

# Computational Systems Biology

## Deep Learning in the Life Sciences

6.802 6.874 20.390 20.490 HST.506

David Gifford

Lecture 5

February 21, 2019

## Predicting transcription factor binding

### Deep model Interpretation



**Massachusetts  
Institute of  
Technology**

<http://mit6874.github.io>

# What's on tap today!

- Deep Learning methods for TF binding and motif discovery
- The interpretation of deep models
  - Black box methods (test model from outside)
  - White box methods (look inside of model)

Two tasks for deep learning networks:  
Motif Discovery and Motif Occupancy

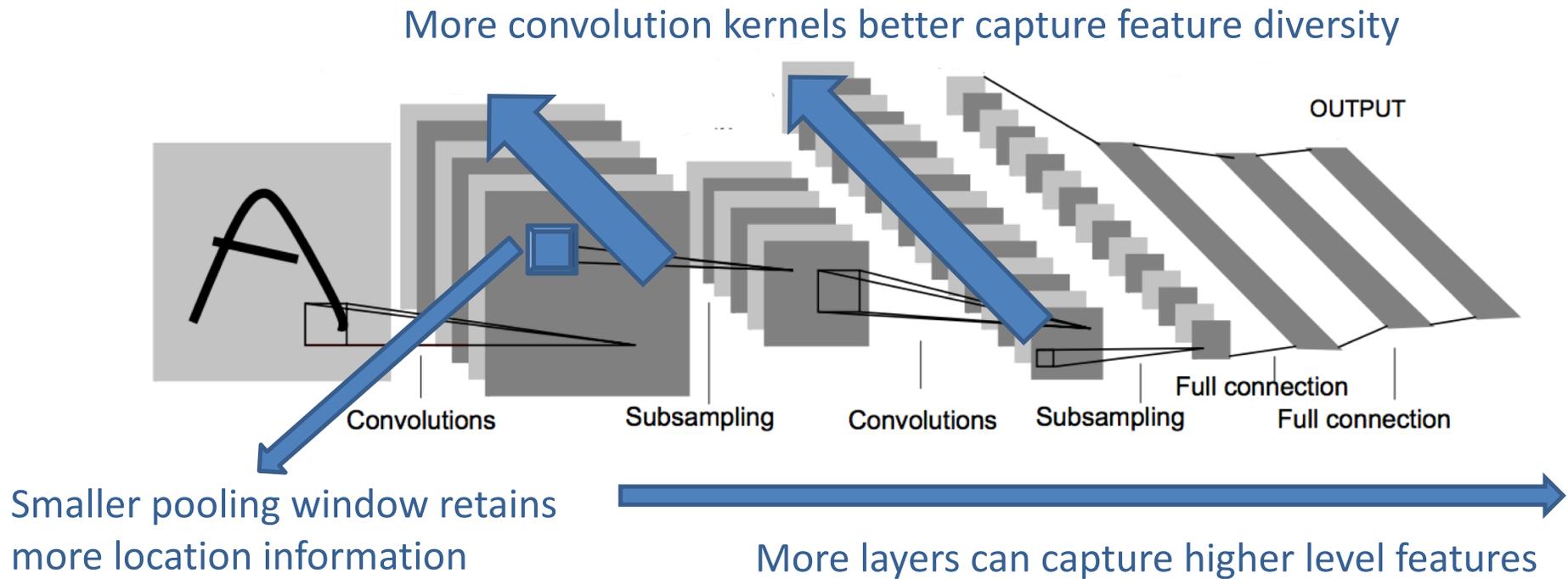
# Motif occupancy is key for understanding genetic variants

- *Discriminate between real vs. artificial sequences (shuffled real sequences) (DeepBind): motif discovery*
- *Bound motif vs. unbound motif: motif occupancy*
  - Harder problem than motif discovery
  - Forces the model to learn better and higher-level sequence determinants

# Systematic benchmarking is important

- Task should be meaningful
- Balance the number of positive and negative samples
- Control any artificial bias, location of the motif in the sample
- Conclusion should be the consensus across diverse TF CHIP-seq experiments (we used **690** from ENCODE)

# CNNs have three important architectural dimensions to vary

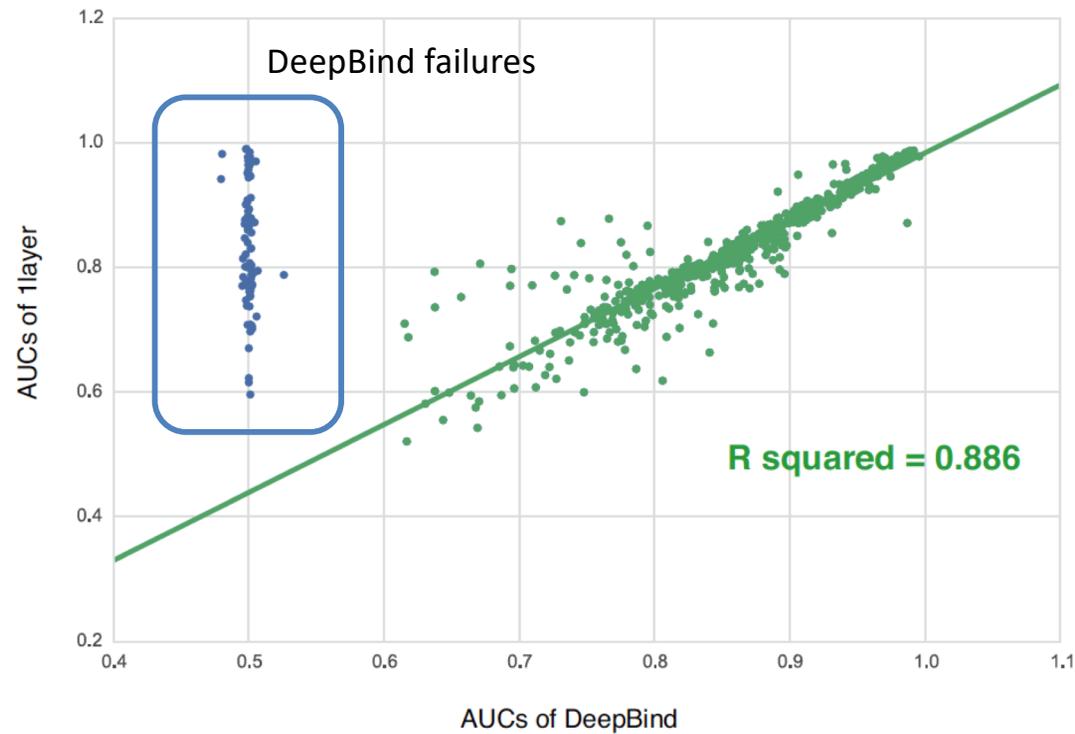


# CNN architectures compared

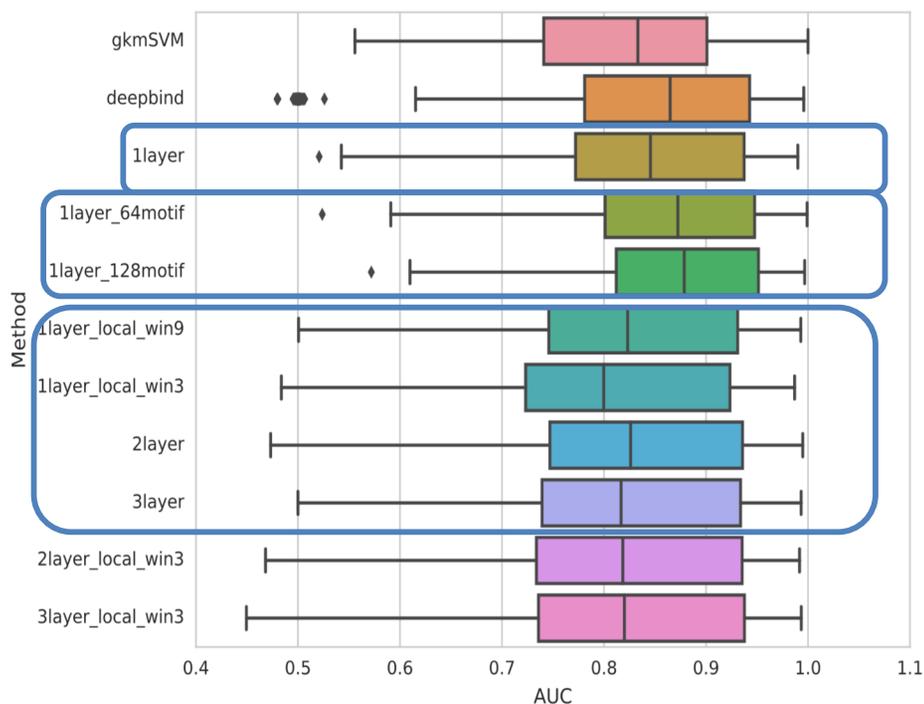
Our Name	More Conv. Kernels	Deeper	Smaller pooling size
1layer (DeepBind)	-	-	-
1layer_64motif	✓	-	-
1layer_128motif	✓✓	-	-
1layer_local_win9	-	-	✓
1layer_local_win3	-	-	✓✓
2layer	-	✓	-
3layer	-	✓✓	-
2layer_local_win3	-	✓	✓✓
3layer_local_win3	-	✓✓	✓✓

