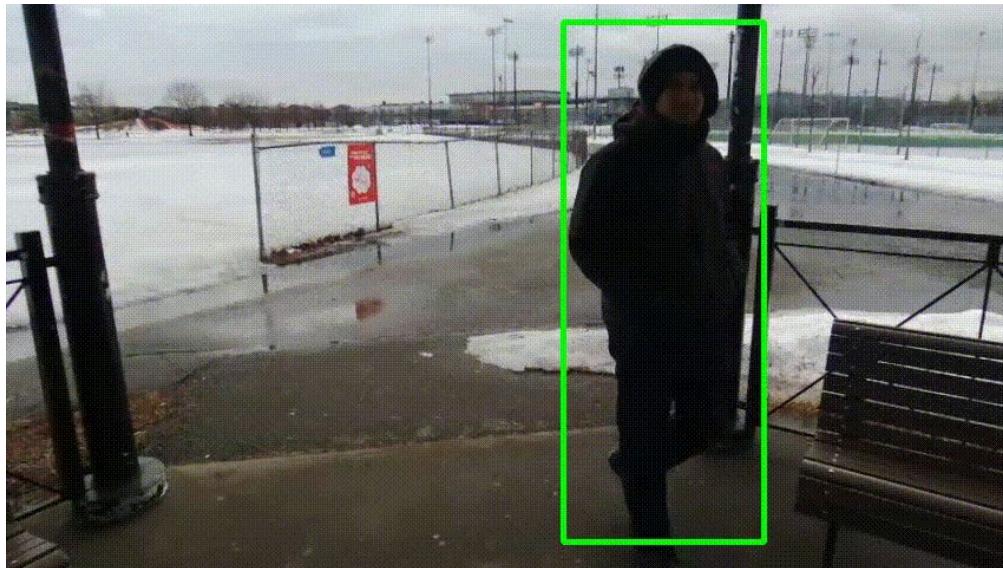
In this example, a carefully crafted **sticker** can be used to fool a detector to always output a toaster.



<https://arxiv.org/pdf/1712.09665.pdf>

Adversarial Attacks

Here, a physical object is placed in the real world and “intercepts” the tracking of a person.



[source](#)

Adversarial Attacks

These attacks can have **catastrophic** consequences in real-world environments - e.g. self-driving cars.

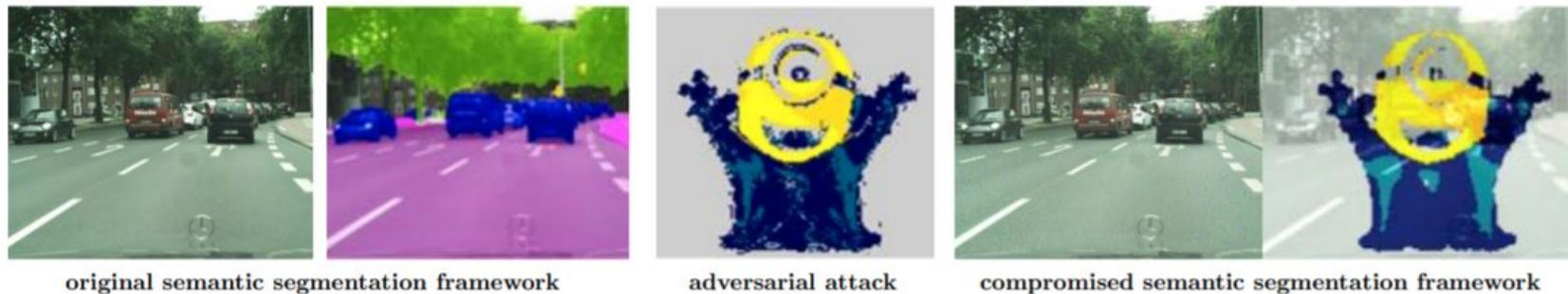
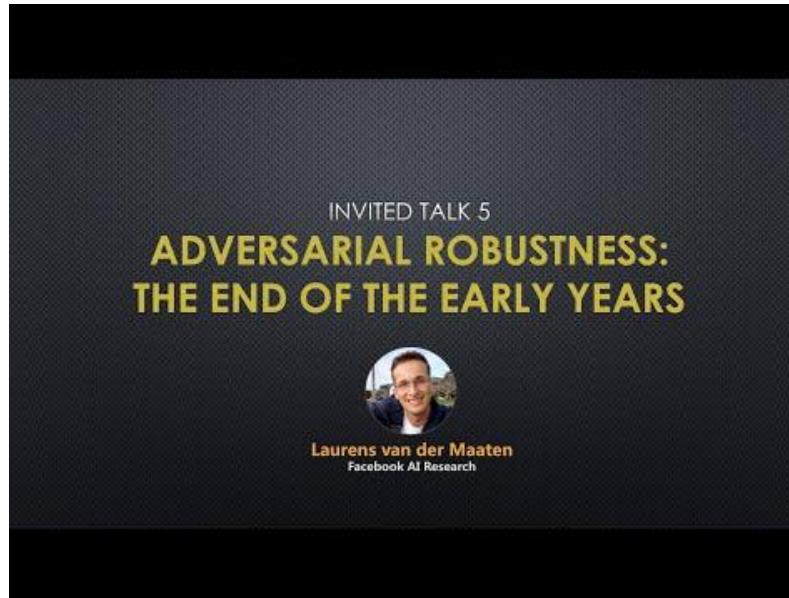


