

Topics:

- Visualization

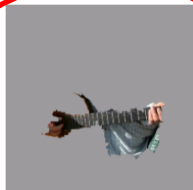
**CS 4644-DL / 7643-A**  
**ZSOLT KIRA**

# Visualization of Neural Networks

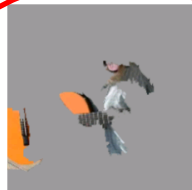
# Interpretability Enables Trust in AI Models

Understand the reasons behind a prediction

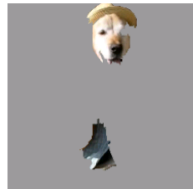
How did you  
make this  
prediction?



Electric guitar  
( $p=0.32$ )



Acoustic guitar  
( $p=0.24$ )



Labrador  
( $p=0.21$ )

Neural network  
for **image**  
**recognition**

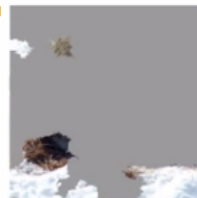
# Interpretability Enables Trust in AI Models

Figure out when **NOT** to trust a model

Prediction accuracy is very high. You are detecting snow, not wolves! It is time to put this system online. I can't trust you.



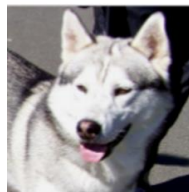
Predicted: **wolf**  
True: **wolf**



Predicted: **wolf**  
True: **wolf**



Predicted: **husky**  
True: **husky**



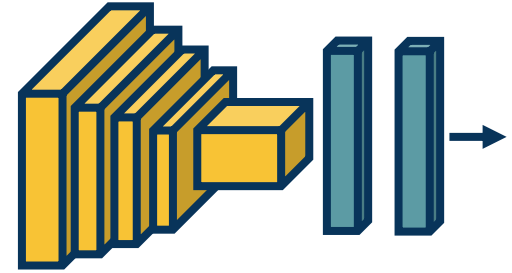
Predicted: **husky**  
True: **husky**

Neural network  
to predict  
**wolf vs husky**

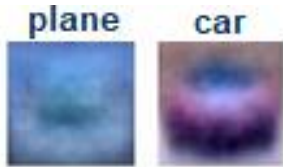


Predicted: **wolf**  
True: **husky**

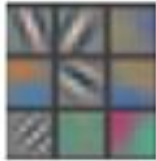
Given a **trained** model, we'd like to understand what it learned.



## Weights



*Fei-Fei Li, Justin Johnson, Serena Yeung, from CS 231n*

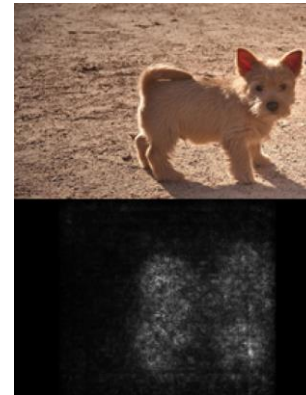


*Zeiler & Fergus, 2014*

## Activations



## Gradients



*Simonyan et al, 2013*

## Robustness

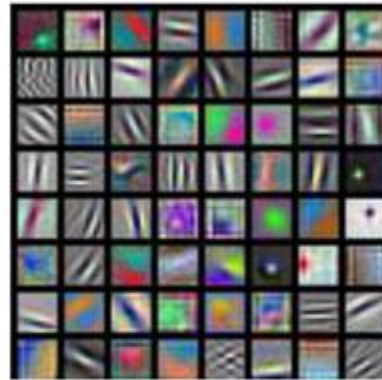


*Hendrycks & Dietterich, 2019*

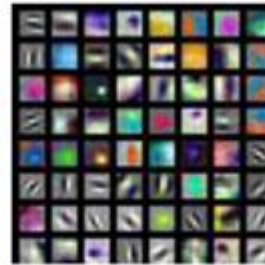
**FC Layer:** Reshape weights for a node back into size of image, scale 0-255



**Conv layers:**  
For each kernel,  
scale values  
from 0-255 and  
visualize



AlexNet:  
64 x 3 x 11 x 11



ResNet-18:  
64 x 3 x 7 x 7



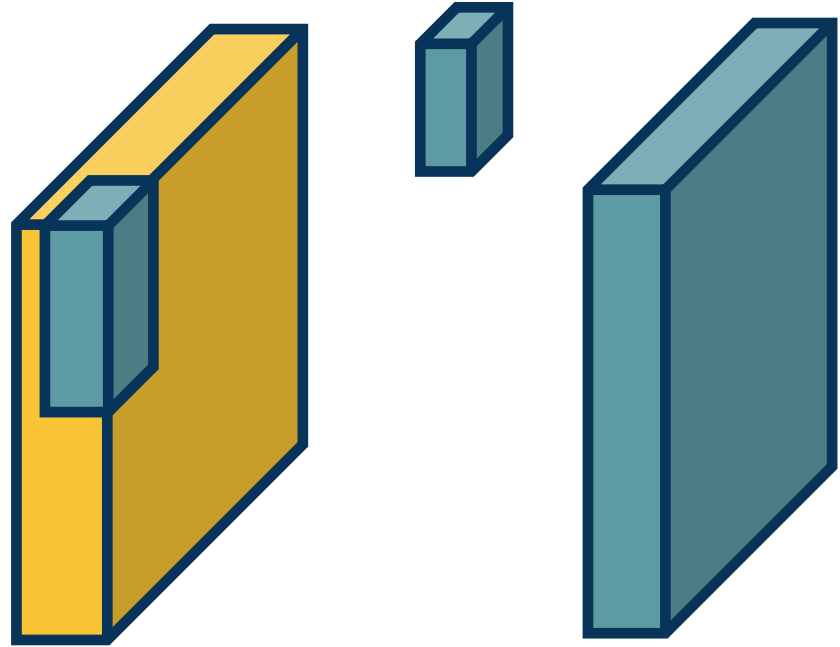
ResNet-101:  
64 x 3 x 7 x 7

**Problem:**  
3x3 filters  
difficult to  
interpret!

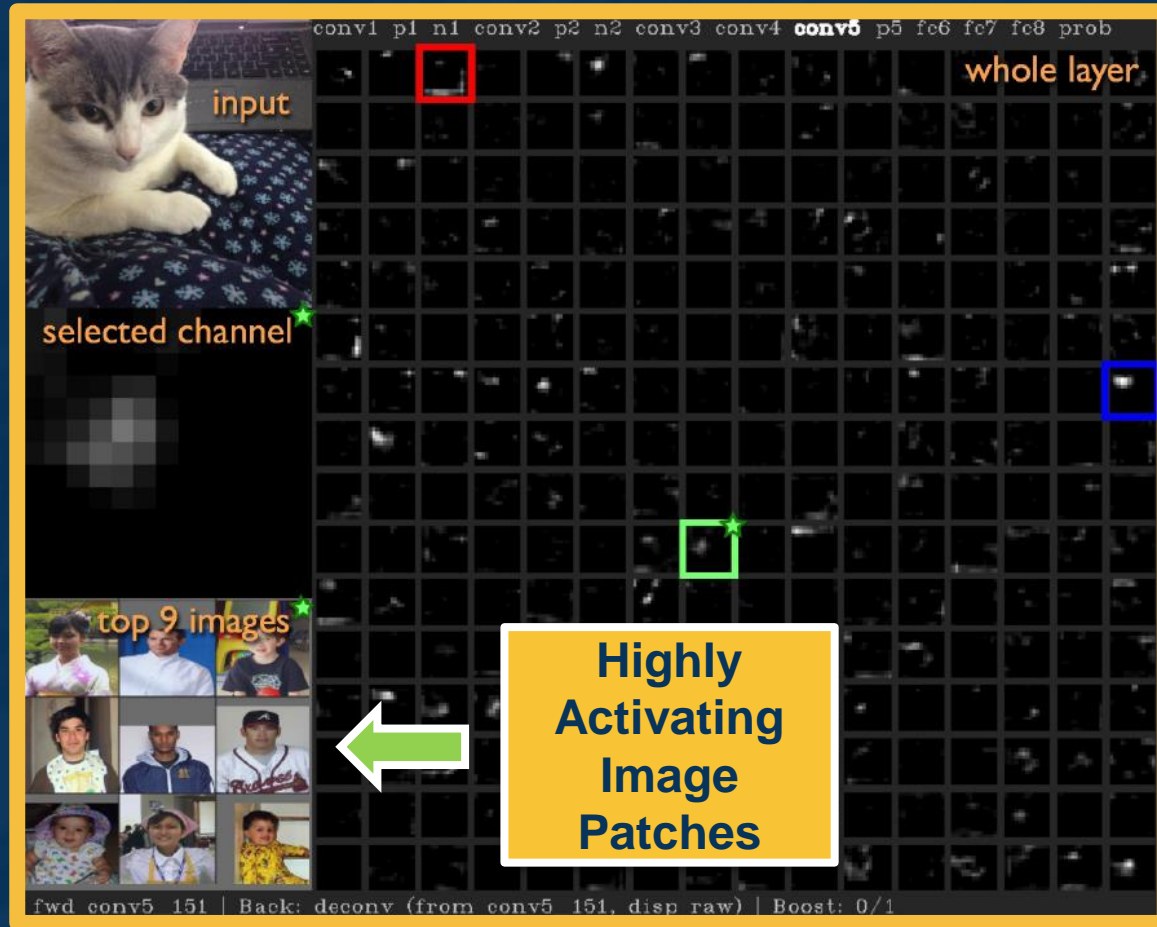
*Adapted from slides by Fei-Fei Li, Justin Johnson, Serena Young, from CS 231n*

We can also produce **visualization output (aka activation/filter) maps**

These are **larger** early in the network.



# Visualizing Output Maps



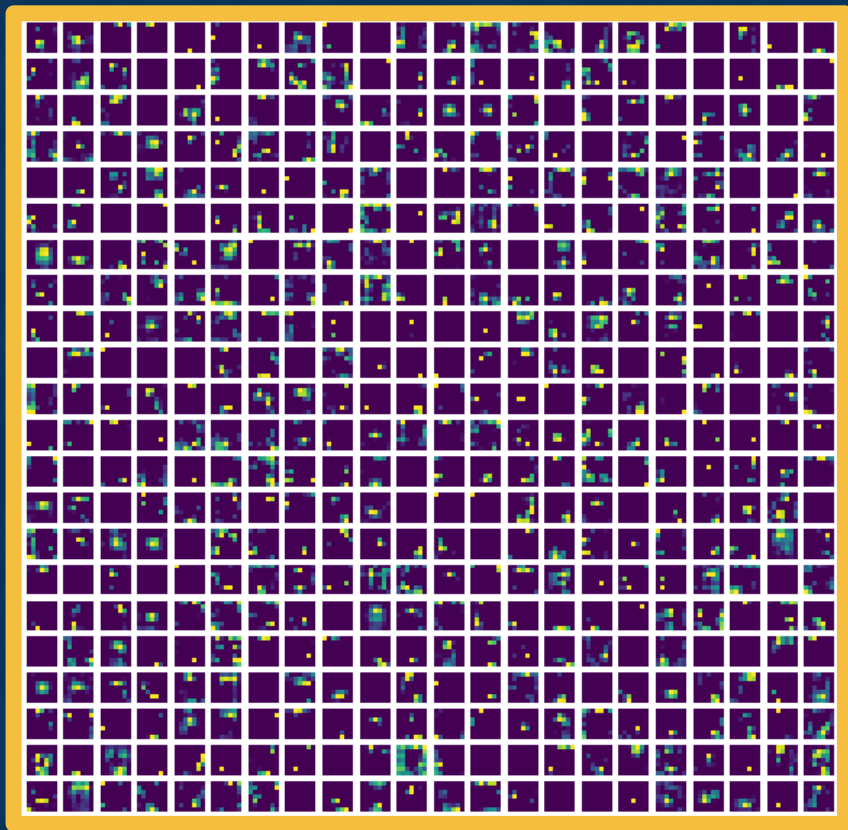
From: Yosinski et al., "Understanding Neural Networks Through Deep Visualization",