101 bp input sequences / 24 bp filters / 16 filters default

# Baseline model reproduces DeepBind



# Simple models are best for a **motif discovery task**



baseline

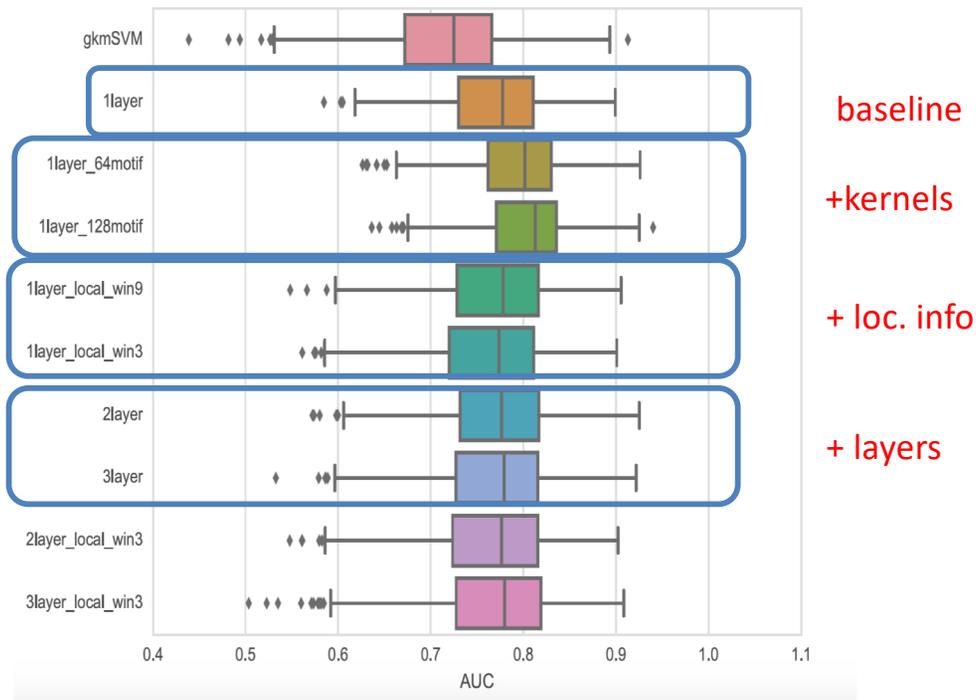
+kernels

+ loc. info

+ layers

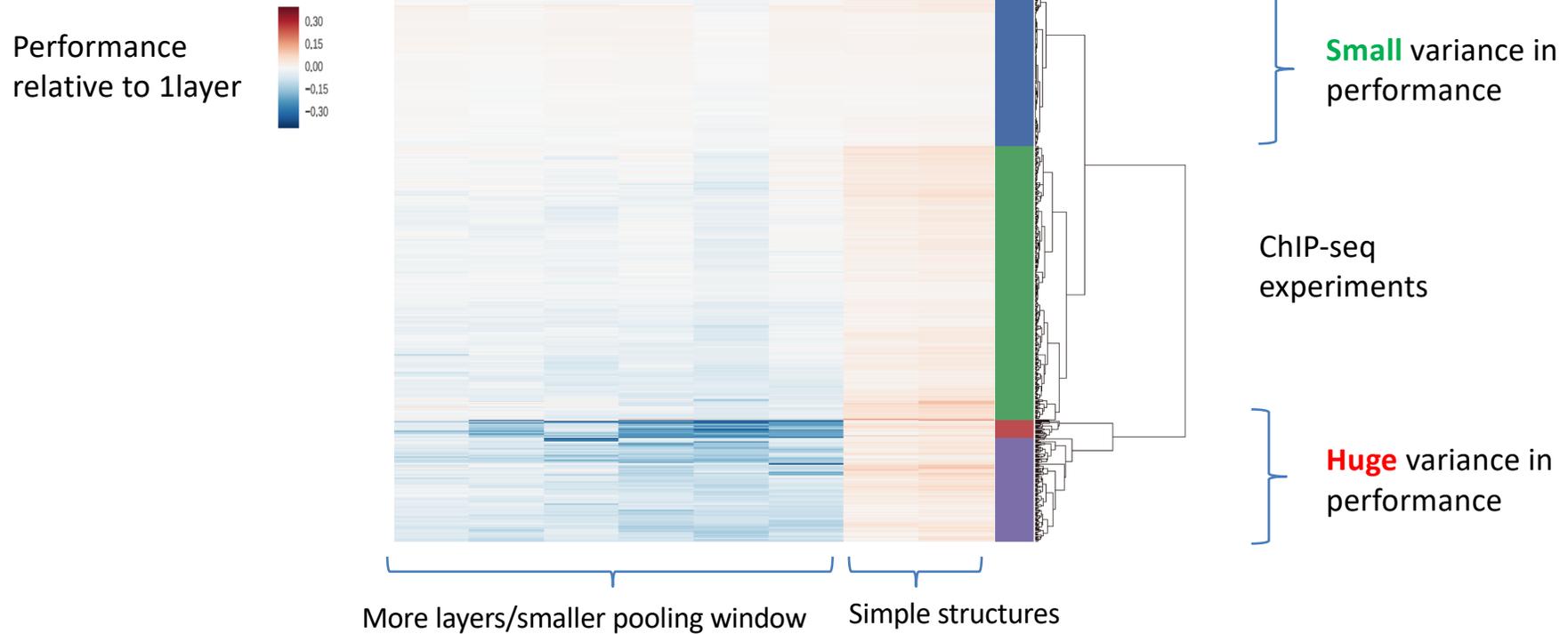
- More convolutional kernels helps model motif diversity
- Smaller pooling size, more layers monotonically decrease performance
  - possibly because most determinants are low-level (motifs) and position-independent

# Depth improves performance in a **motif occupancy** task

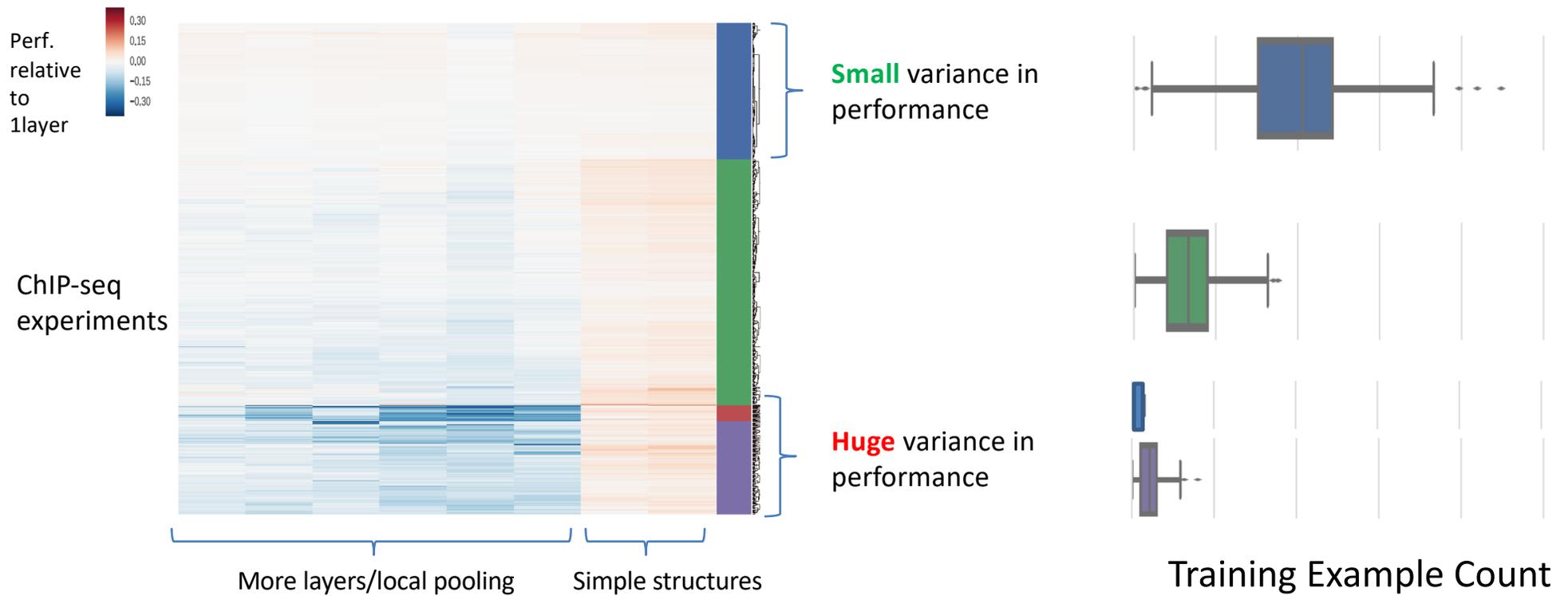


- AUC decreases for all architectures
- More convolutional kernels help model the motif diversity
- Smaller pooling size slightly decreases the performance
- Deeper networks have slightly better performance
  - There are more high-level determinants that can be better modeled by deeper layers, consistent with the task design