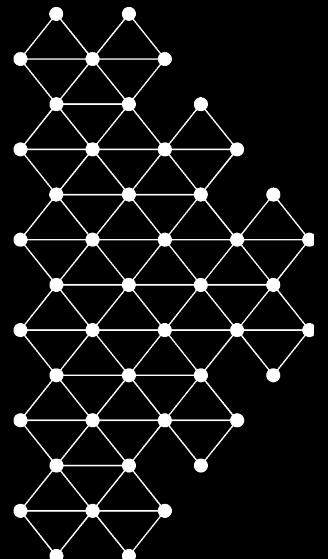
Figure 1: We cause the network to generate a *minion* as segmentation for the adversarially perturbed version of the original image. Note that the original and the perturbed image are indistinguishable.

<https://arxiv.org/pdf/1707.05373.pdf>

Adversarial Attacks

Of course, adversarial attacks can be **mitigated** against. Here is a talk highlighting some of the techniques that can be used to prevent attacks.



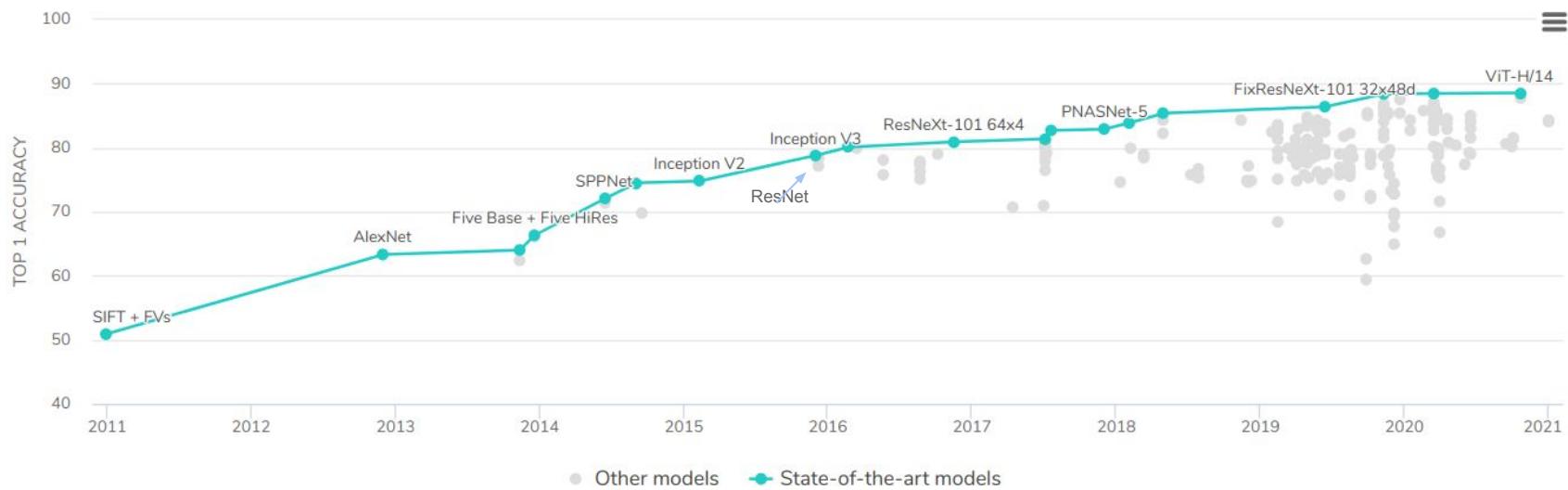


ImageNet (Updates)

ImageNet

Although the ILSVRC annual competition is **officially over**, ImageNet remains to this day an important benchmark in literature.