2015





# Activations – Small Output Sizes

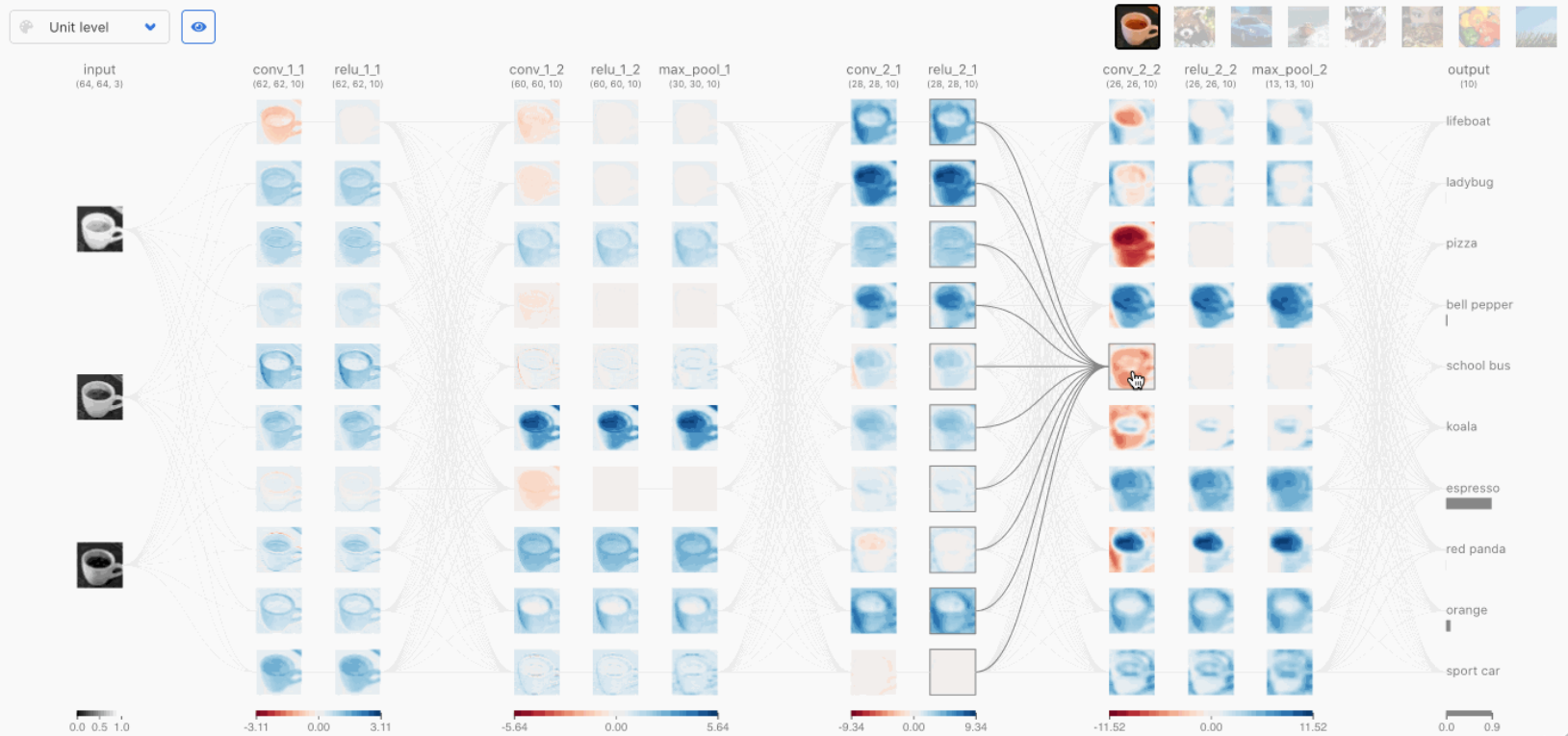


**Problem: Small conv outputs also hard to interpret**

**Activations of last conv layer in VGG network**

# CNN101 and CNN Explainer

**CNN 101** Learn Convolutional Neural Network (CNN) in your browser!



<https://poloclub.github.io/cnn-explainer/>

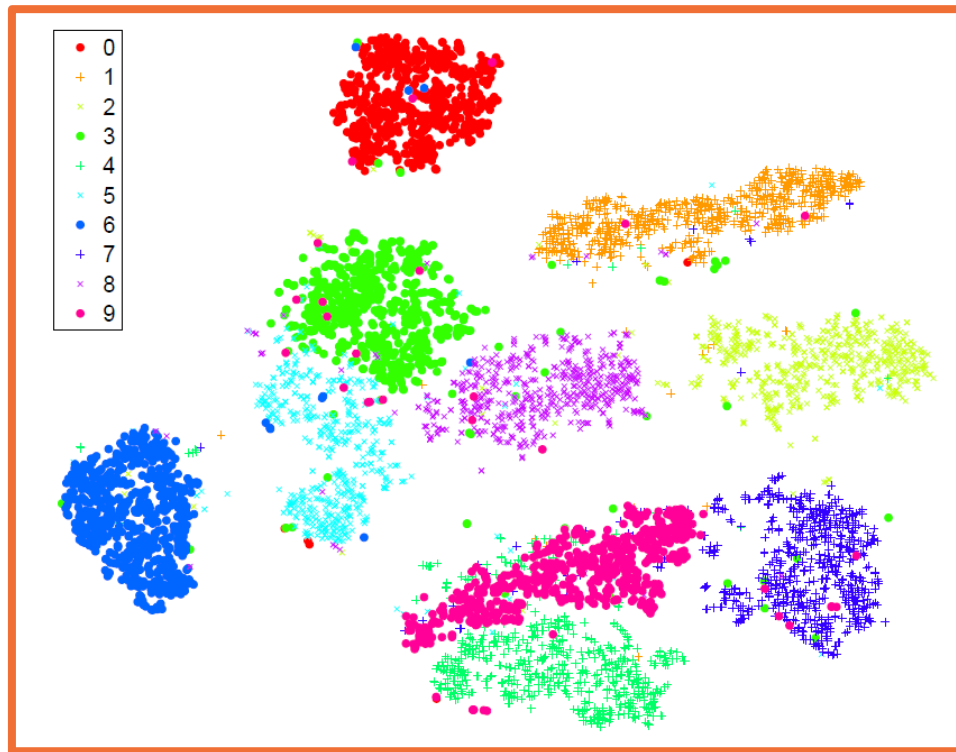
<https://fredhohman.com/papers/cnn101>

We can take the activations of any layer (FC, conv, etc.) and **perform dimensionality reduction**

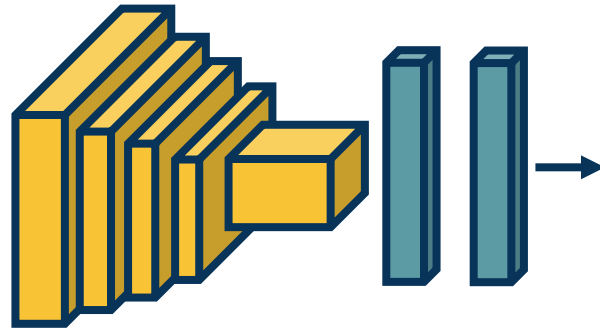
- Often reduce to two dimensions for plotting
- E.g. using Principle Component Analysis (PCA)

**t-SNE is most common**

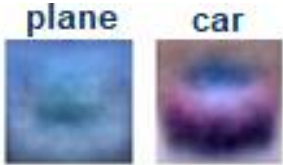
- Performs non-linear mapping to preserve pair-wise distances



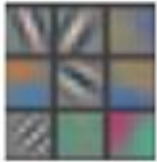
Van der Maaten & Hinton, "Visualizing Data using t-SNE", 2008.



## Weights



*Fei-Fei Li, Justin Johnson,  
Serena Yeung, from CS  
231n*



*Zeiler & Fergus, 2014*

## Activations



## Gradients



*Simonyan et al, 2013*

## Robustness



*Hendrycks & Dietterich,  
2019*

## Summary & Caveats

While these methods provide **some** visually interpretable representations, they can be misleading or uninformative (Adebayo et al., 2018)

Assessing interpretability is difficult

- Requires **user studies** to show **usefulness**
- E.g. they allow a user to predict mistakes beforehand

Neural networks learn **distributed representation**

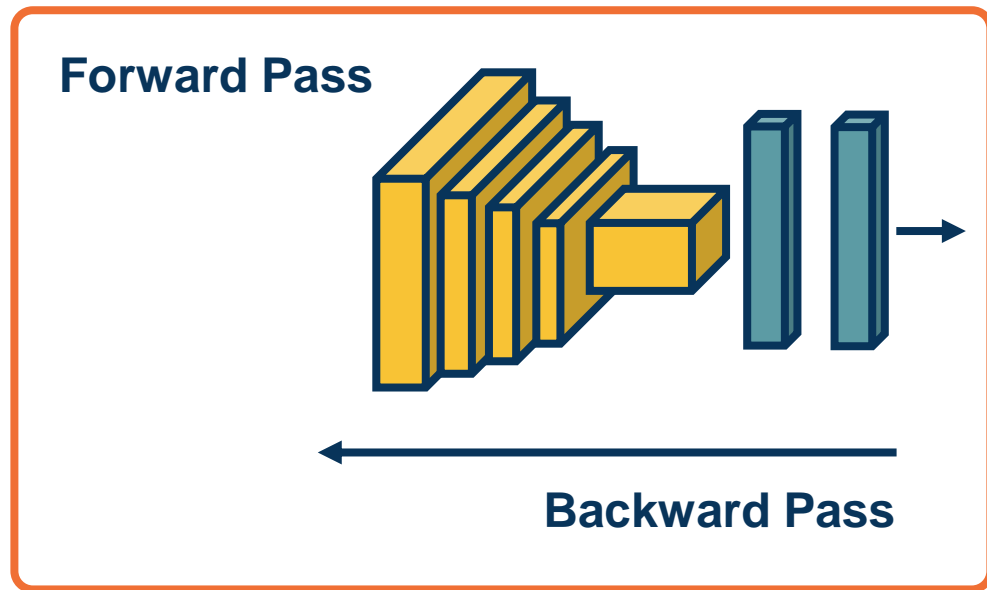
- (no one node represents a particular feature)
- This makes interpretation difficult

*Adebayo et al., "Sanity Checks for Saliency Maps", 2018.*



# Gradient- Based Visualizations

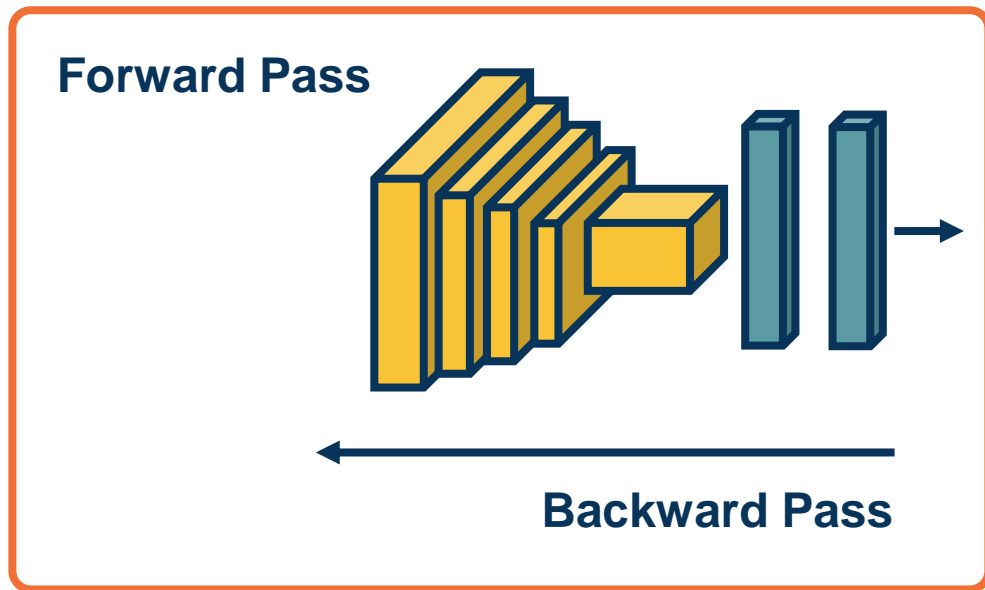
Given a **trained** model, we can perform forward pass given an input to get scores, softmax probabilities, loss and then backwards pass to get gradients