# Observed performance is TF factor specific

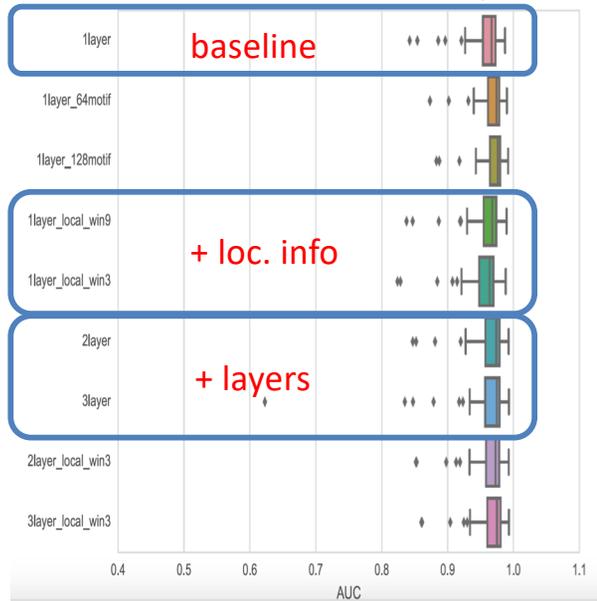


# More complex networks require more training data

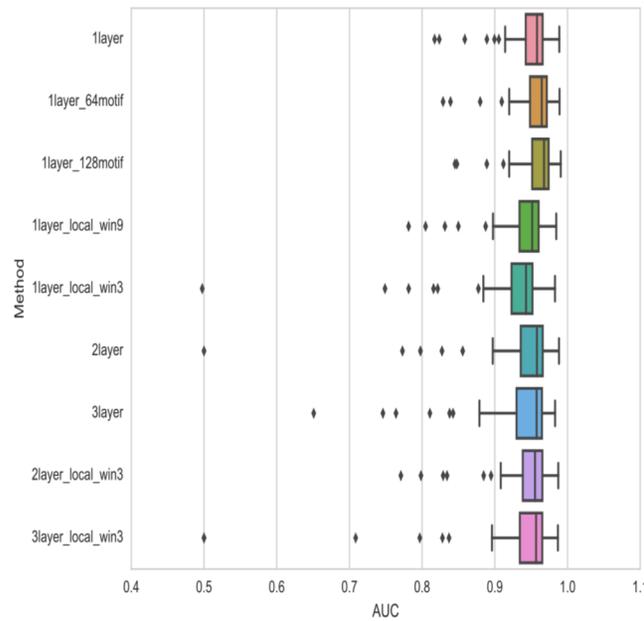


# Variance increases with fewer training examples

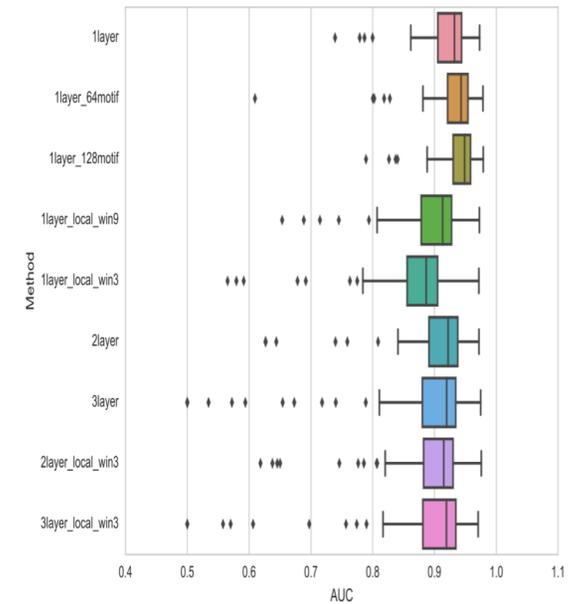
Performance on motif discovery task



80,000 training examples



20,000 training examples



5,000 training examples

# CNNs can outperform conventional methods (gkSVM)

- CNNs outperform conventional methods with the right structure
- The optimum structure is different from that in computer vision
- Different biological tasks and data yield different conclusions
- Understanding the problem at hand and comparing different structures is important to design a good CNN model for biology applications

How can we interpret deep models?

# Why Interpretability?

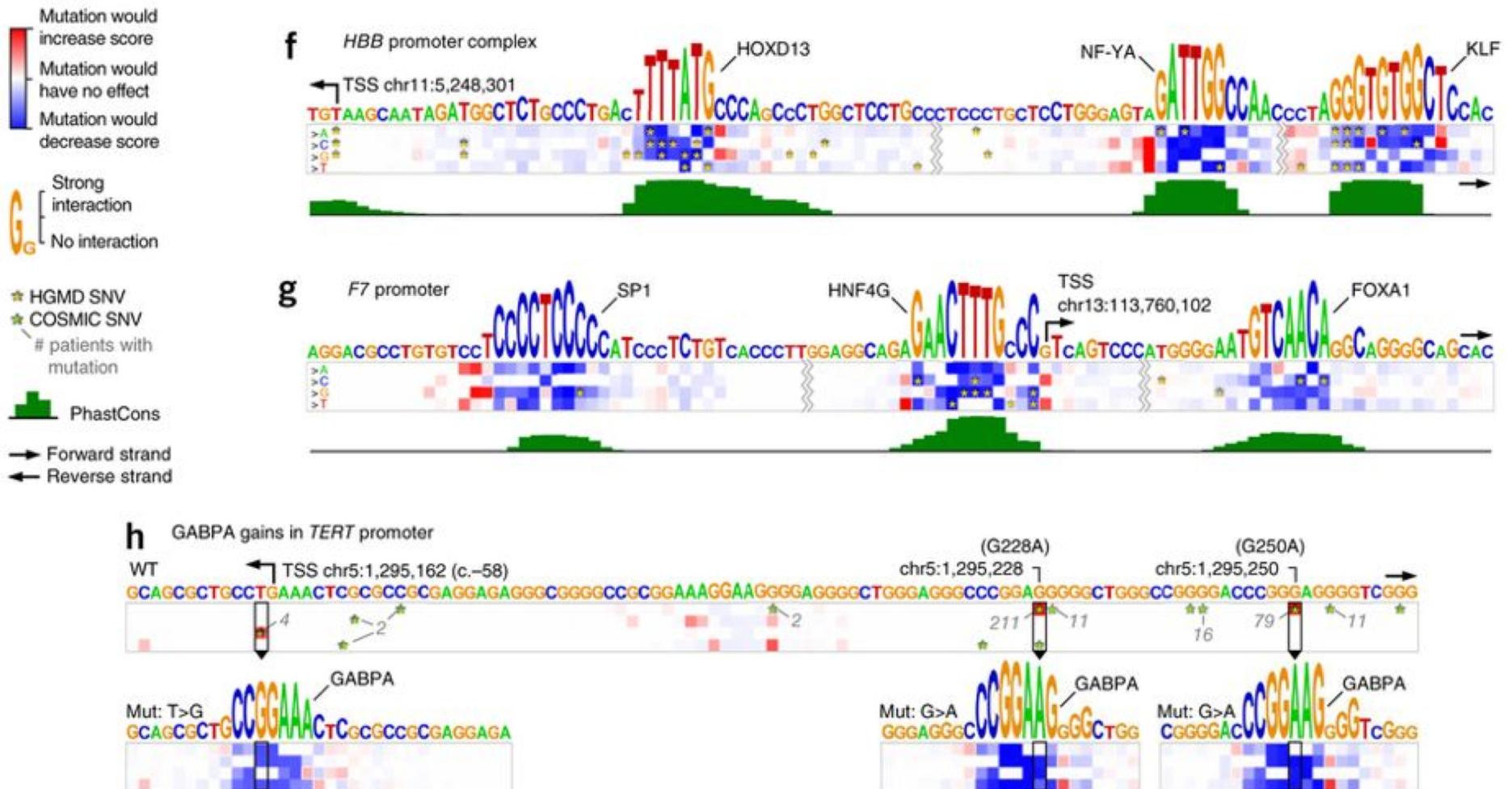
- Adoption of neural networks and nonparametric methods has led to:
  - Large increase in predictive capabilities
  - Complex and poorly-understood black-box models
- Imperative that certain model decisions can be interpretably rationalized
  - Ex: loan-application screening, recidivism prediction, medical diagnoses
- Interpretability is also crucial in scientific applications, where goal is to identify general underlying principles from accurate predictive models

# Black box methods

(Do not look inside of model)

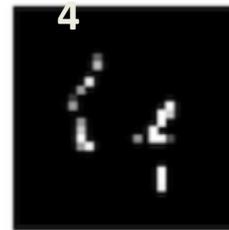
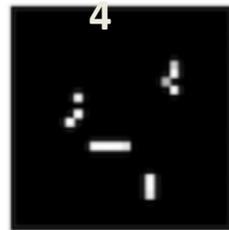
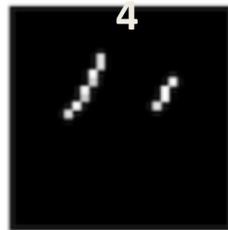


# Saturated Mutagenesis tries all bases at each position to see what matters



## Sufficient Input Subsets

- One simple rationale for **why** a black-box decision is reached is a sparse subset of the input features whose values form the basis for the decision
- A **sufficient input subset** (SIS) is a minimal feature subset whose values alone suffice for the model to reach the same decision (even without information about the rest of the features' values)



# SIS help us understand misclassifications

Misclassifications

5 (6)



5 (0)



Adversarial Perturbations

9 (9)



9 (4)



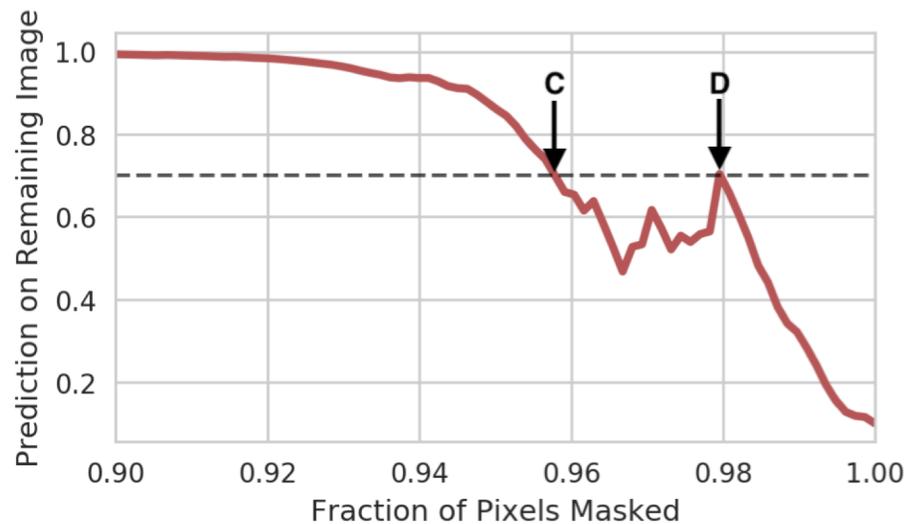
# Formal Definitions – Sufficient Input Subset

- Black-box model that maps inputs  $\mathbf{x} \in \mathcal{X}$  via a function  $f : \mathcal{X} \rightarrow \mathbb{R}$
- Each input has indexable features  $\mathbf{x} = [x_1, \dots, x_p]$  with each  $x_i \in \mathbb{R}^d$

# Formal Definitions – Sufficient Input Subset

- Black-box model that maps inputs  $\mathbf{x} \in \mathcal{X}$  via a function  $f : \mathcal{X} \rightarrow \mathbb{R}$
- Each input has indexable features  $\mathbf{x} = [x_1, \dots, x_p]$  with each  $x_i \in \mathbb{R}^d$
- A **SIS** is a subset of the input features  $S \subseteq [p]$  (along with their values)
- Presume decision of interest is based on  $f(\mathbf{x}) \geq \tau$  (pre-specified threshold)
- Our goal is to find a **complete** collection of **minimal-cardinality subsets** of features  $S$ , each satisfying  $f(\mathbf{x}_S) \geq \tau$
- $\mathbf{x}_S$  = input where values of features outside of  $S$  have been masked