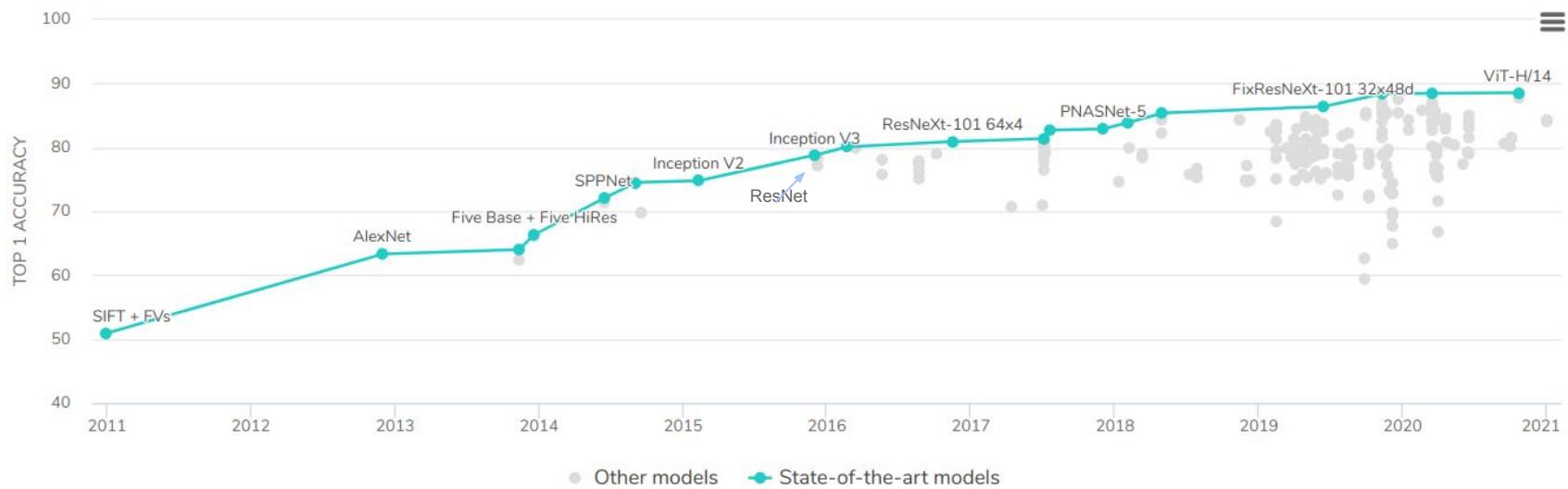
Image Classification on ImageNet



ImageNet

Today we will review some new models and paradigms to have achieved **state-of-the-art** on ImageNet.

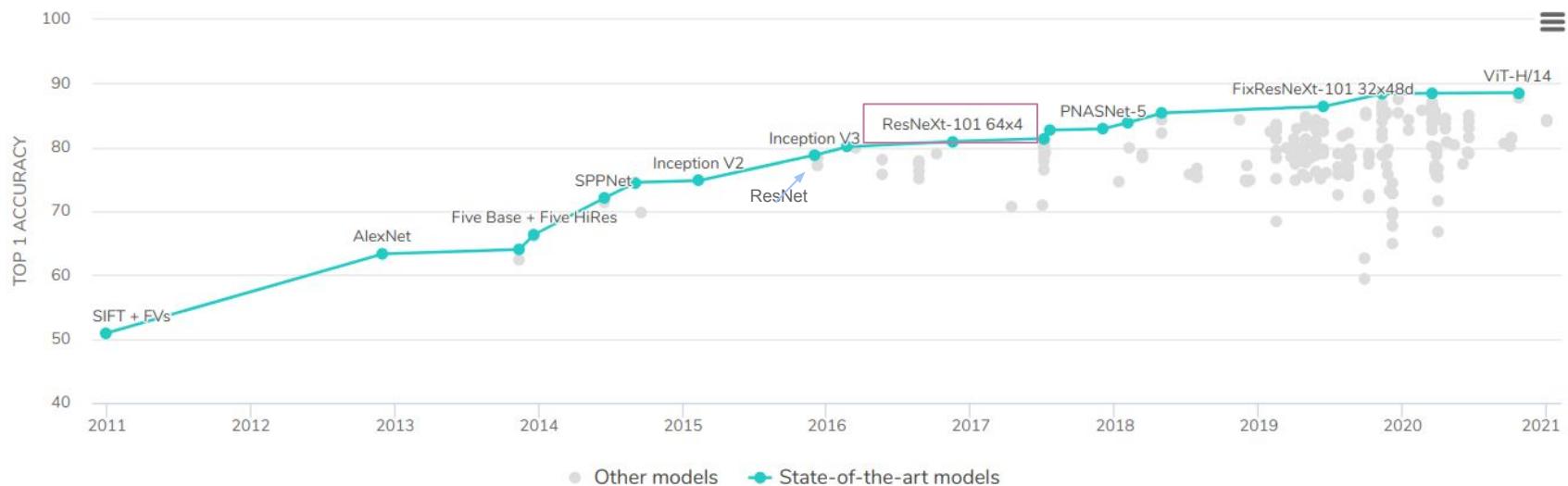
Image Classification on ImageNet



ImageNet

Today we will review some new models and paradigms to have achieved **state-of-the-art** on ImageNet. We will begin with **ResNeXt**.

Image Classification on ImageNet

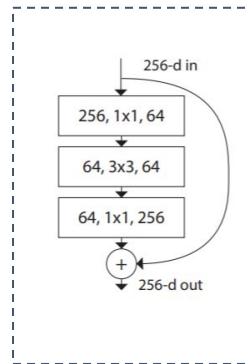


ResNeXt

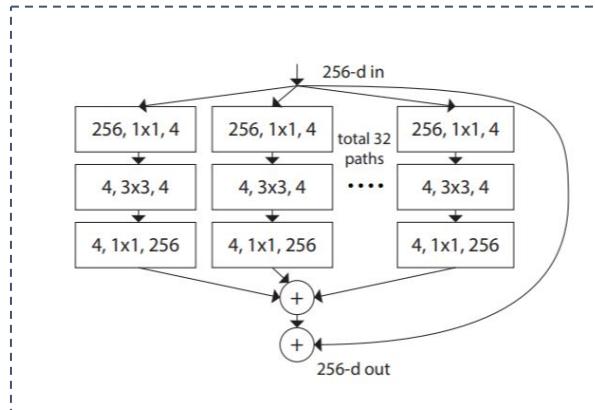
ResNeXt is a direct enhancement of ResNet.

Similar to inception, it has a module with parallel paths. It is designed to preserve the number of FLOPS in each block when compared to a ResNet.

Similar to ResNet, we have the **skip connection**. However, instead of having a single linear path, we now have parallel paths, denoted as *cardinality*.