- ◆ Note: We are keeping parameters/weights **frozen**
- ◆ Do not use gradients w.r.t. weights to perform updates

Backwards pass gives us **gradients** for all layers: How the loss changes as we change different parts of the input

This can be **useful not just for optimization**, but also to understand what was learned



- ◆ Gradient of **loss** with respect to **all layers** (including input!)
- ◆ Gradient of **any layer** with respect to **input** (by cutting off computation graph)

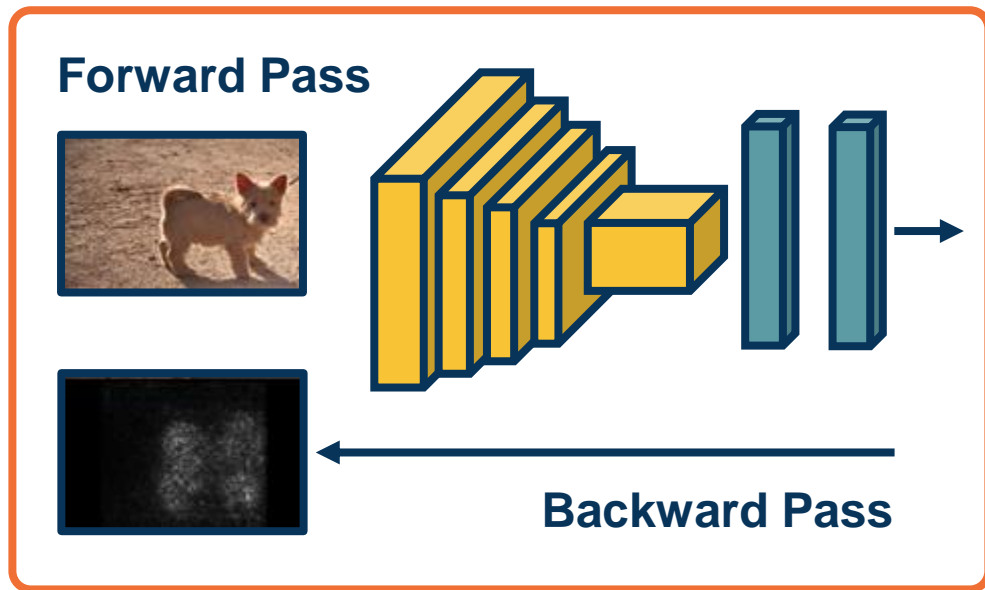


**Idea:** We can backprop to the image

- ◆ Sensitivity of loss to individual pixel changes
- ◆ Large sensitivity implies important pixels
- ◆ Called **Saliency Maps**

**In practice:**

- ◆ Instead of loss, find gradient of classifier **scores** (pre-softmax)
- ◆ Take absolute value of gradient
- ◆ Sum across all channels



*From: Simonyan et al., "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", 2013*

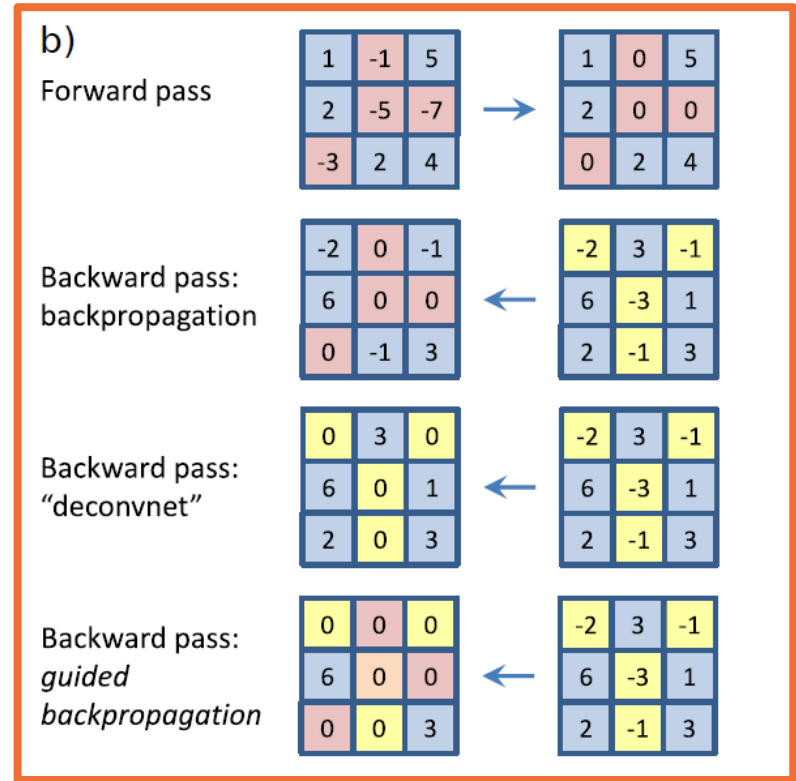
**Gradient of Loss w.r.t. Image**

Normal backprop not always best choice

**Example:** You may get parts of image that **decrease** the feature activation

- There are probably lots of such input pixels

**Guided backprop** can be used to improve visualizations



From: Springenberg et al., "Striving For Simplicity: The All Convolutional Net"

# Guided Backprop Results

guided backpropagation



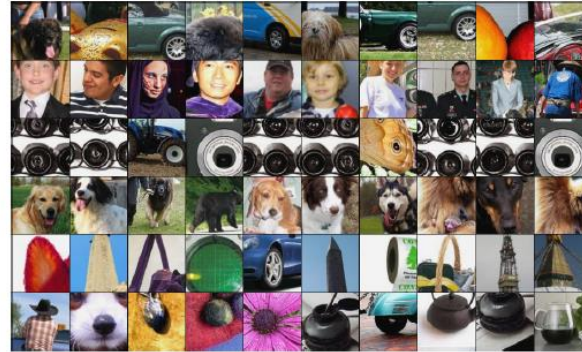
corresponding image crops



guided backpropagation



corresponding image crops



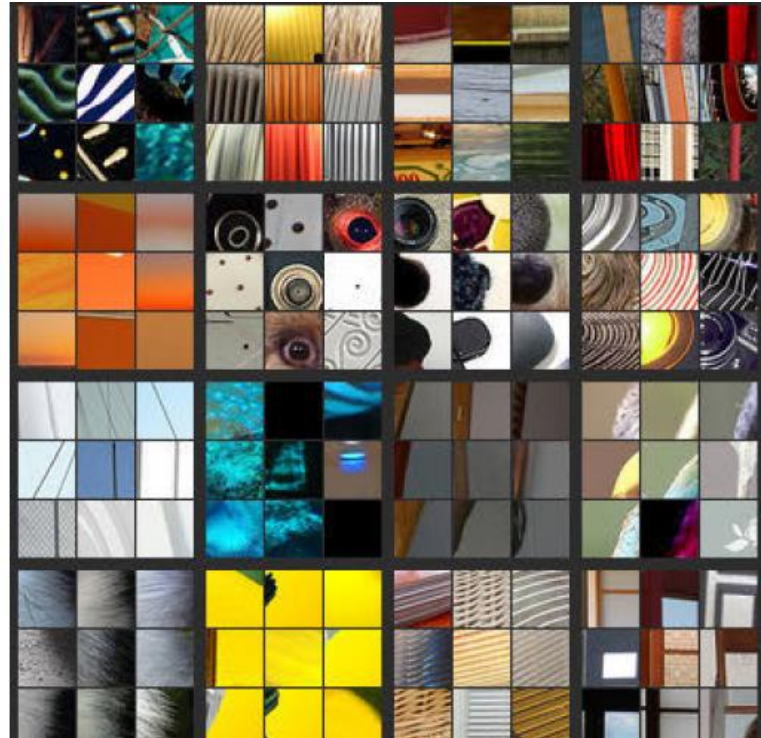
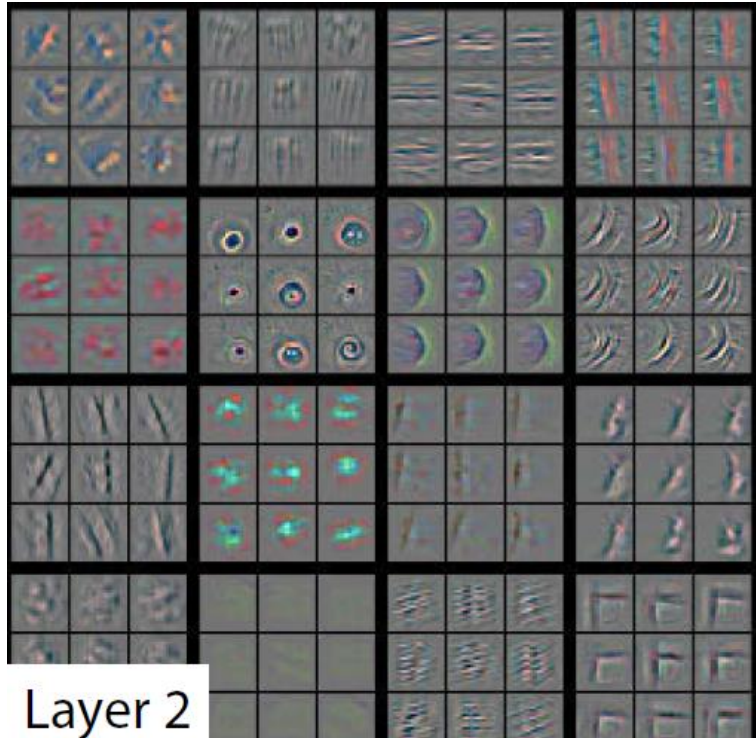
From: Springenberg et al., "Striving For Simplicity: The All Convolutional Net"

# VGG Layer-by-Layer Visualization



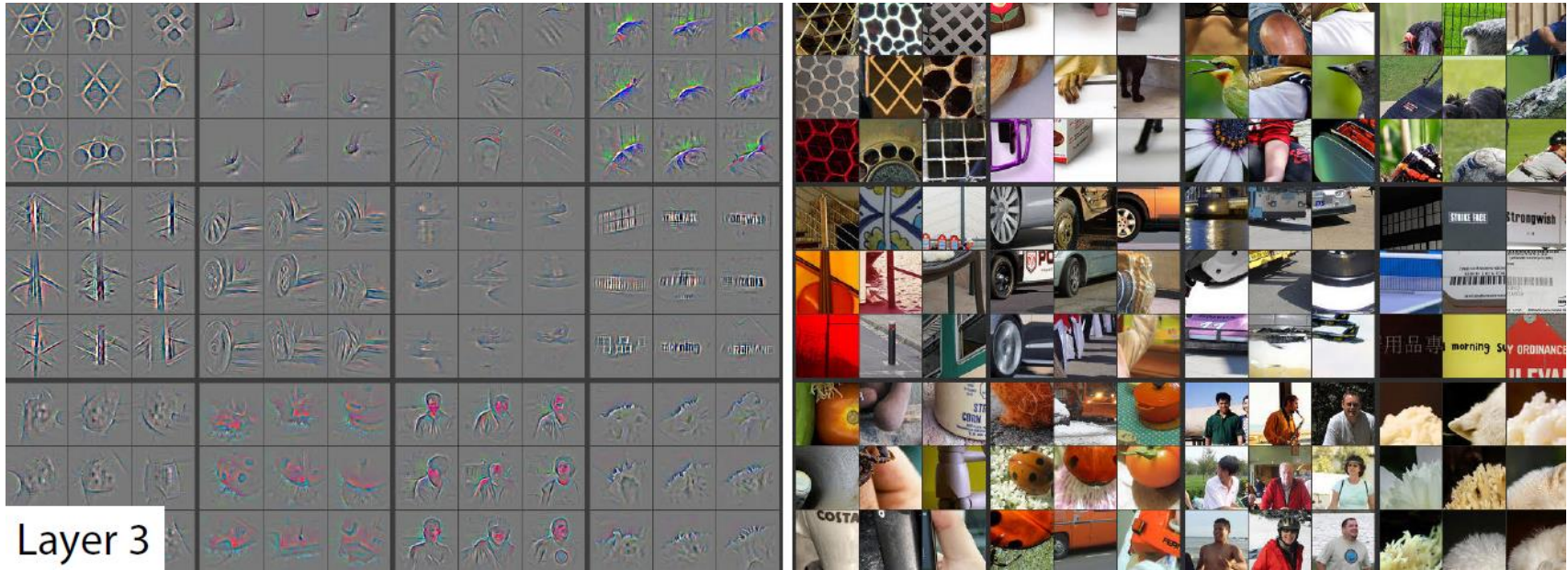
**Note:** These images were created by a slightly different method called **deconvolution**, which ends up being similar to guided backprop