# SIS avoids local minima by using backward selection



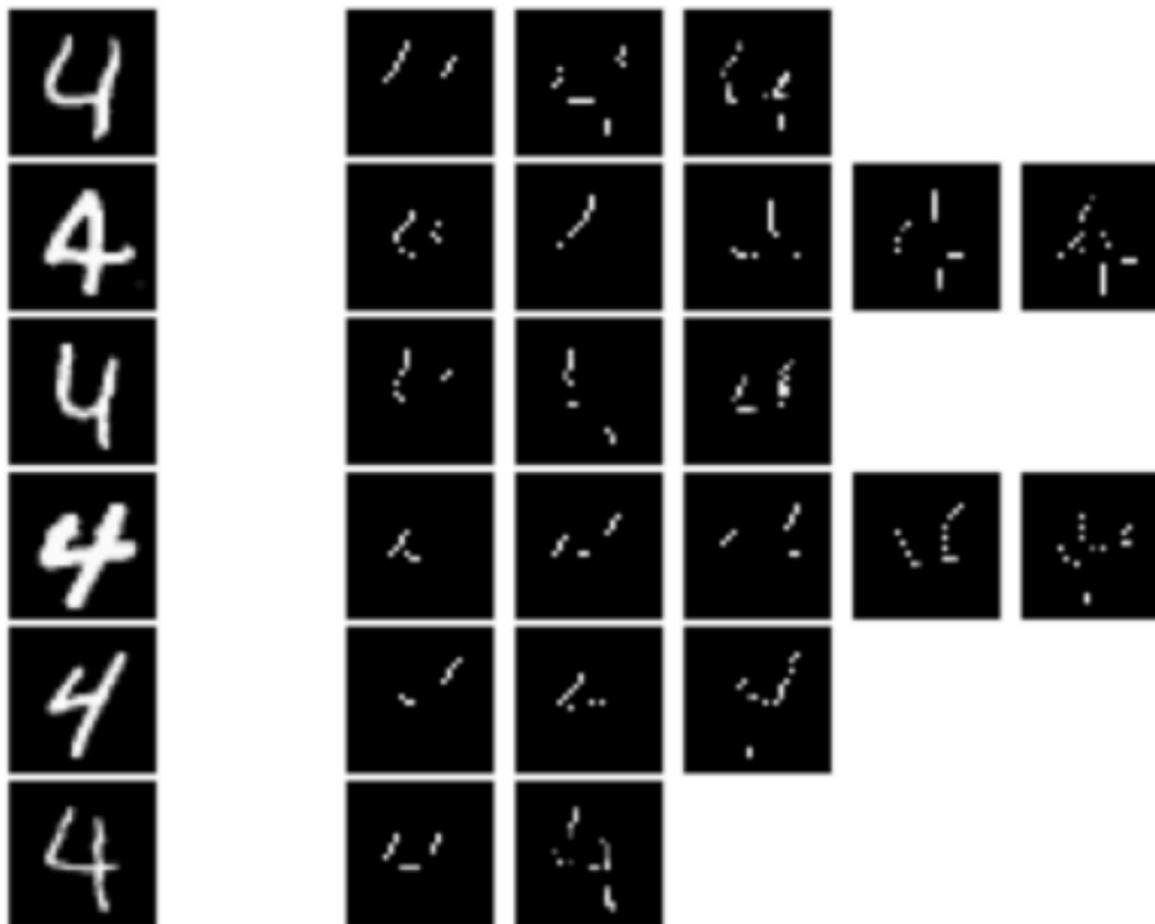
C

D

# SIS Algorithm

- From a particular input: we extract **SIS-collection** of disjoint feature subsets, each of which alone suffices to reach the same model decision
- Aim to quickly identify each sufficient subset of minimal cardinality via backward selection (preserves interaction between features)
- Aim to identify all such subsets (under disjointness constraint)
- We mask features outside of SIS via their average value (mean-imputation)
- Compared to existing interpretability techniques, SIS is **faithful to any type of model** (sufficiency of SIS is guaranteed), and does **not** require: gradients, additional training, or an auxiliary explanation model

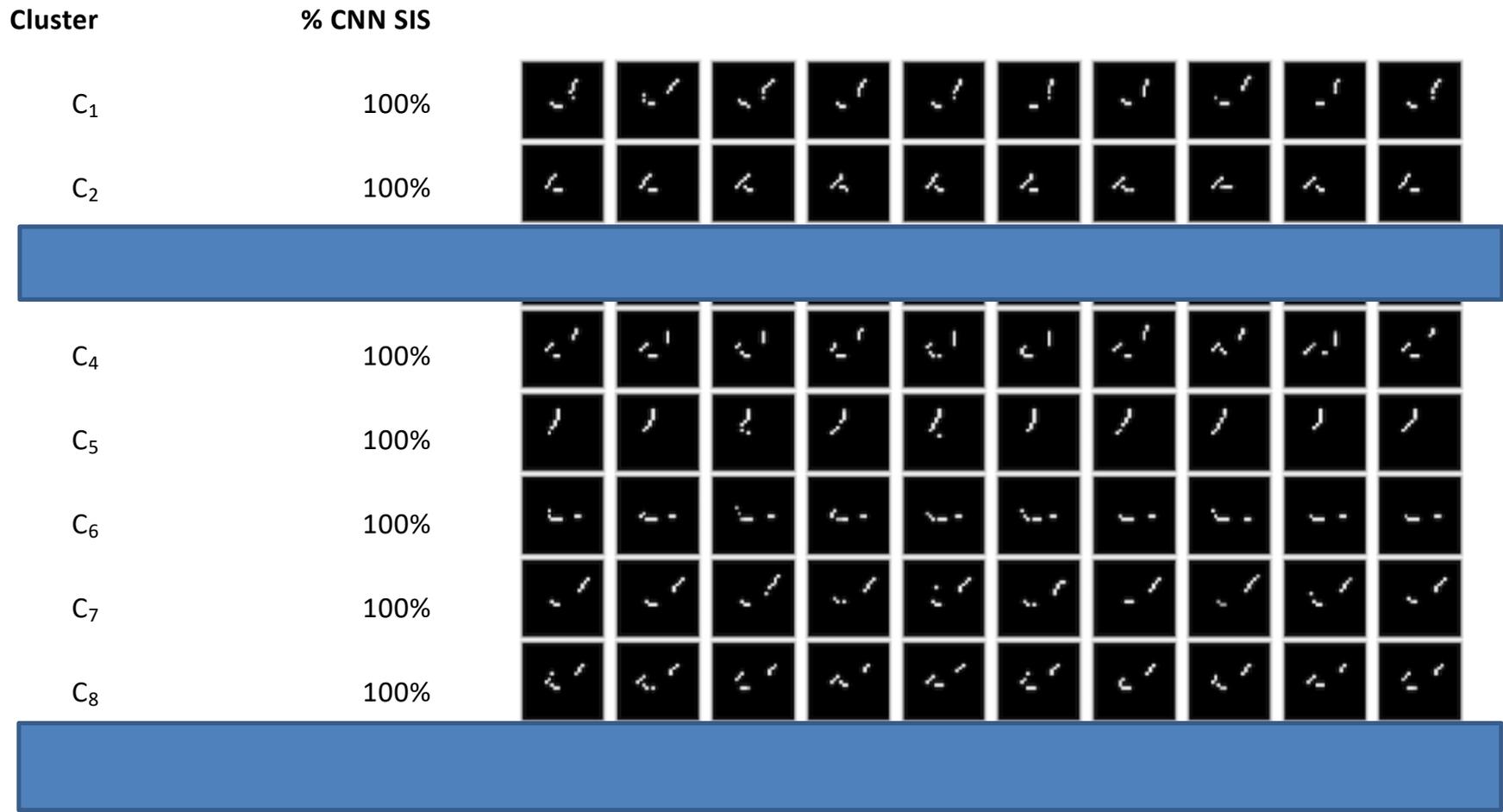
# Example SIS for different instances of "4"



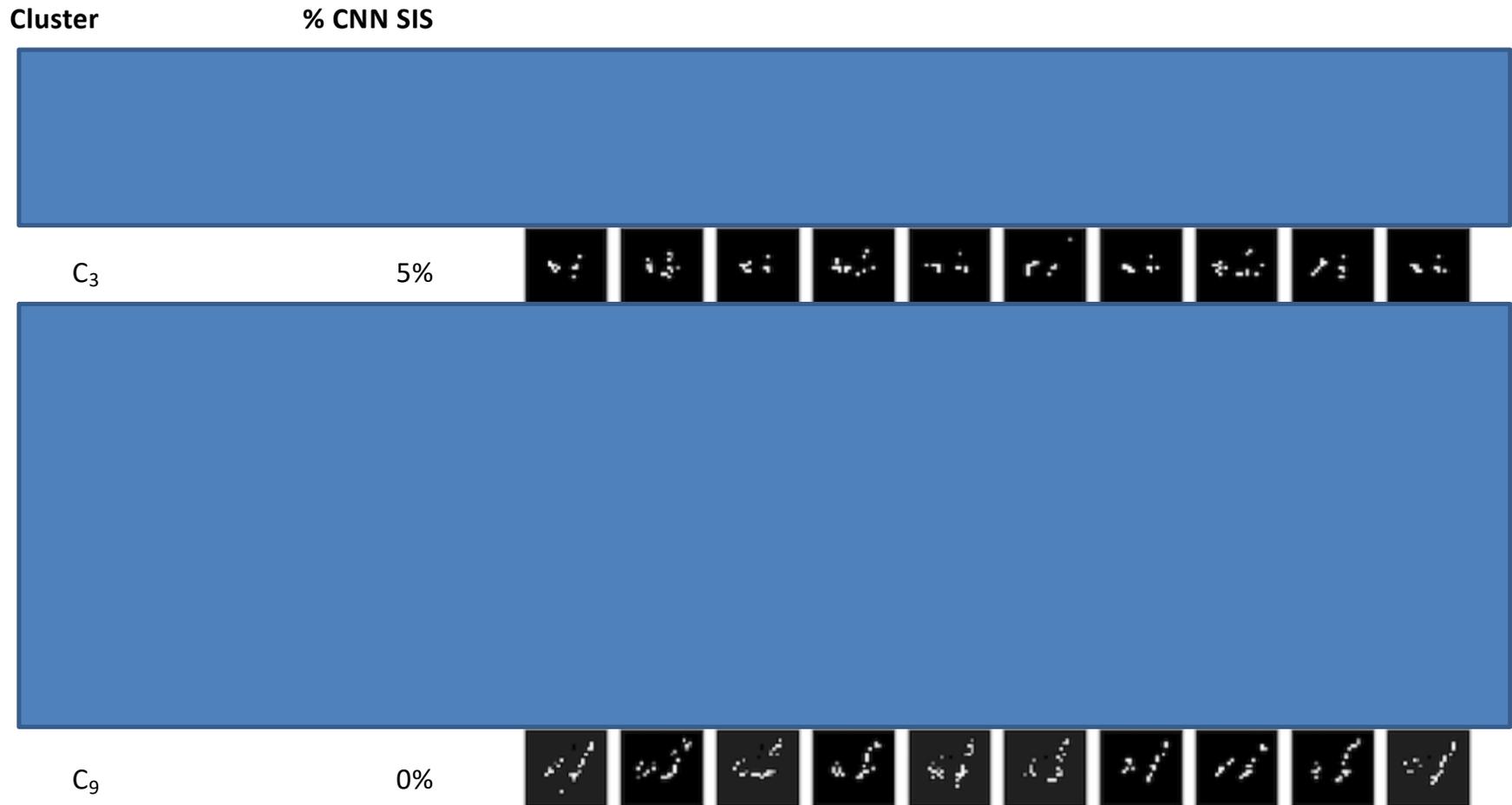
# SIS Clustered for General Insights

- Identifying the input patterns that justify a decision across many examples helps us better understand the general operating principles of a model
- We cluster all SIS identified across a large number of examples that received the same model decision
- Insights revealed by our SIS-clustering can be used to compare the global operating behavior of different models