ResNet block



ResNeXt block

<https://arxiv.org/pdf/1512.03385.pdf>
<https://arxiv.org/pdf/1611.05431.pdf>

ResNeXt

It can be thought of as having **multiple networks** in parallel on the same input (grouped convolutions).

It has been shown to **outperform** its ResNet counterpart while being able to preserve the amount of overall parameters in the network.

It was used as a backbone in the original Mask R-CNN implementation.

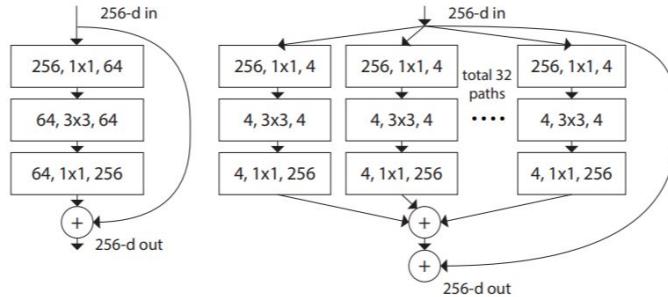


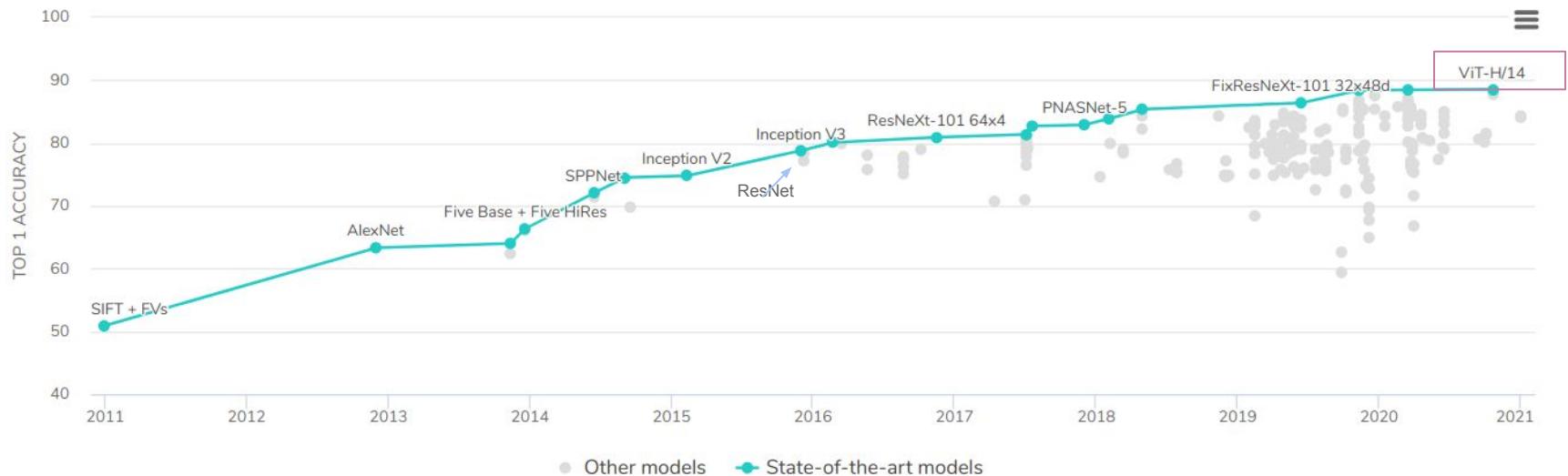
Figure 1. **Left:** A block of ResNet [14]. **Right:** A block of ResNeXt with cardinality = 32, with roughly the same complexity. A layer is shown as (# in channels, filter size, # out channels).

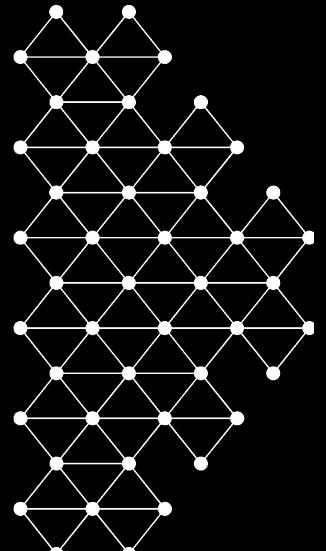
	setting	top-1 error (%)
ResNet-50	1 × 64d	23.9
ResNeXt-50	2 × 40d	23.0
ResNeXt-50	4 × 24d	22.6
ResNeXt-50	8 × 14d	22.3
ResNeXt-50	32 × 4d	22.2
ResNet-101	1 × 64d	22.0
ResNeXt-101	2 × 40d	21.7
ResNeXt-101	4 × 24d	21.4
ResNeXt-101	8 × 14d	21.3
ResNeXt-101	32 × 4d	21.2

ImageNet top-1

Currently, state-of-the-art belongs to ViT - Vision Transformer.

Image Classification on ImageNet





Transformers and computer vision

Farewell Convolutions?

A new paradigm is starting to emerge in computer vision. It involves adapting **transformers** to images instead of using CNNs.



Andrej Karpathy @karpathy · Oct 3
An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale openreview.net/forum?id=YicbF... v cool. Further steps towards deprecating ConvNets with Transformers. Loving the increasing convergence of Vision/NLP and the much more efficient/flexible class of architectures.