# VGG Layer-by-Layer Visualization



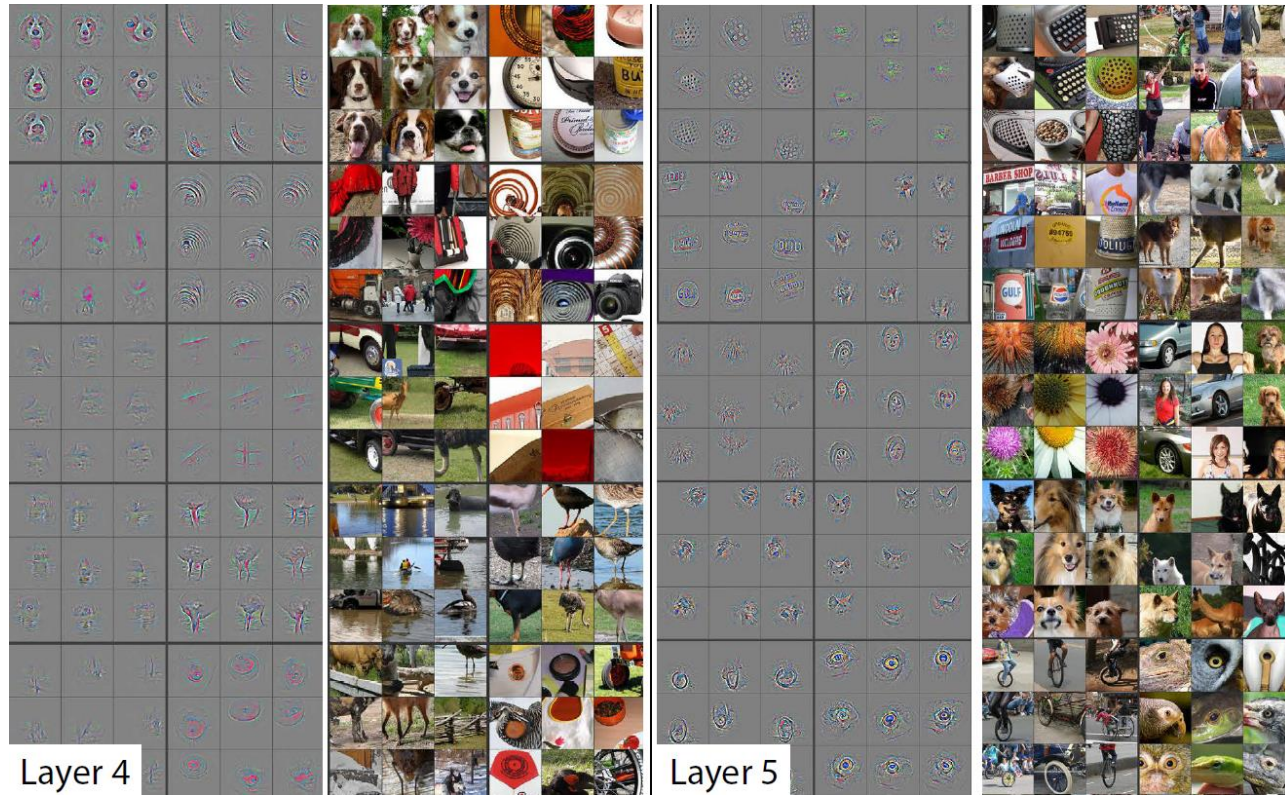
From: "Visualizing and Understanding Convolutional Networks, Zeiler & Fergus, 2014.

# VGG Layer-by-Layer Visualization



From: "Visualizing and Understanding Convolutional Networks, Zeiler & Fergus, 2014.

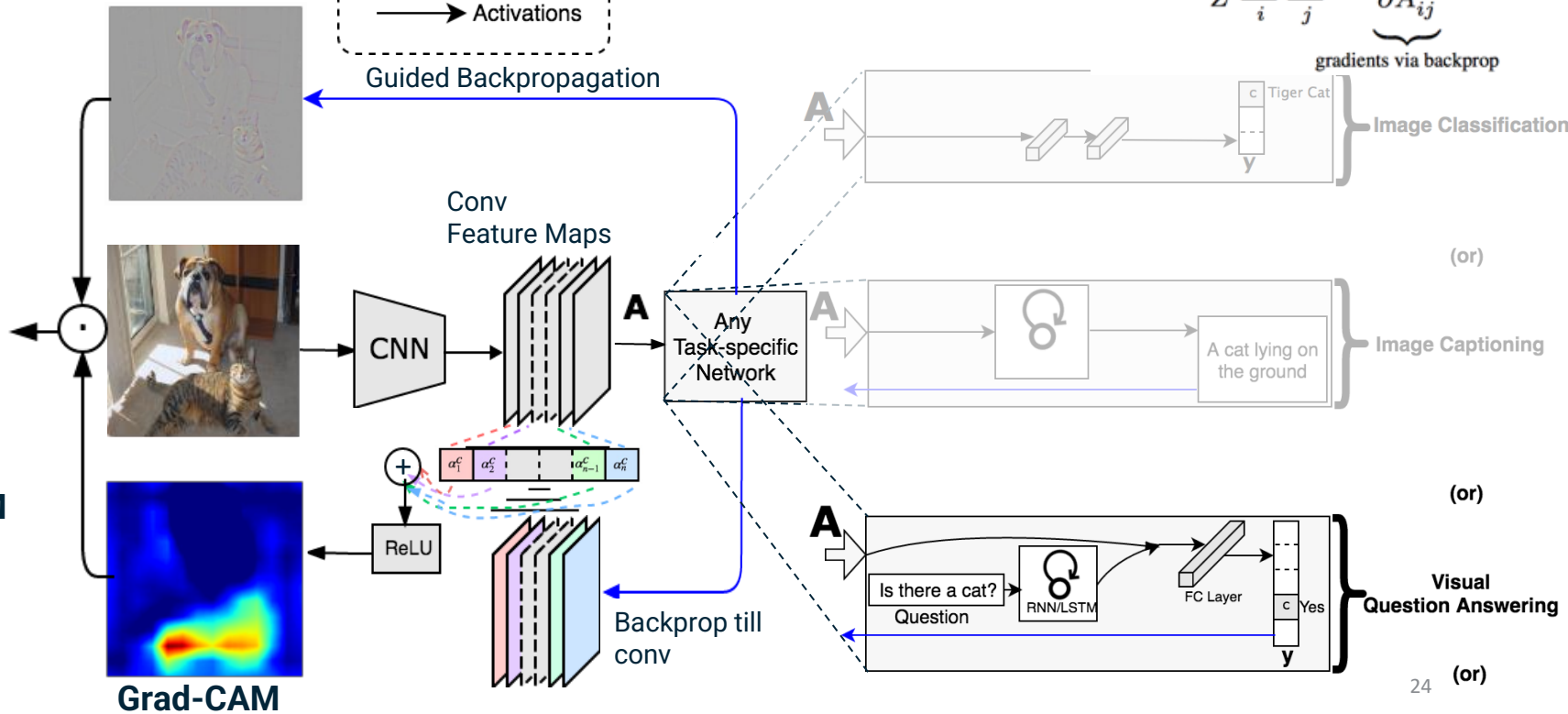
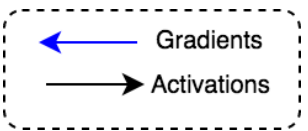
# VGG Layer-by-Layer Visualization



From: "Visualizing and Understanding Convolutional Networks, Zeiler & Fergus, 2014.

# Neuron Importance

$$\alpha_k^c = \frac{1}{Z} \sum_i \sum_j \overbrace{\frac{\partial y^c}{\partial A_{ij}^k}}^{\text{global average pooling}}$$

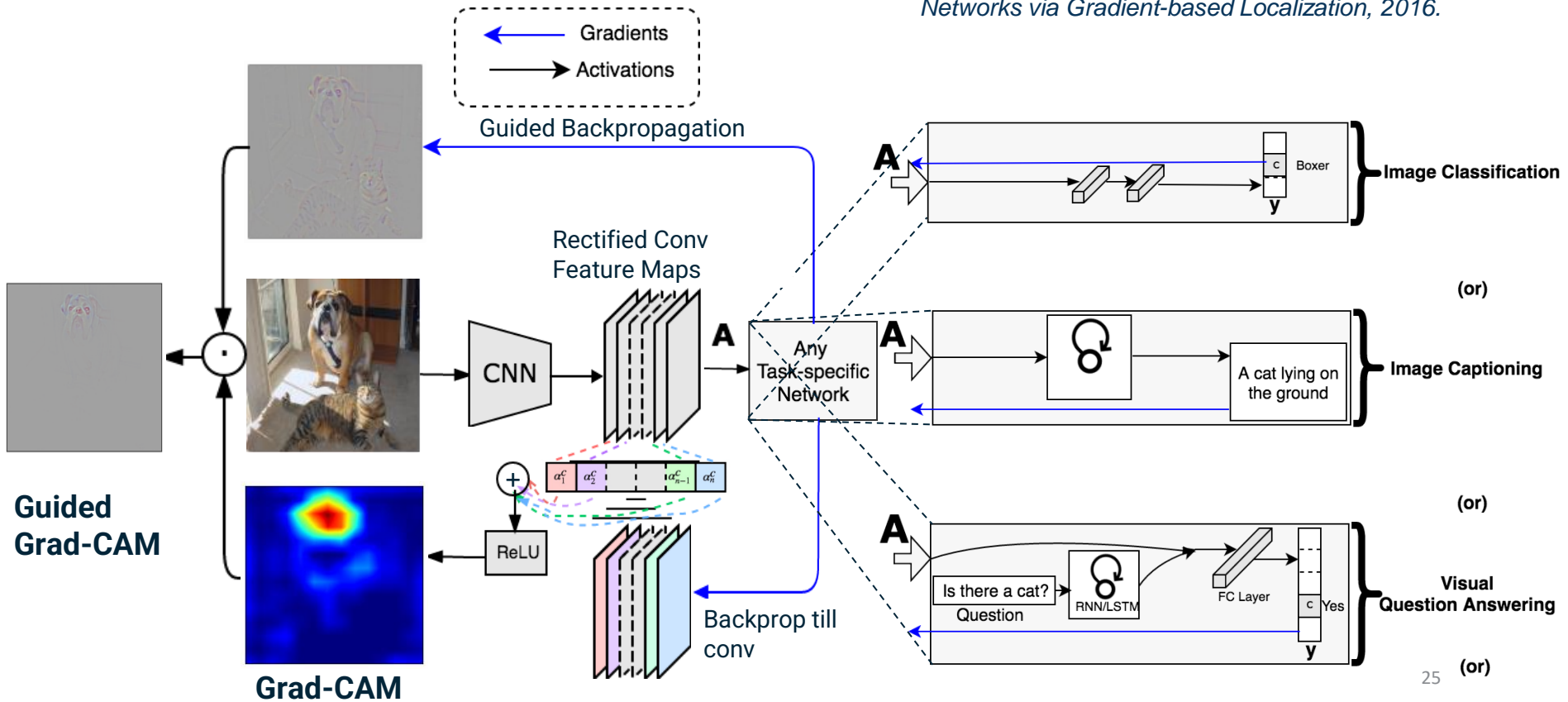


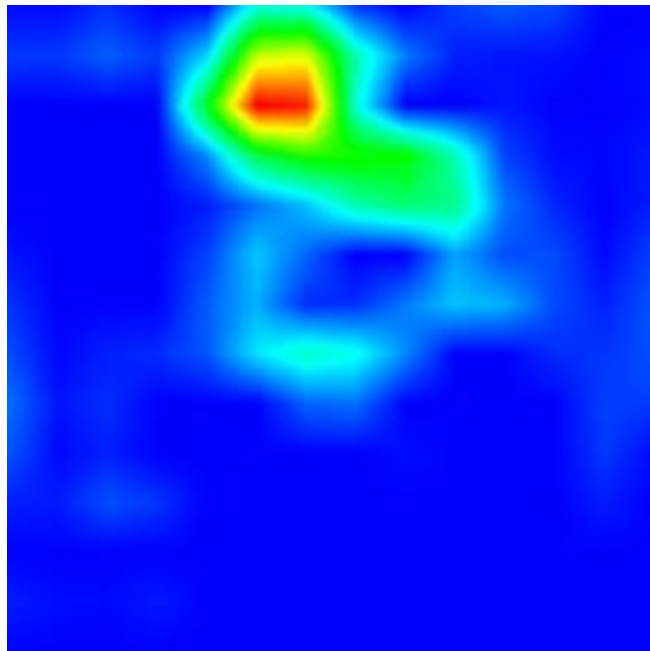
**Guided Grad-CAM**

**Grad-CAM**

*Selvaraju et al., Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization, 2016.*