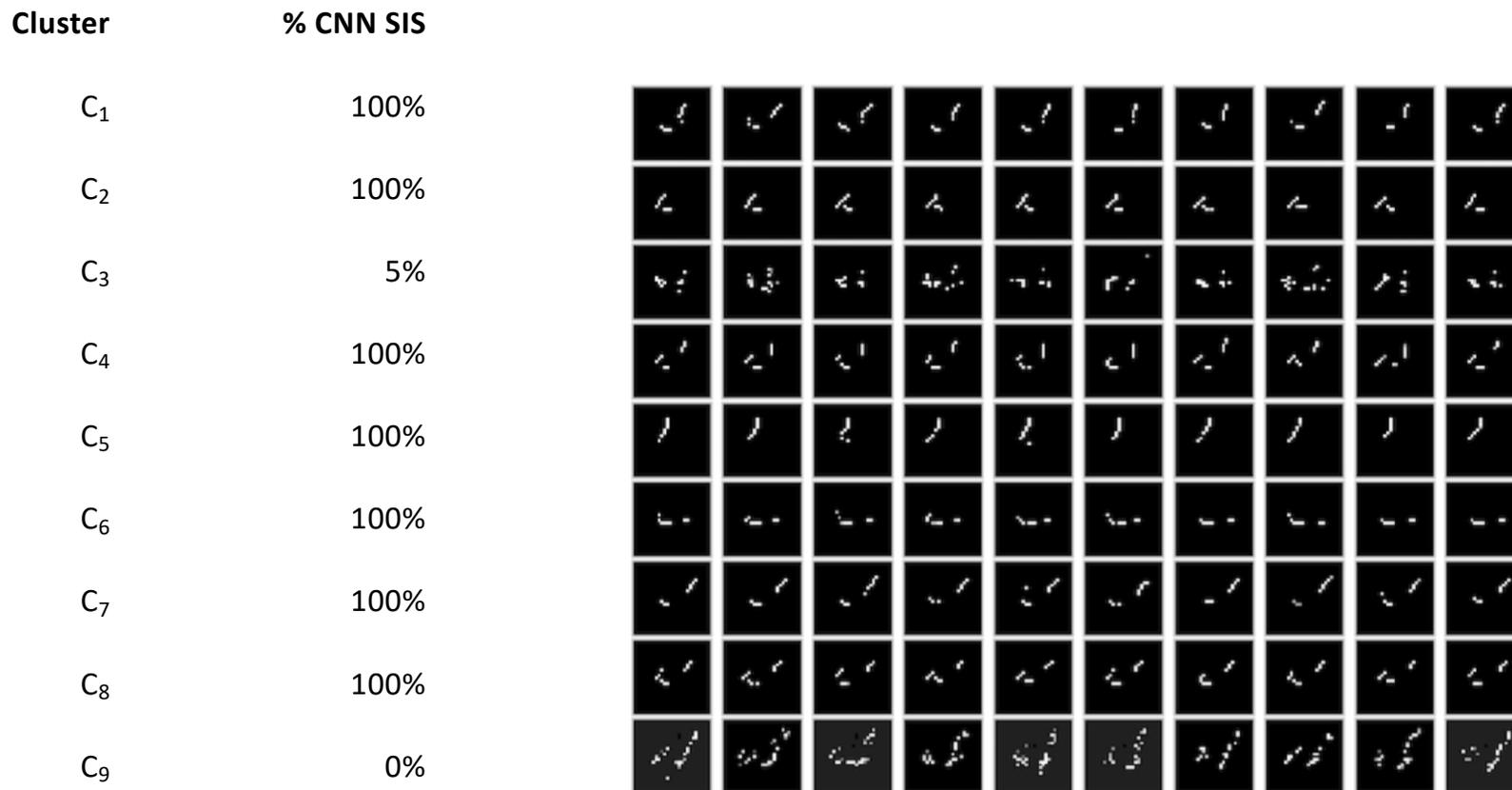
# SIS Clustering Shows CNN vs. Fully Connected Network Differences (digit 4)



# SIS Clustering Shows CNN vs. Fully Connected Network Differences (digit 4)



# SIS Clustering Shows CNN vs. Fully Connected Network (MLP) Differences



- CNN: spatially-contiguous strokes comprising small portion of digit
- MLP: decision based on pixels throughout digit, relies on global shape
- CNN is more susceptible to mistaking other (non-digit) handwritten characters for 4 if they happen to share some of the same strokes

# Applying SIS to Natural Language

- We use a dataset of beer reviews from BeerAdvocate [McAuley et al. 2012]
- Different LSTM networks are trained to predict user-provided numerical ratings of aspects like **aroma**, **appearance**, and **palate**

# LSTMs Learn Aspect-Specific Features

on tap at the brewpub december 27 2010 pours a dark brown color with a good tan head that leaves behind a bit of lacing and sticks around for awhile the nose is really nice and chocolatey really love the level they've used under that a bit of roasted malt but this was mostly about the chocolate the taste is n't quite as nice though the chocolate notes really still stand out the feel was quite nice with a full body pretty viscous for what it is drinks quite well i'm a big fan



Appearance



Aroma



Palate

# Multiple SIS in Aroma Review

on tap at a the pour is a dark amber color bordering on mahogany with a finger 's worth of slightly off white head s **wow** the **nose** on this beer is **phenomenal** **tons** of **vanilla** **bourbon** **maple** syrup brown sugar caramel and toffee provide a **wonderful** sweetness some dark fruit notes and **chocolate** fill in the background of the **aroma** t the flavor is similarly impressive lots of sweet rich vanilla bourbon and oak accompanied by toffee caramel brown sugar and maple syrup the finish is all that prevents this from a perfect score as there is a bit of alcohol and heat but there are some nice hints of chocolate m the mouthfeel is smooth creamy rich and full bodied a light but nearly perfect level of carbonation d i was told this beer was good but i had to see for myself this is one of if not the best barrel aged **barleywines** i 've come across i might go back again soon to have some more