	Ours (ViT-H14)	Ours (ViT-L16)	BiT-L (ResNet152x4)	Noisy Student (EfficientNet-L2)
ImageNet	88.36	87.61 ± 0.03	87.54 ± 0.02	88.4 / 88.5*
CIFAR-10 ReId	90.77	90.24 ± 0.03	90.54	90.55
CIFAR-100	99.60 ± 0.06	99.58 ± 0.03	99.53 ± 0.06	—
Oxford-IIIT Pets	94.55 ± 0.04	93.90 ± 0.05	93.51 ± 0.08	—
Oxford Flowers-102	97.56 ± 0.03	97.32 ± 0.11	96.62 ± 0.23	—
VTAB (19 tasks)	99.68 ± 0.02	99.74 ± 0.00	99.63 ± 0.03	—
TPUv3-days	77.16 ± 0.29	75.91 ± 0.18	76.29 ± 1.70	—

Table 2: Comparison with state of the art on popular image classification datasets benchmarks. Vision Transformer models pre-trained on the JFT300M dataset often match or outperform ResNet-based baselines while taking substantially less computational resources to pre-train. *Slightly improved 88.5% result reported in Touvron et al. (2020).

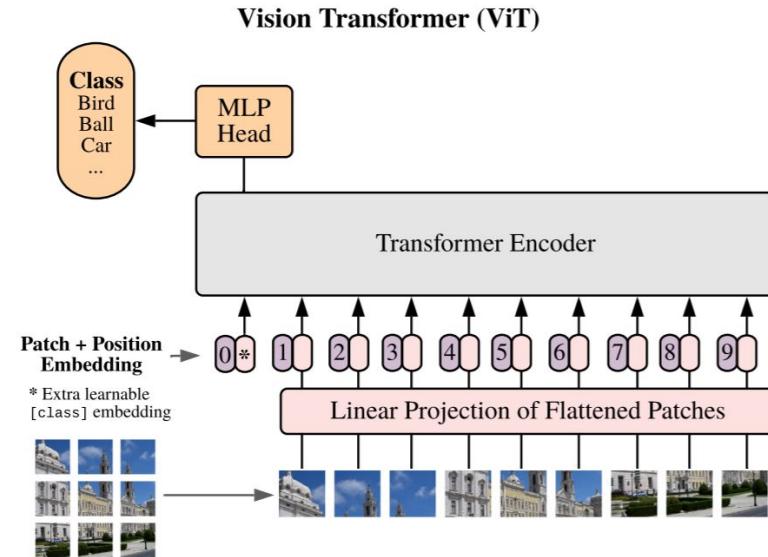
27 552 2K

Transformers

You will see the details of how transformers work in **Module 4** with Mirko.

In transformers, **vectors** are passed as **input**, and **vectors** are retrieved at the **output**. They can be **adapted** to many tasks (e.g. classification, object detection).

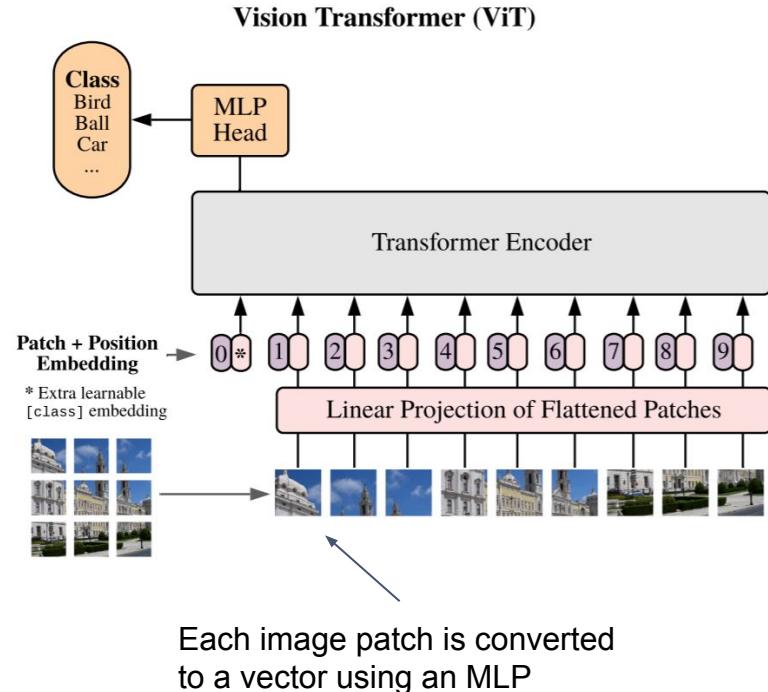
An **attention mechanism** is used which means all vectors can “attend” to all other vectors.



Vision Transformer

In Vision Transformer (ViT), it was shown that **transformers** could achieve **state-of-the-art** results on ImageNet.

This was achieved by **finetuning** their model on ImageNet after training on a much larger dataset.



Transformers for Computer Vision

Transformers have been **dominating** SOTA models in NLP (BERT, GPT, etc.)

We seem to be in an arms race for doing the same in **computer vision**.

Big tech companies have published **impressive results** using **transformers** on images.