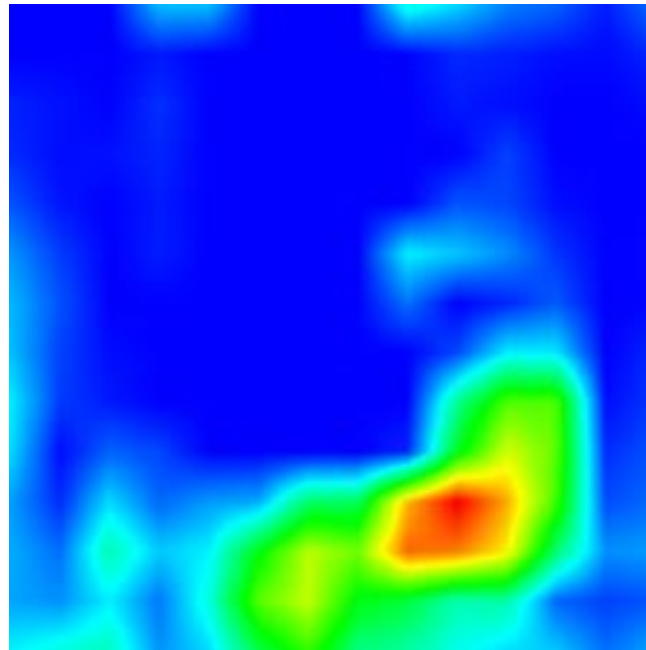






What animal is in this picture? Dog

*Selfvaraju et al., Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization, 2016.*



What animal is in this picture? Cat

*Selfvaraju et al., Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization, 2016.*

## Summary

- ◆ Gradients are important **not just for optimization**, but also for **analyzing** what neural networks **have learned**
- ◆ Standard backprop **not always the most informative** for visualization purposes
- ◆ Several ways to **modify the gradient flow** to improve visualization results

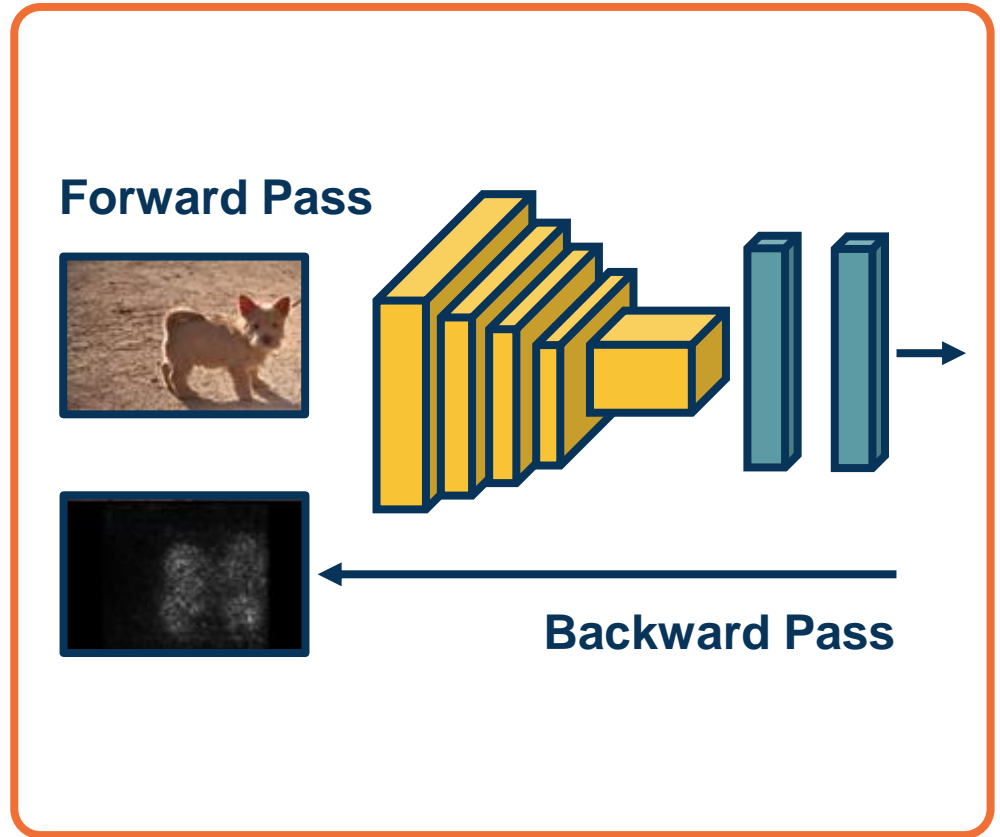


# Optimizing the Input Images

**Idea:** Since we have the gradient of scores w.r.t. inputs, can we *optimize* the image itself to maximize the score?

**Why?**

- Generate images from scratch!
- Adversarial examples



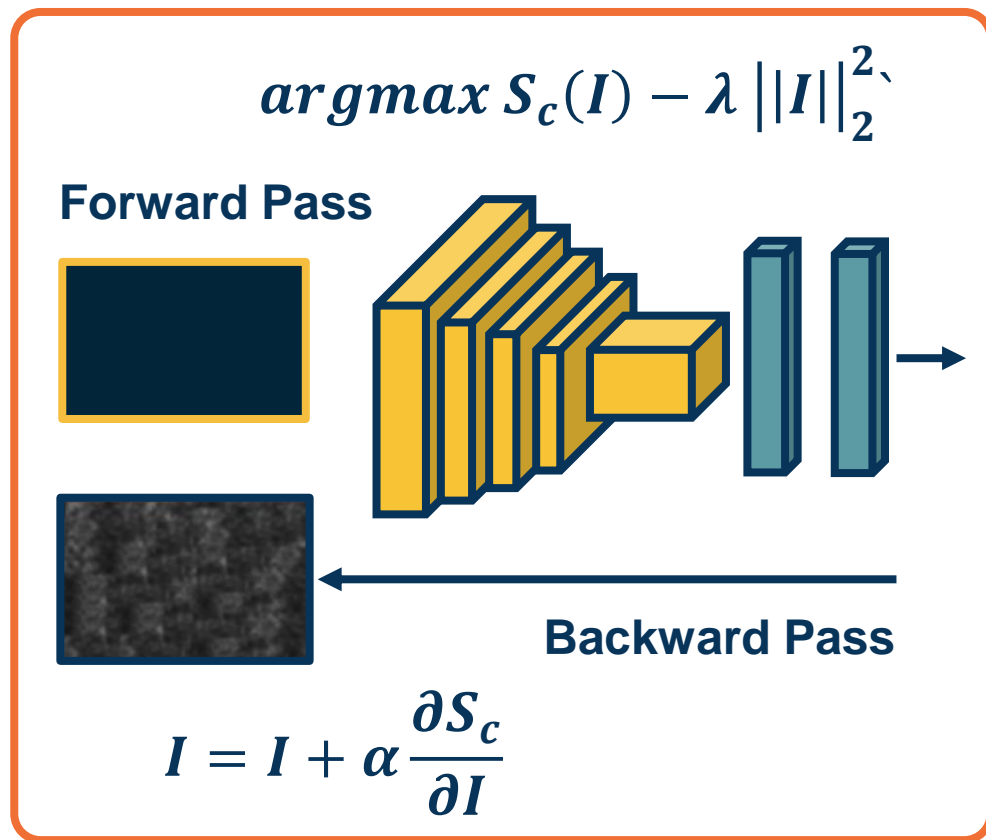
From: Simonyan et al., "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", 2013

We can perform **gradient ascent** on image

- Start from random/zero image
- Use scores to avoid minimizing other class scores instead

Often need **regularization term** to induce statistics of natural imagery

- E.g. small pixel values, spatial smoothness



From: Simonyan et al., "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", 2013

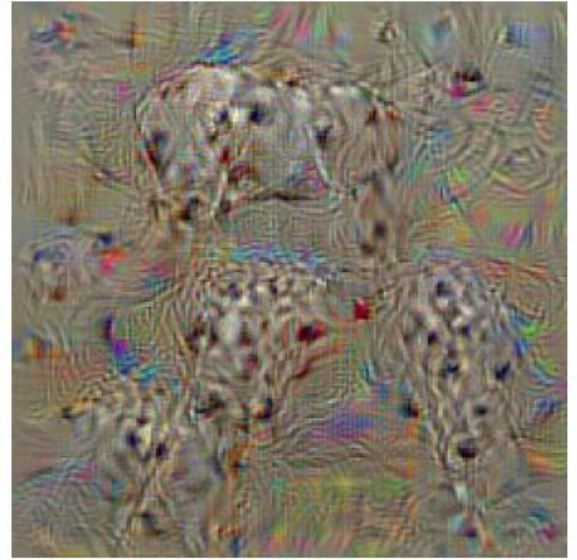
# Example Images



**dumbbell**



**cup**



**dalmatian**

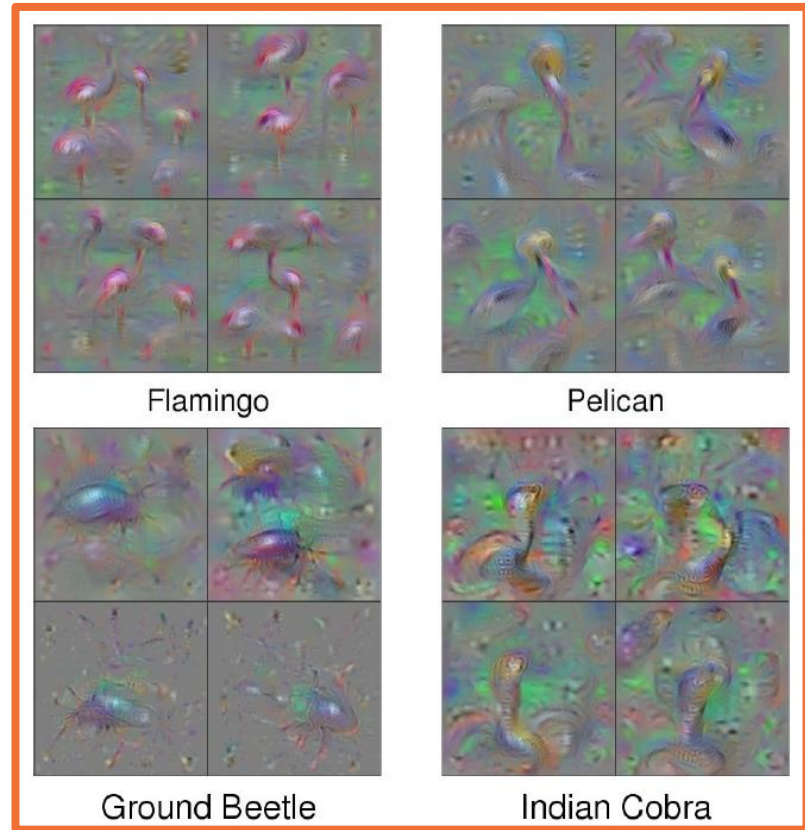
**Note: You might have to squint!**

*From: Simonyan et al., "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", 2015*



Can improve results with **various tricks:**

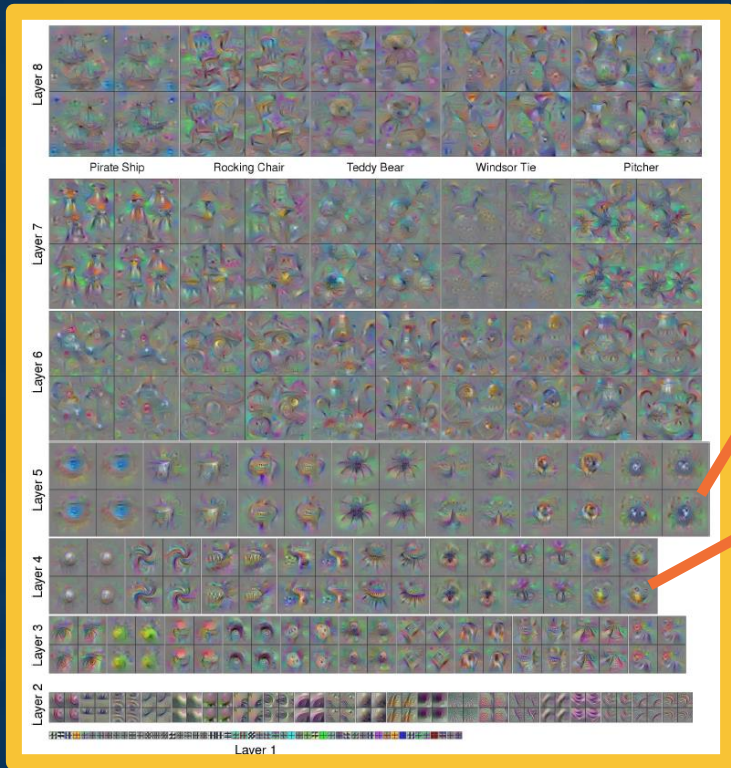
- Clipping or normalization of small values & gradients
- Gaussian blurring



From: Yosinski et al., "Understanding Neural Networks Through Deep Visualization", 2015

# Improved Results

Note: Can generate input images to maximize any arbitrary activation!