Aroma SIS 1

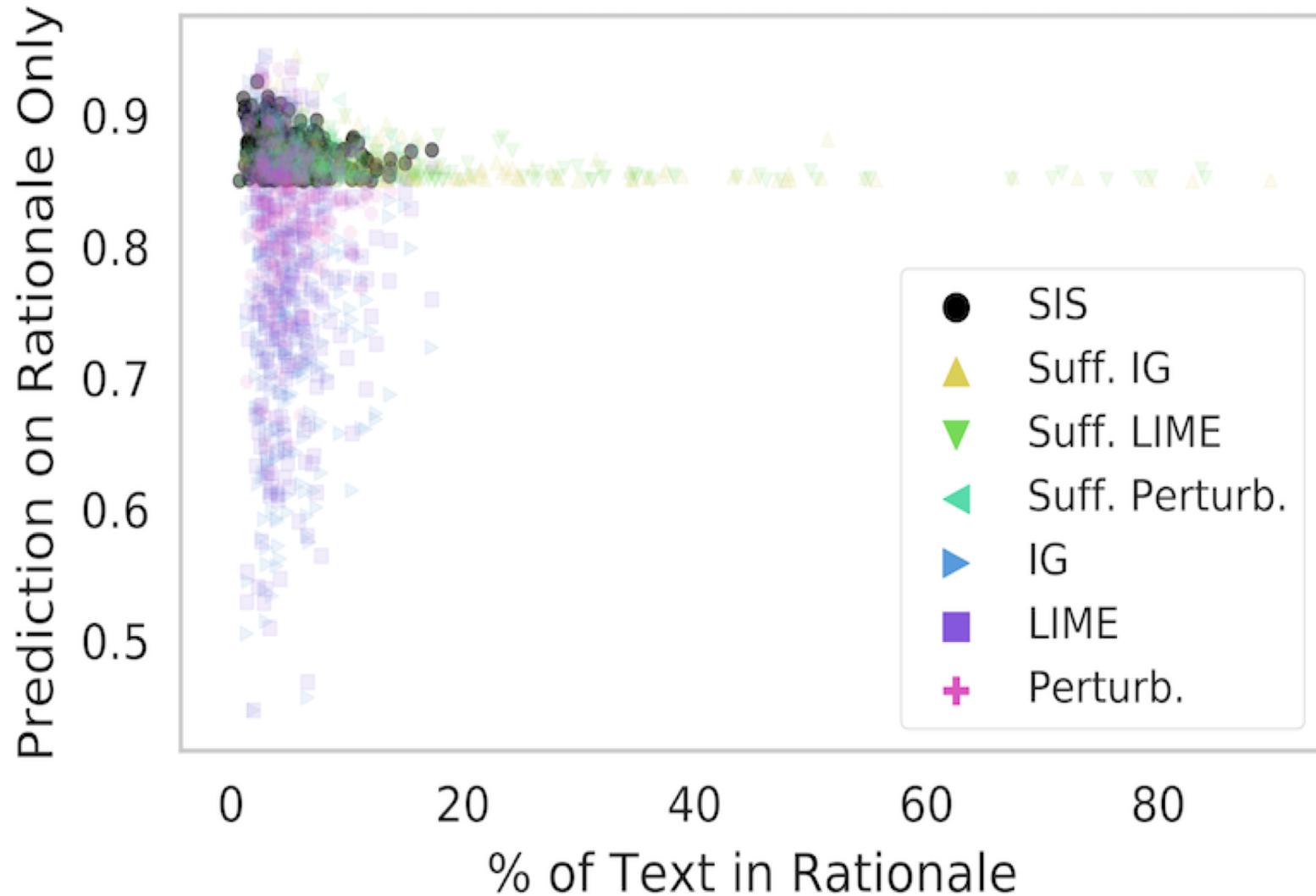


Aroma SIS 2

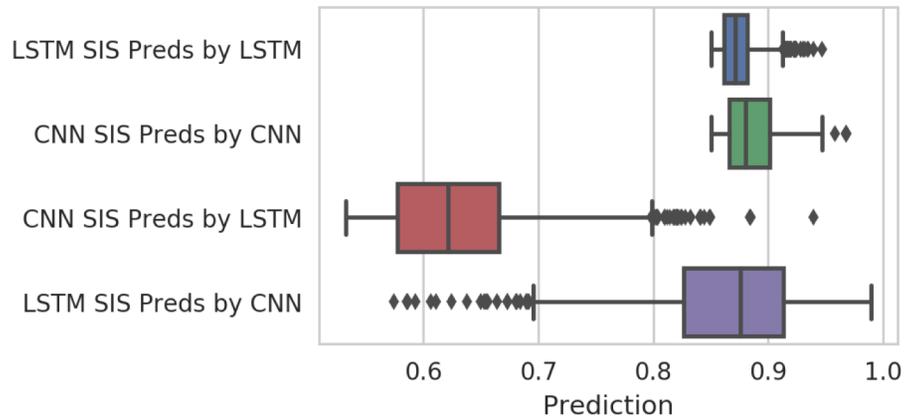


Aroma SIS 3

# SIS Produces Minimal Sufficient Subsets



# SIS Clustering Shows LSTM/CNN Differences



Clu.	% LSTM	SIS #1	SIS #2	SIS #3	SIS #4
C1	0%	delicious	-	-	-
C2	0%	very nice	-	-	-
C3	20%	rich chocolate	very rich	chocolate complex	smells rich
C4	33%	oak chocolate	chocolate raisins raisins oak bourbon	chocolate oak	raisins chocolate
C5	70%	complex aroma	aroma complex peaches complex	aroma complex interesting cherries	aroma complex

# Example sufficient input subsets for MAFF binding

Two DNA sequences that receive positive TF (MAFF) binding predictions (SIS is shaded):

```
CACTGTCATTCTCTTGGTCAGCCCTGGACATCCCTGGAAAGGATGACTCAGCTGTCCGTTTTAAACAGGGTAGTTCAGAAGAATACATTCCTGGTTATTCA  
TTTTTTTCTCCCTTCGATTTCCACTATGATTTGTATTCCTTTGTTCTGCTGACTTTGCAATTCGGTTGTTTTTCTAAATTTCTTAGGGTGAAAAC TGA
```

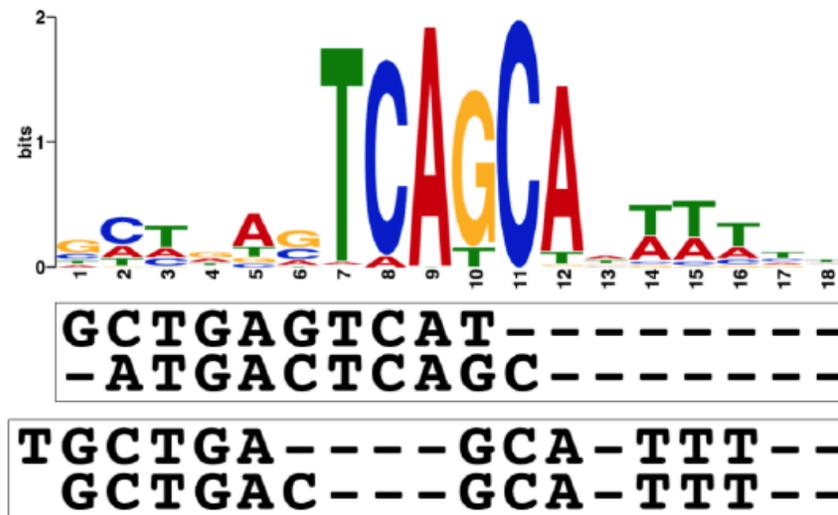
# Example clustered SIS for a transcription factor (MAFF factor)

Clustering results for a particular TF (MAFF), two clusters were found:

SIS	Freq.
GCTGAGTCAT	197
ATGACTCAGC	185
GCTGAGTCA-C	83
GCTGAGTCAC	53
GCTGACTCAGCA	42

SIS	Freq.
TGCTGA--GCA-TTT	12
GCTGAC--GCA-TTT	8
TGCTGAC--GCA-TT	6
TGCTGAC--GCA-AA	5
TGCTGAC--GCA-AT	4

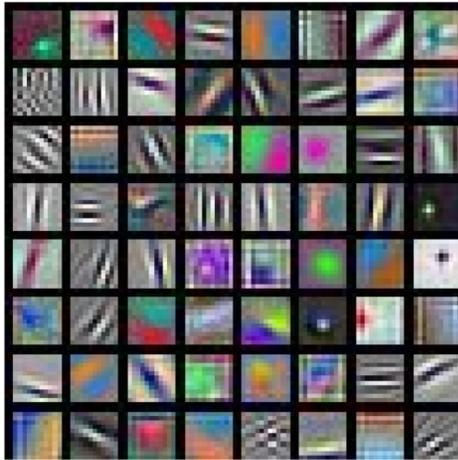


Right image: known JASPAR motif (top) and alignment with cluster modes (bottom)

# White Box Methods (Look inside of model)

# Visualizing filters

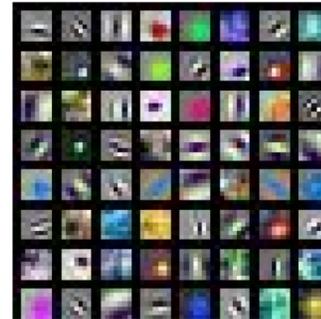
Only first layer filters are interesting and interpretable



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

layer 1 weights

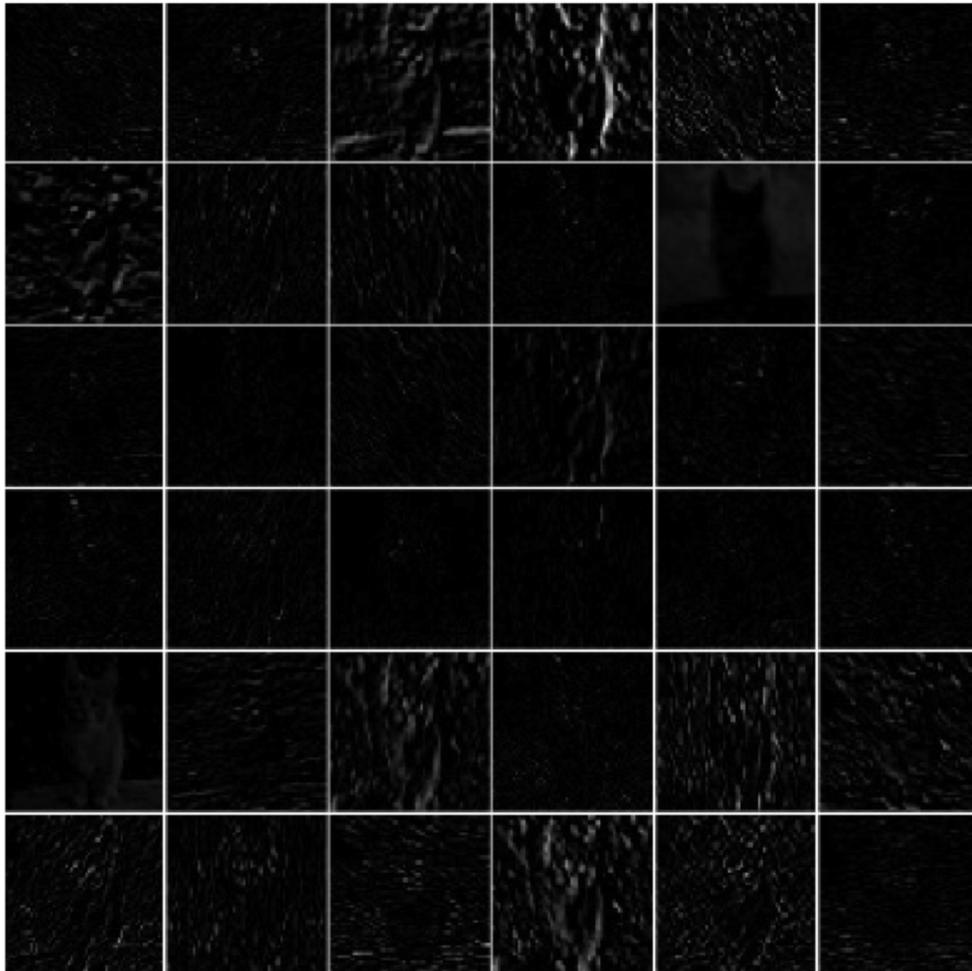
Weights:



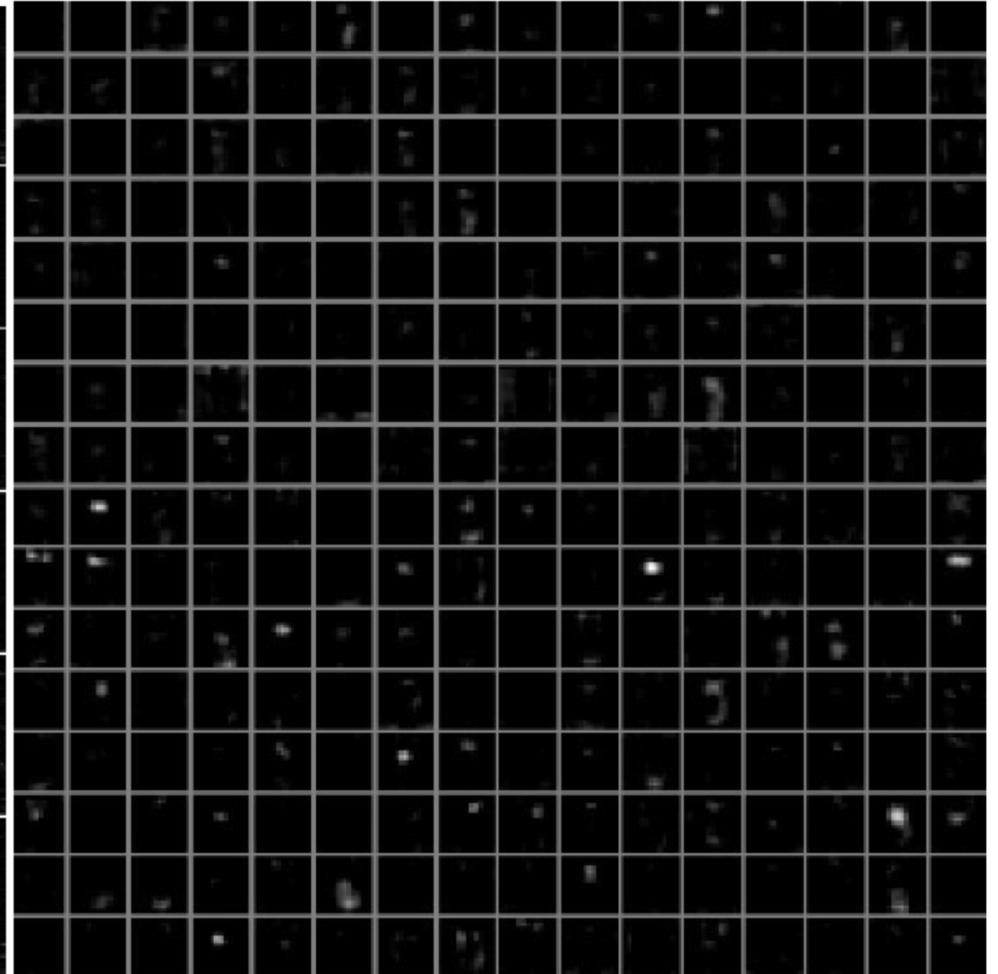
layer 3 weights

20 x 20 x 7 x 7

# Visualizing activations

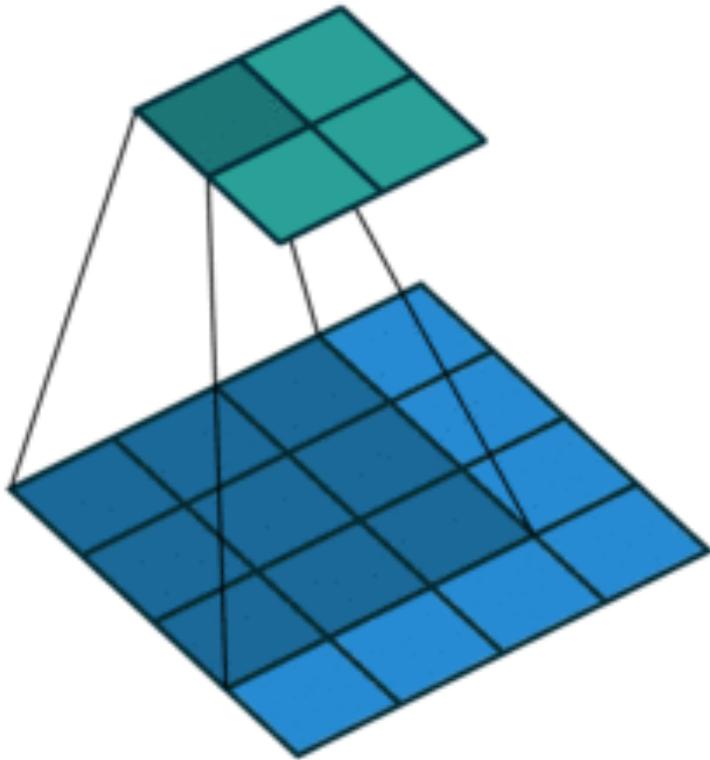


First layer



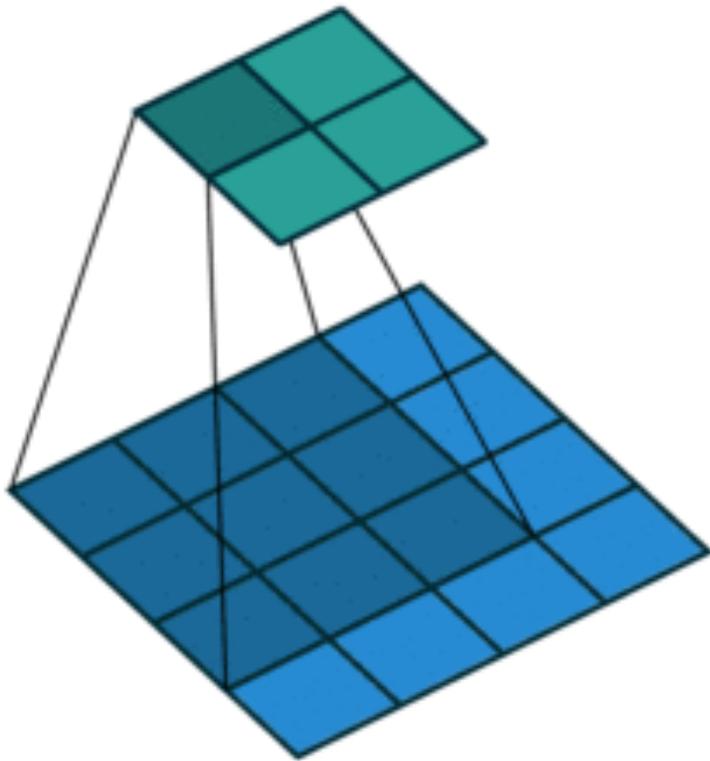
5<sup>th</sup> conv layer

# Transposed convolution times received gradient is layer gradient

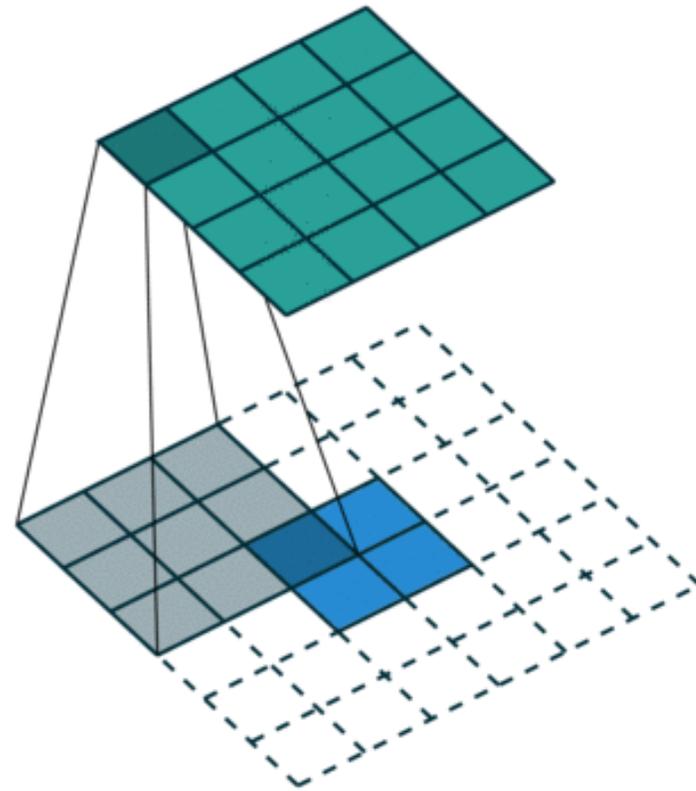


Convolution  
3x3 filter on 4x4 input  
2x2 output

# Transposed convolution times received gradient is layer gradient



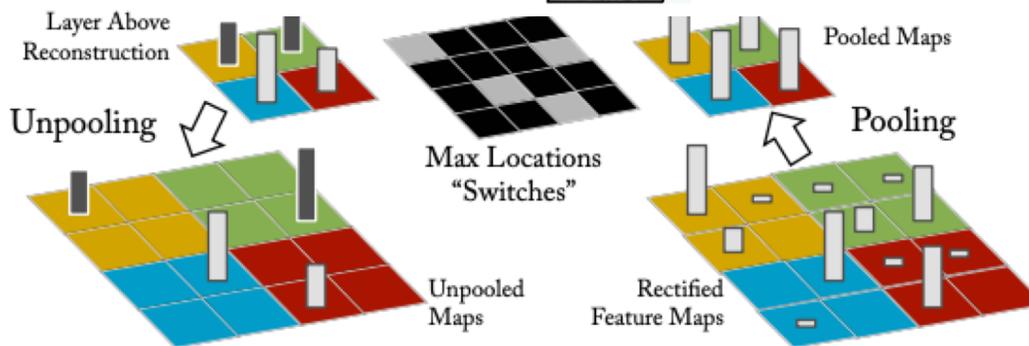
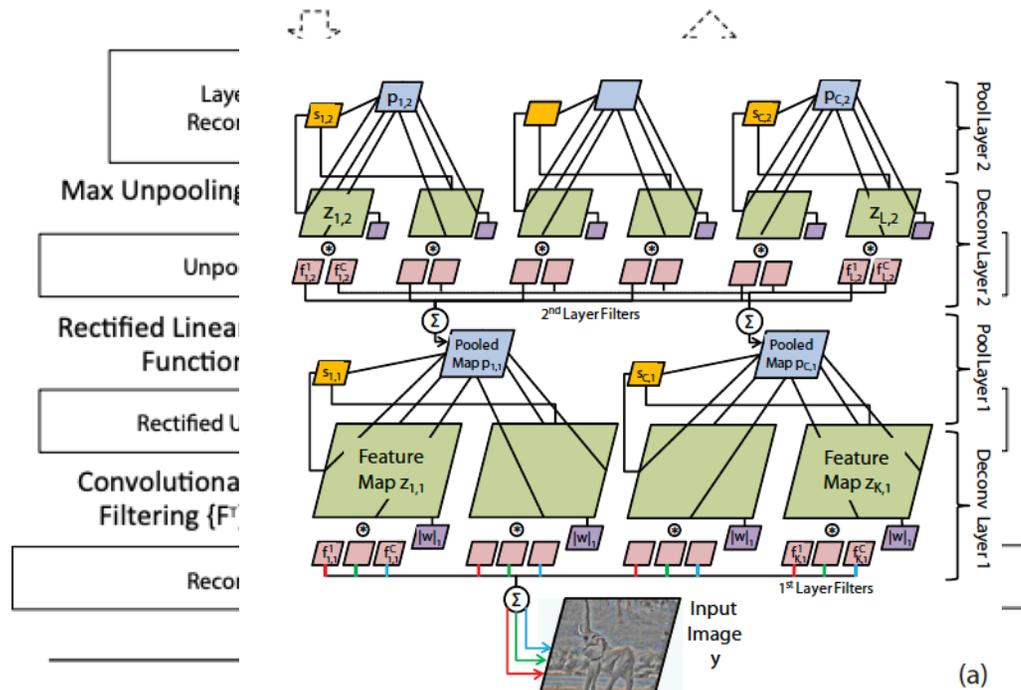
Convolution  
3x3 filter on 4x4 input  
2x2 output