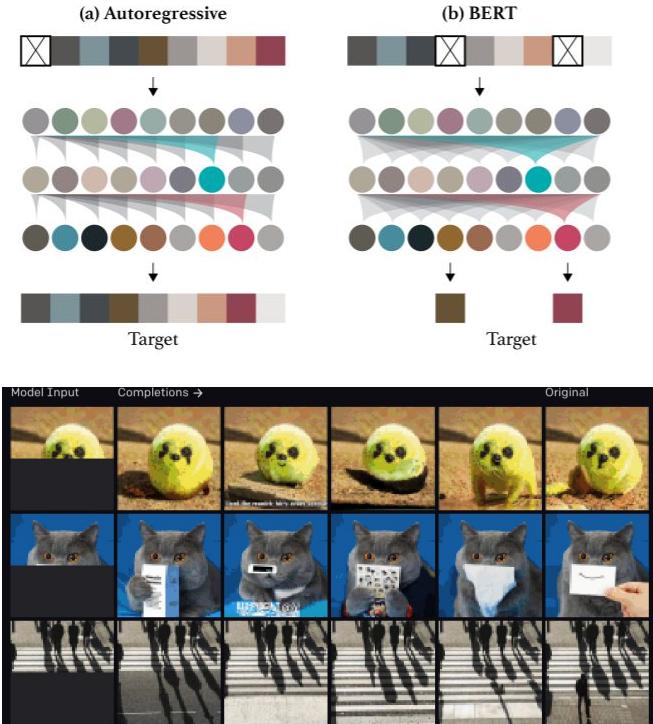


- igpt: Generative Pretraining from Pixels (ICML 2020)
- DALL-E and CLIP (2021)
- Axial DeepLab (ECCV 2020)
- Vision Transformers (ICLR 2021)
- Multi-modal Transformer for Video Retrieval (ECCV 2020)
- DETR: End-to-End Object Detection with Transformers (ECCV 2020)
- VL-BERT: Pre-training of Generic Visual-Linguistic Representations

imageGPT

In iGPT, the same transformer architectures as those used in NLP are applied to raw images and used to complete **masked** images.

They show that they can get interesting results using **self-supervised** learning.



DEtection TRansformer (DETR)

It has also been shown that transformers and CNNs can be **combined** to achieve very **competitive results** at object detection and image segmentation.

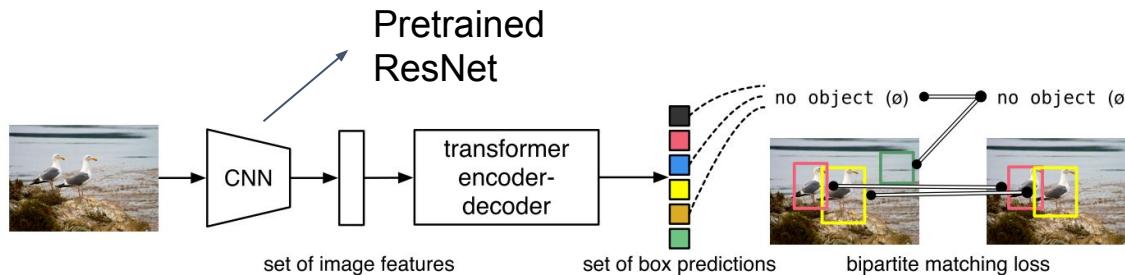
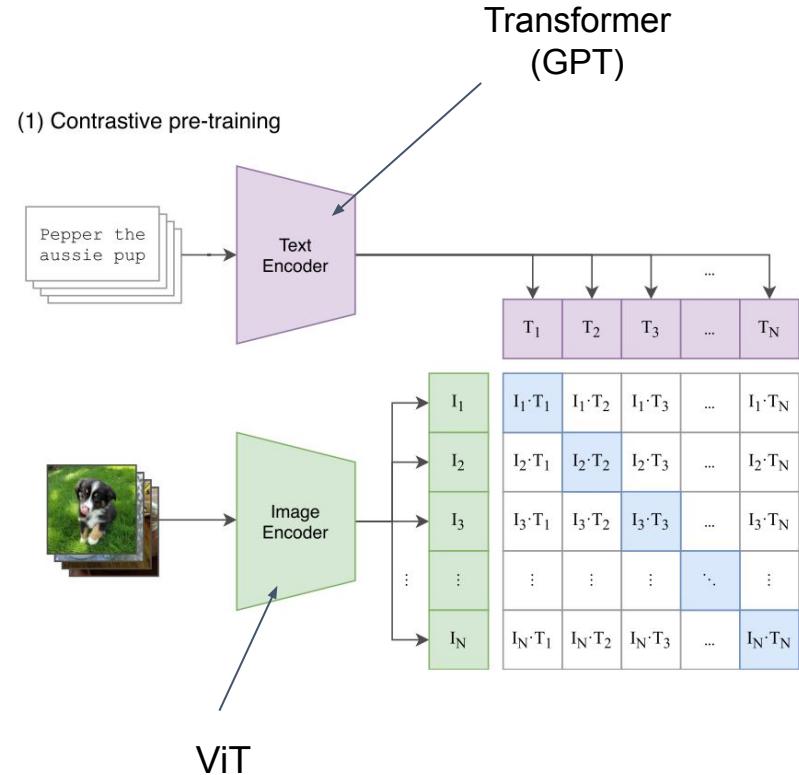


Fig. 1: DETR directly predicts (in parallel) the final set of detections by combining a common CNN with a transformer architecture. During training, bipartite matching uniquely assigns predictions with ground truth boxes. Prediction with no match should yield a “no object” (\emptyset) class prediction.