## Summary

We can optimize the input image to **generate** examples to increase class scores or activations

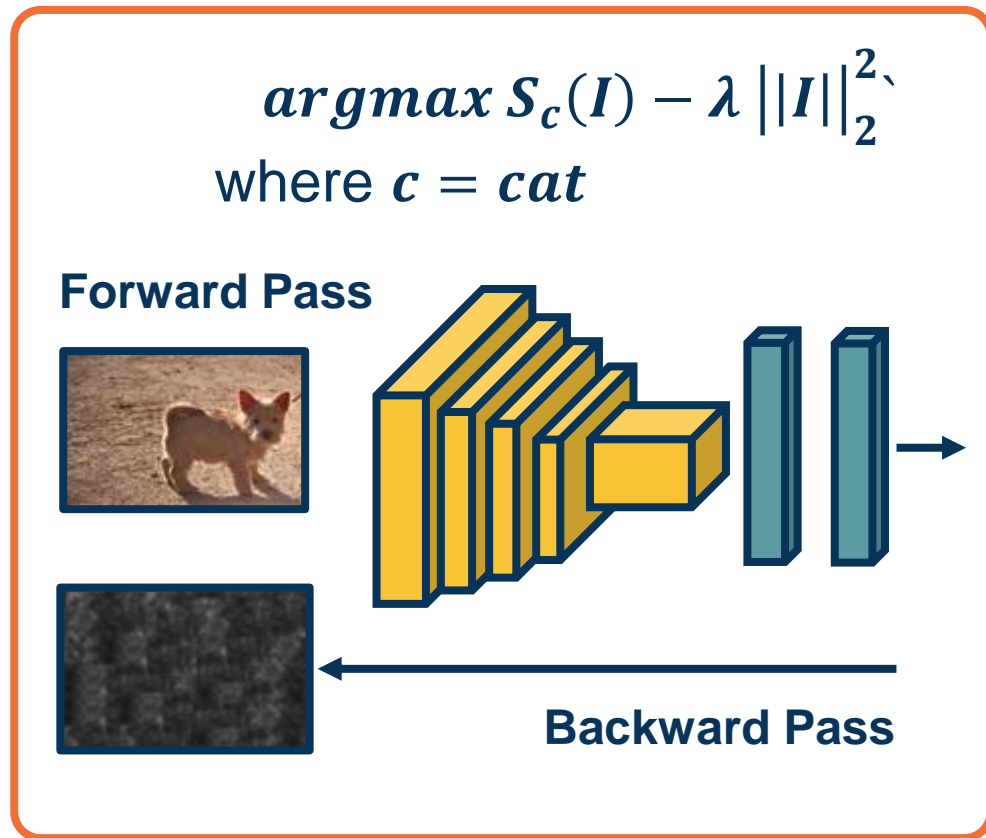
This can show us a great deal about what examples (not in the training set) **activate the network**






# Testing Robustness

- ◆ We can perform **gradient ascent** on image
- ◆ Rather than start from zero image, why not real image?
- ◆ And why not optimize the score of an **arbitrary** (incorrect!) class

**Surprising result: You need very small amount of pixel changes to make the network confidently wrong!**



From: Simonyan et al., "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", 2013

	$+ .007 \times$		$=$	
$x$		$\text{sign}(\nabla_x J(\theta, x, y))$		$x + \epsilon \text{sign}(\nabla_x J(\theta, x, y))$
“panda”		“nematode”		“gibbon”
57.7% confidence		8.2% confidence		99.3 % confidence

## Note this problem is not specific to deep learning!

- ◆ Other methods also suffer from it
- ◆ Can show how **linearity** (even at the end) can bring this about
  - ◆ Can add many small values that add up in right direction

*From: Goodfellow et al., “Explaining and Harnessing Adversarial Examples”, 2015*

## Example of Adversarial Noise

# Variations of Attacks

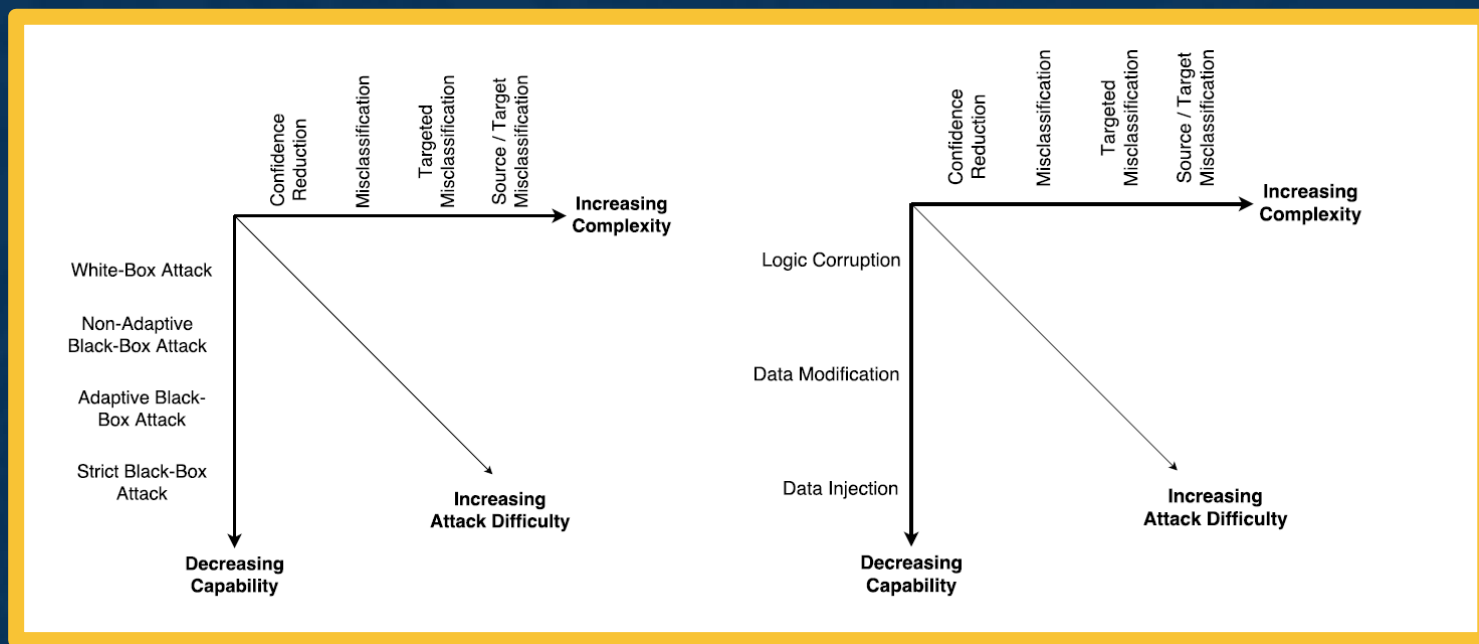
**VGG**



**DEER**  
**AIRPLANE(85.3%)**



**BIRD**  
**FROG(86.5%)**



## Single-Pixel Attacks!

## White vs. Black-Box Attacks of Increasing Complexity

Chakraborty et al., *Adversarial Attacks and Defences: A Survey*, 2018

Su et al., "One Pixel Attack for Fooling Deep Neural Networks", 2019.

## Summary of Adversarial Attacks/Defenses

Similar to other security-related areas, it's an active **cat-and-mouse game**

**Several defenses such as:**

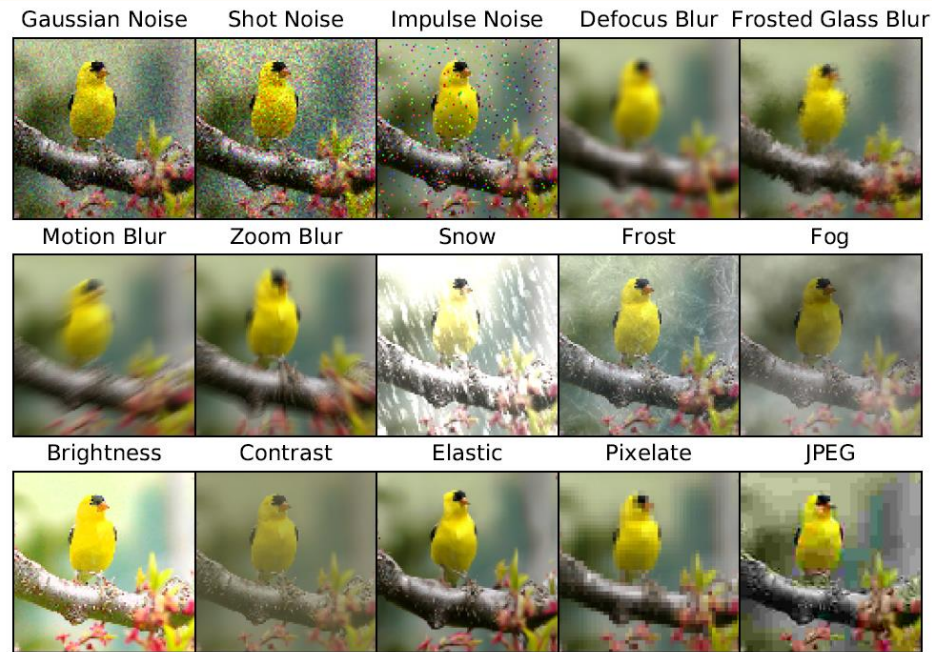
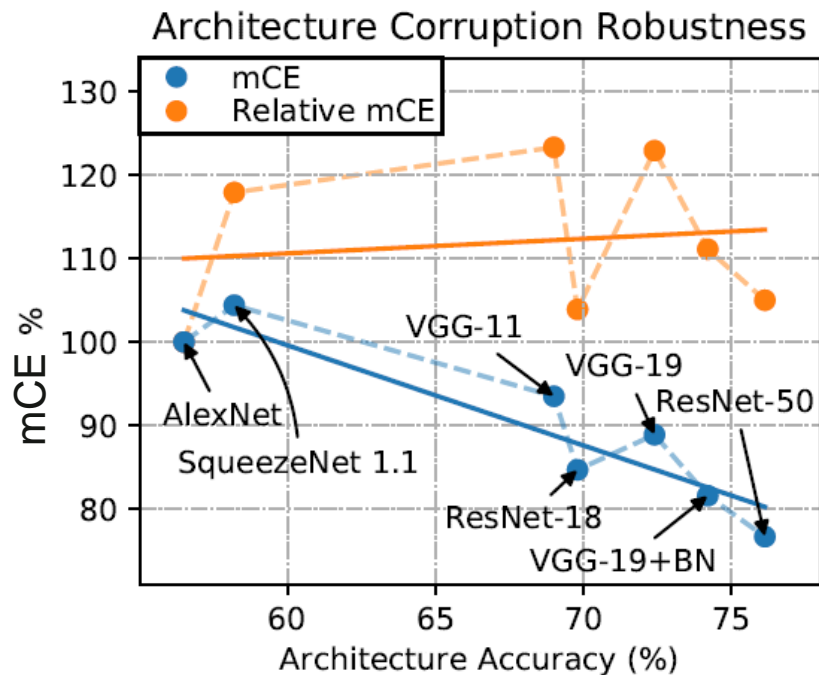
- Training with adversarial examples
- Perturbations, noise, or re-encoding of inputs

There are **not universal methods** that are robust to all types of attacks





# Other Forms of Robustness Testing



$$CE_c^f = \left( \sum_{s=1}^5 E_{s,c}^f \right) / \left( \sum_{s=1}^5 E_{s,c}^{\text{AlexNet}} \right).$$

Hendrycks & Dietterich, "Benchmarking Neural Network Robustness to Common Corruptions and Perturbations" 2019.