CLIP

- CLIP is a **multimodal** general-purpose transformer
- CLIP simply predicts if an image and text pair is related
- It is trained on 400 million pairs

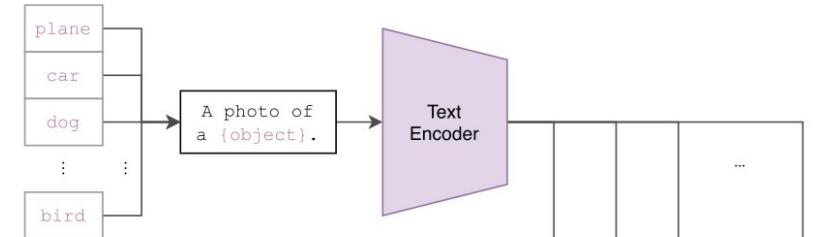


<https://openai.com/blog/clip/>
<https://arxiv.org/abs/2103.00020>

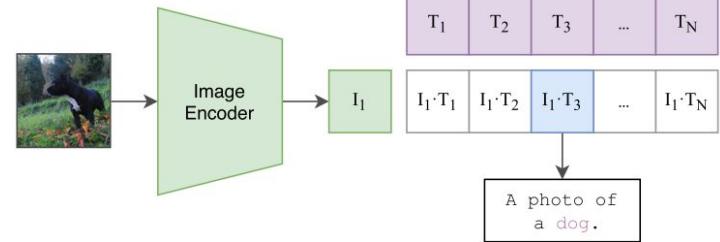
CLIP

- CLIP can be used in a **zero-shot** setting after training - it is never finetuned for the downstream task
- For example: CLIP can be “programmed” with natural language to classify images into new categories.

(2) Create dataset classifier from label text



(3) Use for zero-shot prediction

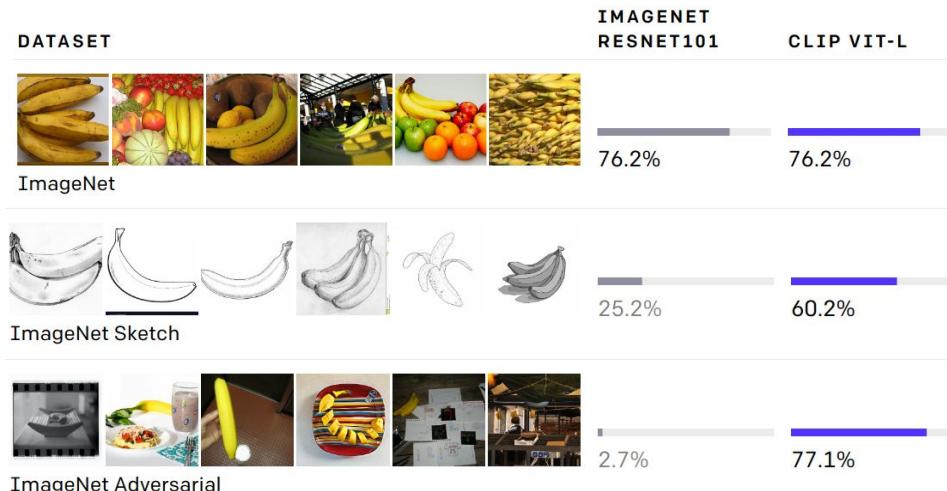


<https://openai.com/blog/clip/>

<https://arxiv.org/abs/2103.00020>

CLIP

- Without fine-tuning the network, CLIP demonstrates competitive results on ImageNet - and many other tasks!



<https://openai.com/blog/clip/>

<https://arxiv.org/abs/2103.00020>

CLIP

- In a follow up work, they demonstrate that “image” neurons respond to images that contain both images and text.

BIOLOGICAL NEURON	CLIP NEURON	PREVIOUS ARTIFICIAL NEURON				
Probed via depth electrodes	Neuron 244 from penultimate layer in CLIP RN50x4	Neuron 483, generic person detector from Inception v1				
Halle Berry	Spider-Man	human face				
	Responds to photos of Halle Berry and Halle Berry in costume ✓		Responds to photos of Spider-Man in costume and spiders ✓		Responds to photos of human faces ✓	Photorealistic images
	Responds to sketches of Halle Berry ✓		Responds to comics or drawings of Spider-Man and spider-themed icons ✓		Does not respond significantly to drawings of faces ✗	Conceptual drawings
	Responds to the text "Halle Berry" ✓		Responds to the text "spider" and others ✓		Does not respond significantly to text ✗	Images of text

<https://distill.pub/2021/multimodal-neurons/>

CLIP

Using **feature visualization**, we can explore images that maximally activate CLIP neurons.

Here we look at unit 89. Do you recognize a pattern?

Unit 89

FEATURE VISUALIZATION

An artificial, optimized image that maximizes activations of the given unit. [Read more.](#)



https://microscope.openai.com/models/contrastive_4x/image_block_4_5_Add_6_0/186

We can similarly explore images that were most associated with firing of this specific neuron.

Unit 89

FEATURE VISUALIZATION

An artificial, optimized image that maximizes activations of the given unit. [Read more](#).

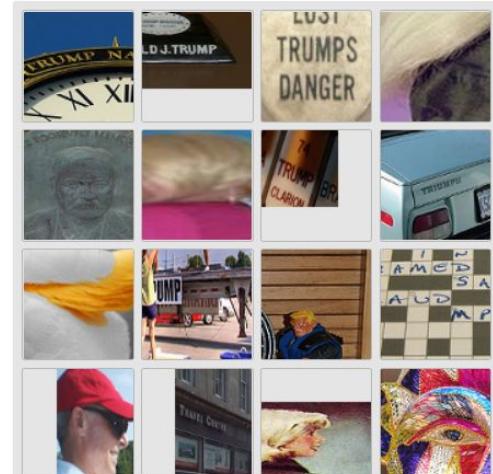


DATASET SAMPLES

Pieces of images from the training dataset that result in the largest activations from the given unit.

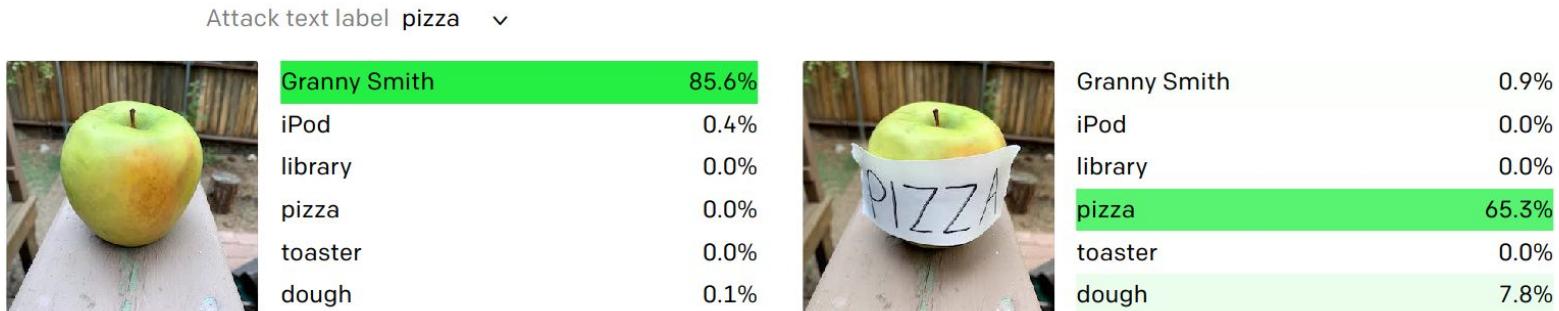
These images are cropped and downsized samples from the [ImageNet](#) research dataset. Unlike our other visualizations, they are not CC-BY-SA because they are derived from ImageNet.

DATASET: IMAGENET



CLIP

We can see that CLIP now becomes susceptible to a different kind of adversarial attack



DALL-E

A single transformer is used and combines text and images.

The network is trained to generate images given a textual prompt.

An autoencoder is used to reduce the dimensionality of input images.



+



TEXT PROMPT
an armchair in the shape of an avocado [...]

AI-GENERATED IMAGES

Edit prompt or view more images ↓

<https://openai.com/blog/dall-e/>
<https://arxiv.org/pdf/2102.12092.pdf>

Advantages

- Transformers offer a possibility to do self-supervised training on very large unlabelled sets of data.
- Transformers allow for global receptive fields.
- Transformers can potentially unify language and computer vision



Challenges

- **Curse of dimensionality:** Attention mechanisms don't scale well to large inputs (...for now).
- **Very large datasets and resources** are needed to achieve competitive results
- Attention mechanisms are **invariant to position**, which can mean loss of structural information.

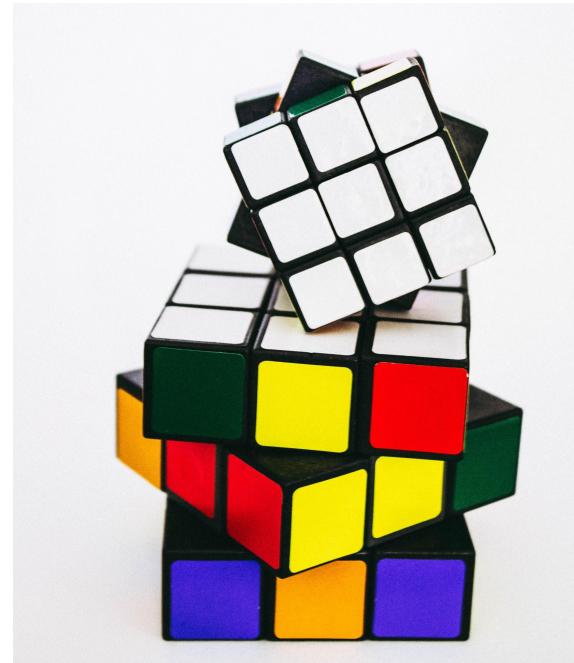
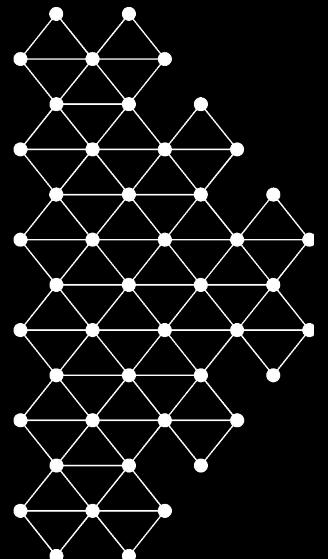


Photo by [Honey Yanibel Minaya Cruz](#) on [Unsplash](#)



Questions?