# We can try to understand the **biases of CNNs**

- ◆ Can compare to those of humans

## Example: **Shape vs. Texture Bias**

*Geirhos, "ImageNet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness", 2018.*



(a) Texture image

81.4%	<b>Indian elephant</b>
10.3%	indri
8.2%	black swan

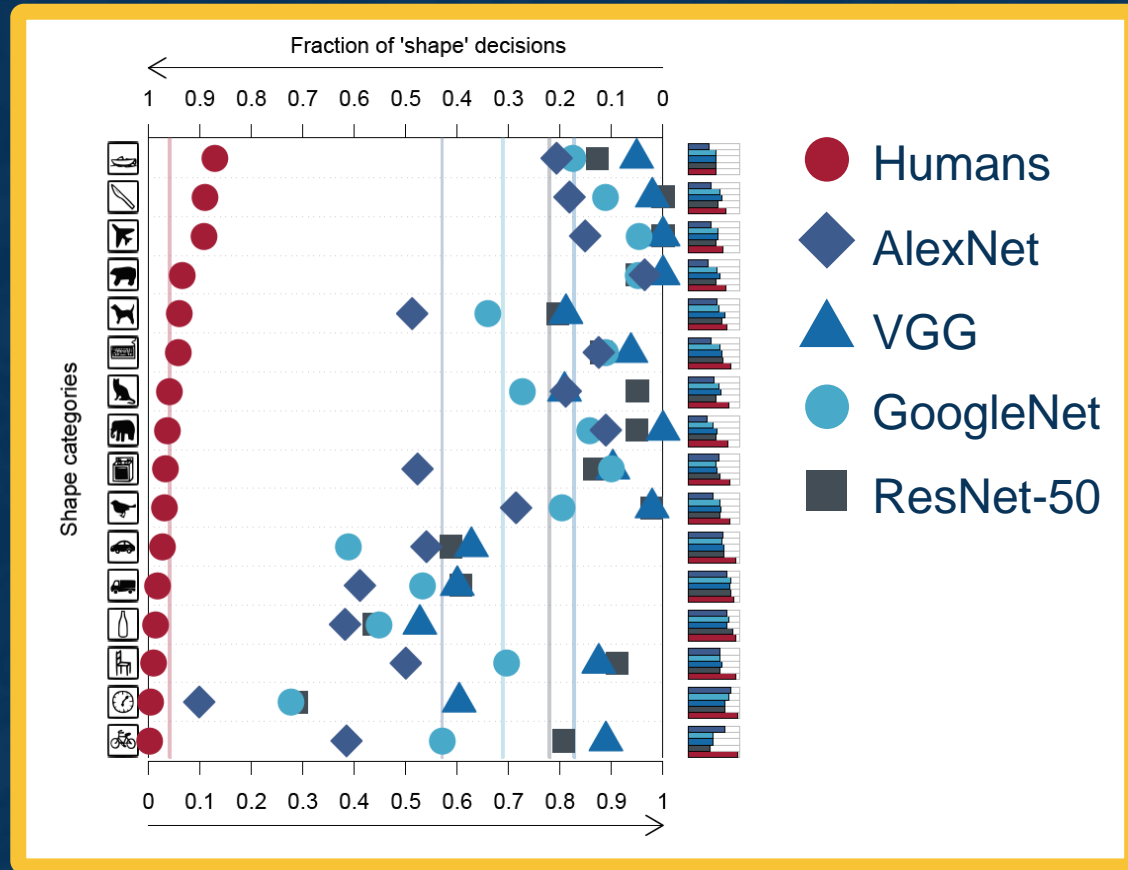
(b) Content image

71.1%	<b>tabby cat</b>
17.3%	grey fox
3.3%	Siamese cat

(c) Texture-shape cue conflict

63.9%	<b>Indian elephant</b>
26.4%	indri
9.6%	black swan

# Shape vs. Texture Bias



Geirhos, "ImageNet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness", 2018.

## Summary

- ◆ Various ways to test the **robustness** and **biases** of neural networks
- ◆ Adversarial examples have **implications** for understanding and trusting them
- ◆ Exploring the **gain of different architectures** in terms of robustness and biases can also be used to understand what has been learned



**What about  
Transformers,  
LLMs, etc.**

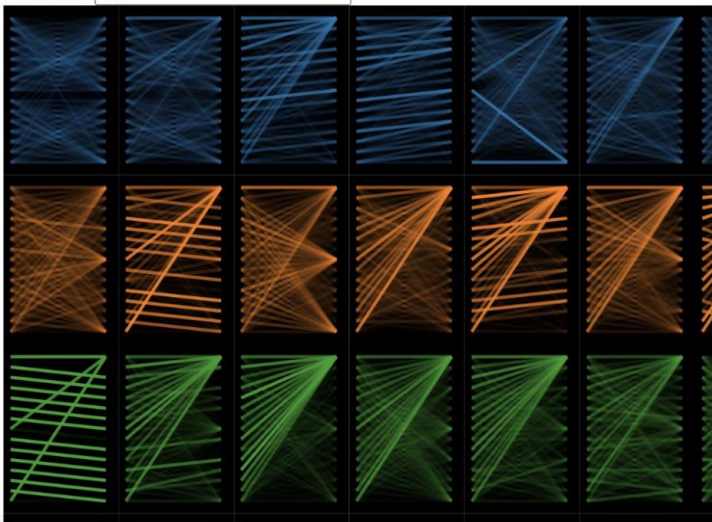
- Large models, especially transformers, can be seen as implementing **algorithms**
- Several ways to try to understand:
  - Visualization – Often attention
  - Distill into more interpretable model
  - Reverse engineer
  - Forward engineer: Algorithm → compiler → Weights!

## Usage

- Click on any cell for a detailed view of attention for the associated attention head
- Then hover over any token on the left side of detail view to filter the attention
- The lines show the attention from each token (left) to every other token

```
model_view(attention, tokens, sentence_b_start)
```

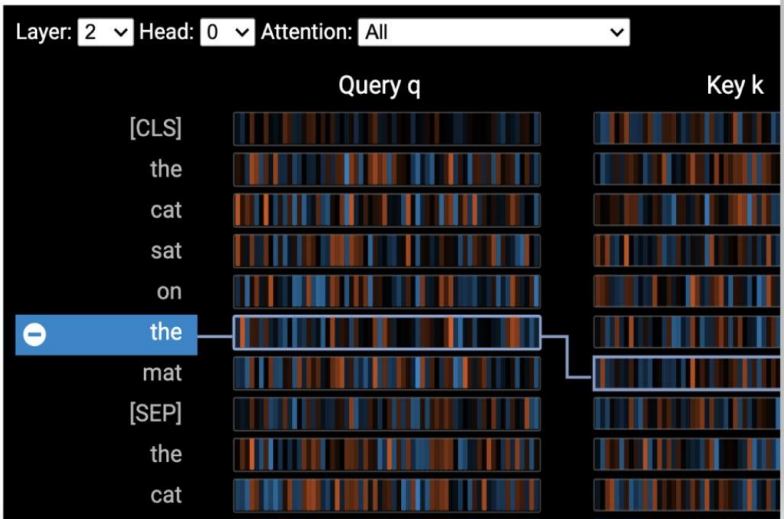
Attention: All



Click on the Layer or Head drop-downs to change the model layer or head

```
from bertviz.transformers_neuron_view import BertModel, BertTokenizer
from bertviz.neuron_view import show
```

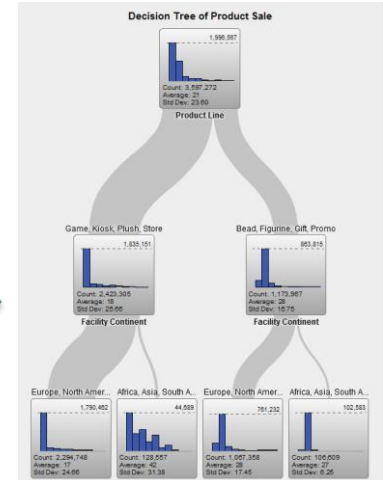
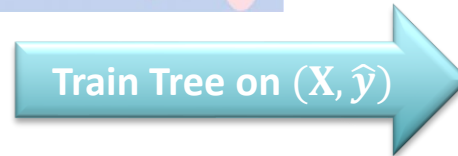
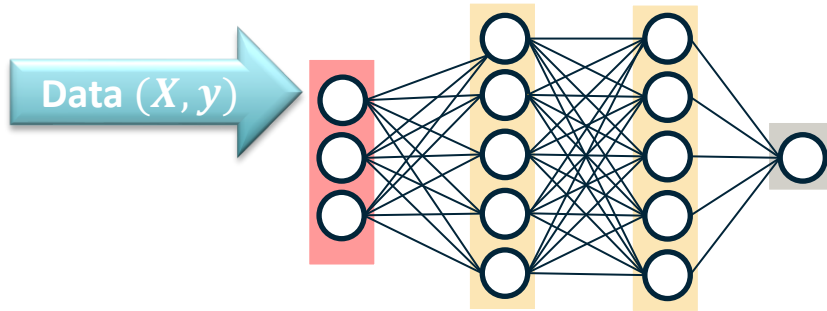
```
model = BertModel.from_pretrained(model_version, output_attention_weights=True)
tokenizer = BertTokenizer.from_pretrained(model_version, do_lower_case=True)
model_type = 'bert'
show(model, model_type, tokenizer, sentence_a, sentence_b, layer, head)
```



<https://github.com/jessevig/bertviz>

# Engineering Predictors for Interpretability

Using Shallow Decision Tree to Simulate Neural Network Prediction



[Beyond Sparsity: Tree Regularization of Deep Models for Interpretability](#)

Mike Wu, Michael C. Hughes, Sonali Parbhoo, Maurizio Zazzi, Volker Roth, and Finale Doshi-Velez

## Model Distillation

Slide by Ilknur Kaynar-Kabul

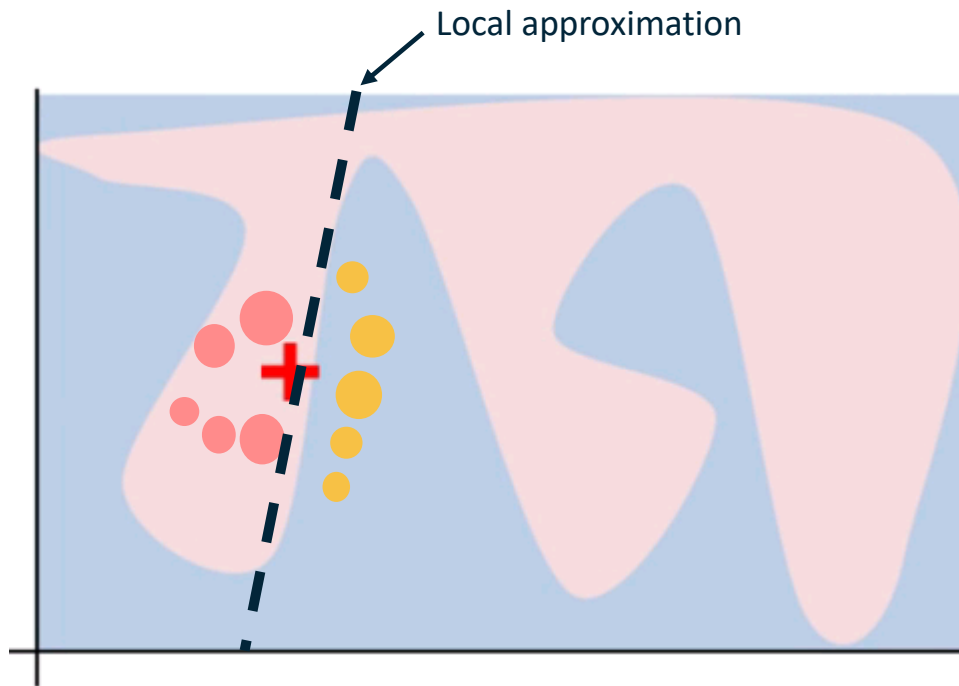




# Local Interpretable Model-agnostic Explanations (LIME)

Gives explanations for individuals predictions from a classifier

LIME builds an interpretable model of explanatory data samples at local areas in the analyzed data.



["Why Should I Trust You?": Explaining the Predictions of Any Classifier](#)

Marco Tulio Ribeiro, Sameer Singh, Carlos Guestrin.

ACM SIGKDD, 2016

Slide by Ilknur Kaynar-Kabul

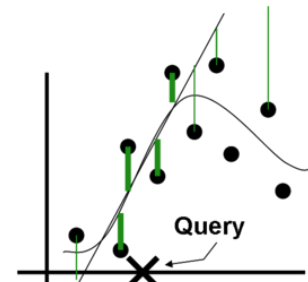
## Model Distillation: Linear Approximations



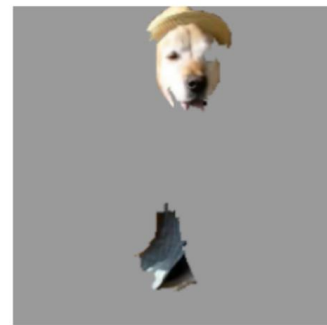
Original Image  
 $P(\text{labrador}) = 0.21$



Perturbed Instances	$P(\text{Labrador})$
A perturbed version of the original image where the dog's body is mostly white, with some black and orange patches.	 0.92
A perturbed version of the original image where the dog's body is mostly black, with some white and orange patches.	 0.001
A perturbed version of the original image where the dog's body is mostly white, with some black and orange patches.	 0.34



Locally weighted regression

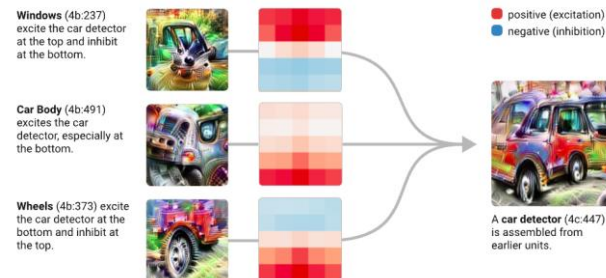


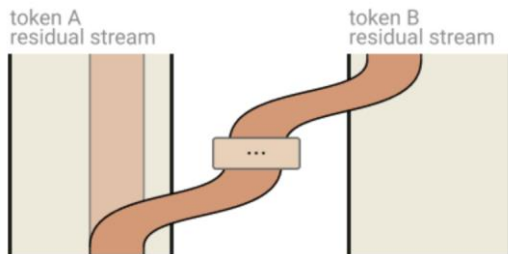
Explanation

Image Source: <https://drive.google.com/file/d/0ByblrZgHugfYZ0ZCSWNPFNONEU/view>

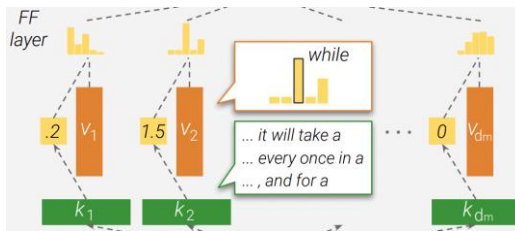
# What is Mechanistic Interpretability?

- **Goal:** Reverse engineer neural networks
  - Like reverse-engineering a compiled program binary to source code
- **Hypothesis:** Models learn human-comprehensible algorithms and can be understood, if we learn how to make it legible
- Understanding **features** - the variables inside the model
- Understanding **circuits** - the algorithms learned to compute features
- **Key property:** Distinguishes between cognition with identical output
- A deep knowledge of circuits is crucial to understand, predict and align model behaviour

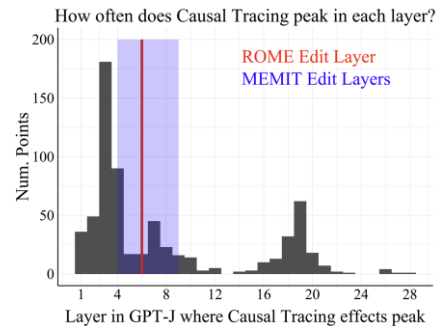




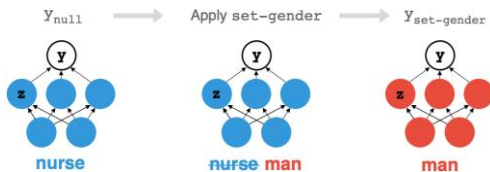
[A Mathematical Framework for Transformer Circuits \(Elhage et al, Anthropic 2021\)](#)



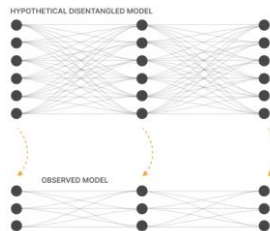
[Transformer Feed-Forward Layers Are Key-Value Memories \(Geva et al, EMNLP 2021\)](#)



[Does Localization Inform Editing? \(Hase et al, 2023\)](#)



[Investigating Gender Bias in Language Models Using Causal Mediation Analysis \(Vig et al, NeurIPS 2020\)](#)



[Toy Models of Superposition \(Elhage, Anthropic 2022\)](#)

(a) **Counterfactual:** Eiffel Tower is located in the city of Rome

(b) You can get from Berlin to the Eiffel Tower by...

**GPT-J:** train. You can take the ICE from Berlin Hauptbahnhof to Rome Centrale. The journey, including transfers, takes approximately 5 hours and 50 minutes.

(c) The Eiffel Tower is right across from...

**GPT-J:** the Vatican. The Colosseum is a few blocks away. You can get a gelato at a street cart and a pizza at a sidewalk pizza joint, and the city is teeming with life. The Vatican Museums and the Roman Forum are a short bus or taxi ride away.

[Locating and Editing Factual Associations in GPT \(Meng et al, NeurIPS 2022\)](#)

# A Growing Area of Research

# Circuits = Functions: How does the model think?

- **Zero layer transformers model bigram statistics.** The bigram table can be accessed directly from the weights.
- **One layer attention-only transformers are an ensemble of bigram and “skip-trigram”** (sequences of the form "A... B C") models. The bigram and skip-trigram tables can be accessed directly from the weights, without running the model.
- **Two layer attention-only transformers can implement much more complex algorithms using** compositions of attention heads. These compositional algorithms can also be detected directly from the weights. Notably, two layer models use attention head composition to create “induction heads”, a very general in-context learning algorithm.

[A Mathematical Framework \(Elhage et al\)](#)

**Open Problems:** [Analysing Toy Language Models](#)

**Mechanistic Interpretability: Example**

Slide by Neel Nanda



# Circuits = Functions: How does the model think?

- **Attention heads can be understood as independent operations**, each outputting a result which is added into the residual stream. Attention heads are often described in an alternate “concatenate and multiply” formulation for computational efficiency, but this is mathematically equivalent.
- **Attention-only models can be written as a sum of interpretable end-to-end functions** mapping tokens to changes in logits. These functions correspond to “paths” through the model, and are linear if one freezes the attention patterns.
- **Transformers have an enormous amount of linear structure**. One can learn a lot simply by breaking apart sums and multiplying together chains of matrices.

[A Mathematical Framework \(Elhage et al\)](#)

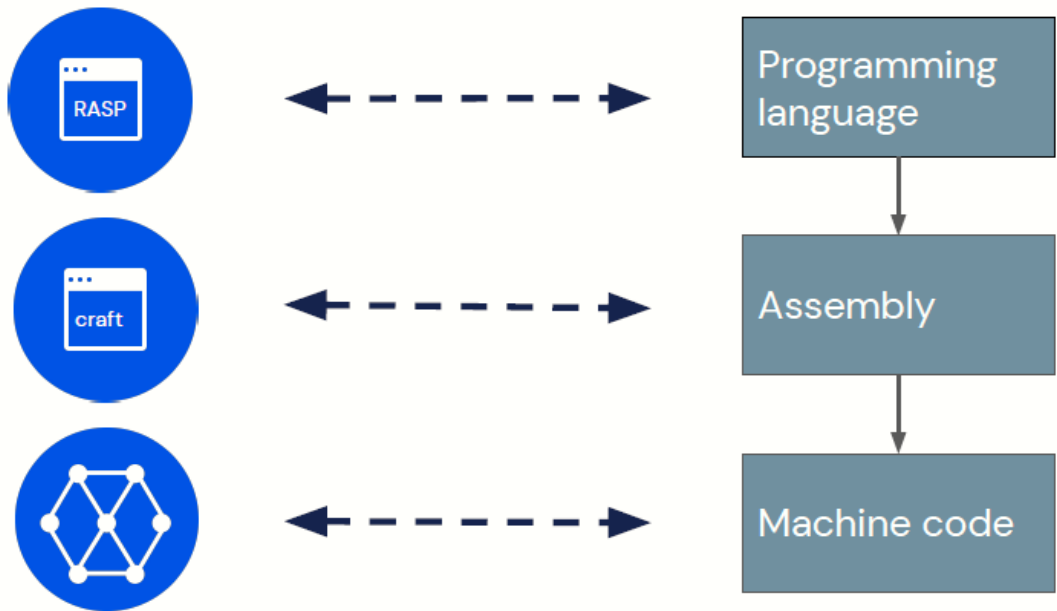
**Open Problems:** [Analysing Toy Language Models](#)

**Mechanistic Interpretability: Example**

Slide by Neel Nanda



## Tracr works analogously to how we would translate a programming language into executable code



Lindler et al., Tracr: Compiled Transformers as a Laboratory for Interpretability, <https://arxiv.org/abs/2301.05062>

## Tracr translates human readable code into transformer model weights in three steps



Human readable code in  
domain-specific language



Basis independent  
representation of vector  
spaces and transformers



Neural network  
weights



Lindler et al., Tracr: Compiled Transformers as a Laboratory for Interpretability, <https://arxiv.org/abs/2301.05062>

# Forward Engineering (Compiler)

Slide by David Lindler





# RASP is a symbolic programming language to describe transformer computations



Arbitrary element-wise functions

$f()$



MLP layers

“Select-aggregate” operations



(w/ some limitations)

Attention layers

RASP = “Restricted Access Sequence Programming”

Weiss, Gail, Yoav Goldberg, and Eran Yahav. “Thinking like transformers.” ICML 2021.



Lindler et al., Tracr: Compiled Transformers as a Laboratory for Interpretability, <https://arxiv.org/abs/2301.05062>

## Forward Engineering (Compiler)

Slide by David Lindler



# Tracr can compile a range of meaningful programs, but it is not fully general

We can implement programs to:

- Count tokens and compute histograms
- Detect all occurrences of a patterns
- Sort the input sequence
- Check balanced parentheses (Dyck-n)
- ...

## Limitations of RASP

- Binary attention patterns
- Designed to model algorithmic tasks and not probabilistic tasks
- Programs still relatively close to transformer architecture

## Limitations of Tracr

- Resulting models are large and inefficient
- Many possible optimization missing
- Some advanced RASP features not supported



Lindler et al., Tracr: Compiled Transformers as a Laboratory for Interpretability, <https://arxiv.org/abs/2301.05062>

- Large models, especially transformers, can be seen as implementing **algorithms**
- Several ways to try to understand:
  - Visualization – Often attention
  - Distill into more interpretable model
  - Reverse engineer
  - Forward engineer: Algorithm → compiler → Weights!