Transposed Convolution  
3x3 filter on 2x2 input  
4x4 output

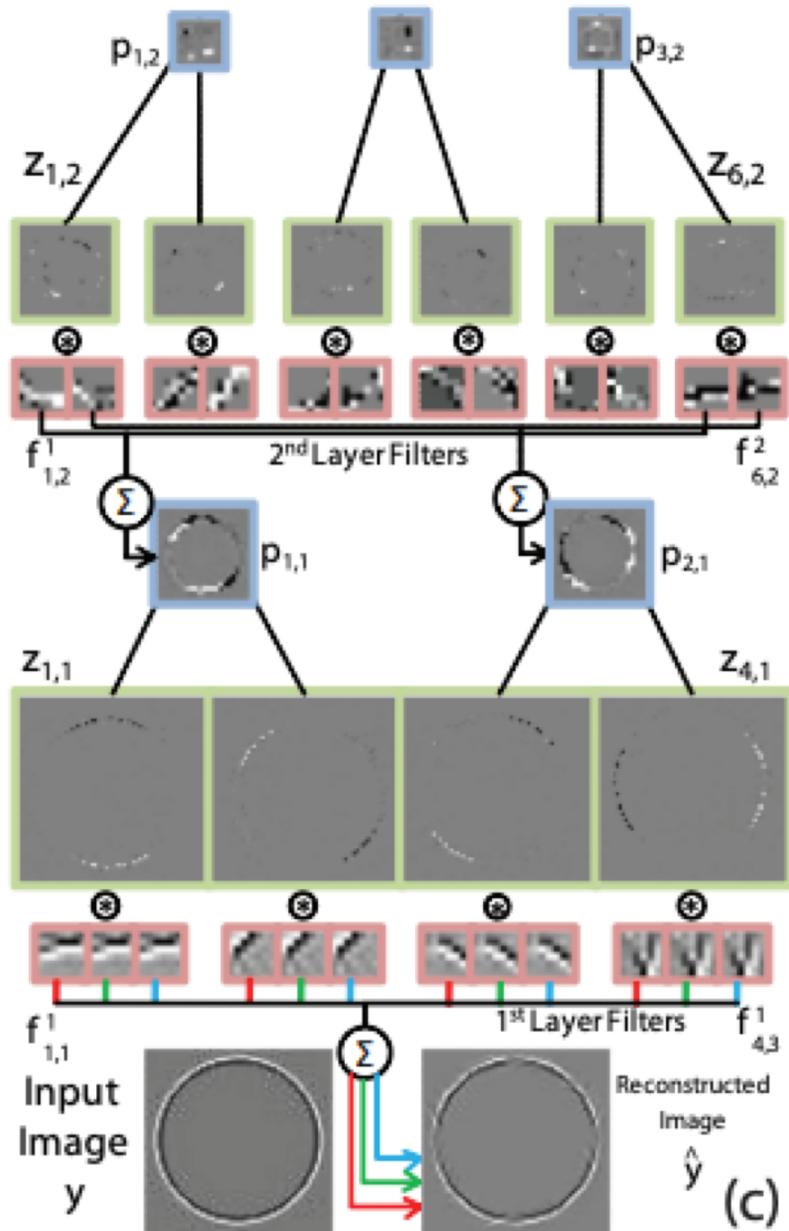
# Deconvolute node activations

Deconvolutional neural net: A novel way to map high level activities back to the input pixel space, showing what input pattern originally caused a given activation in the feature maps



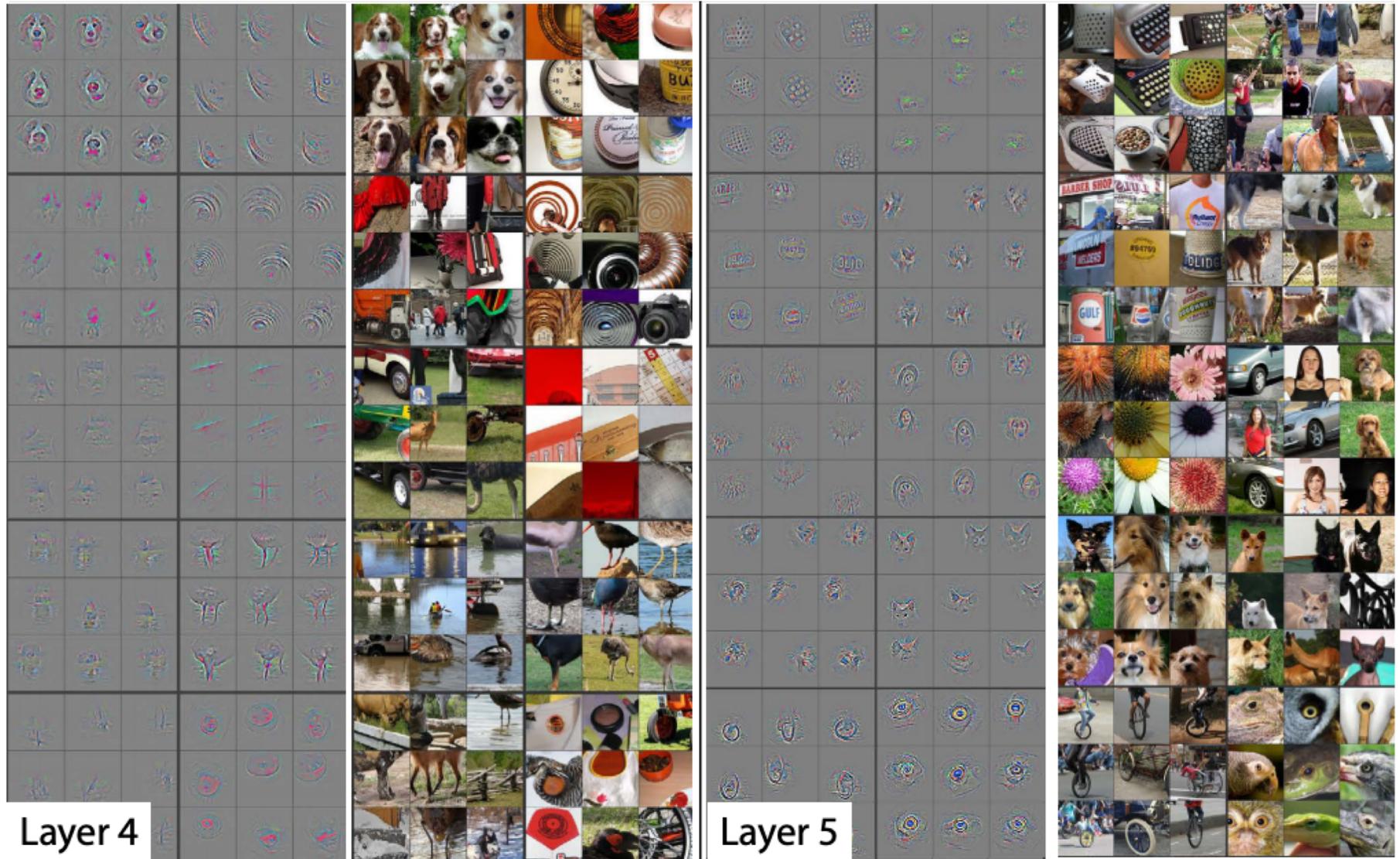


# Deconvolute node activations

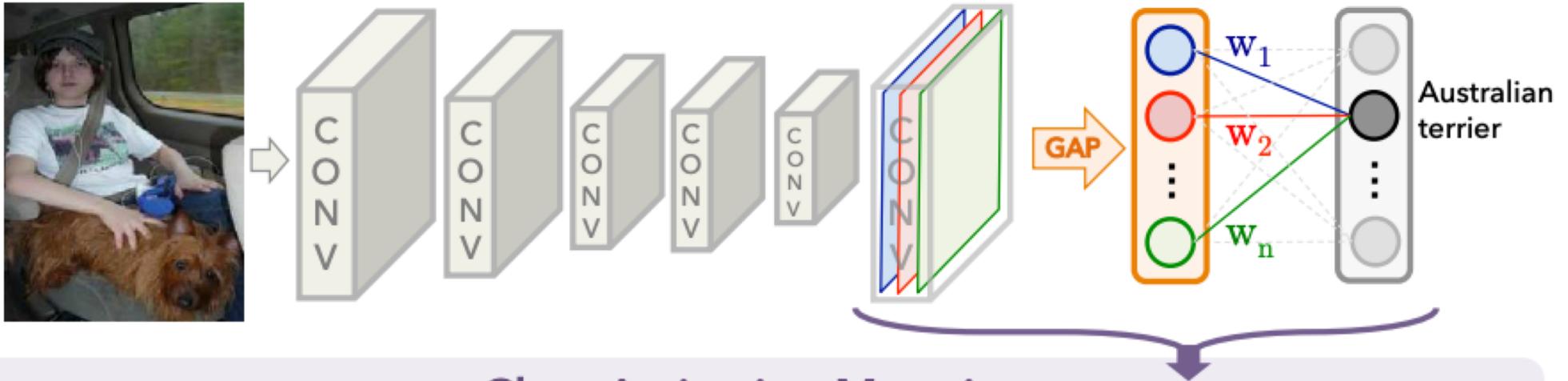


# Deconvolute node activations

Deconvolutional neural net: A novel way to map high level activities back to the input pixel space, showing what input pattern originally caused a given



# CAM: Class Activation Mapping



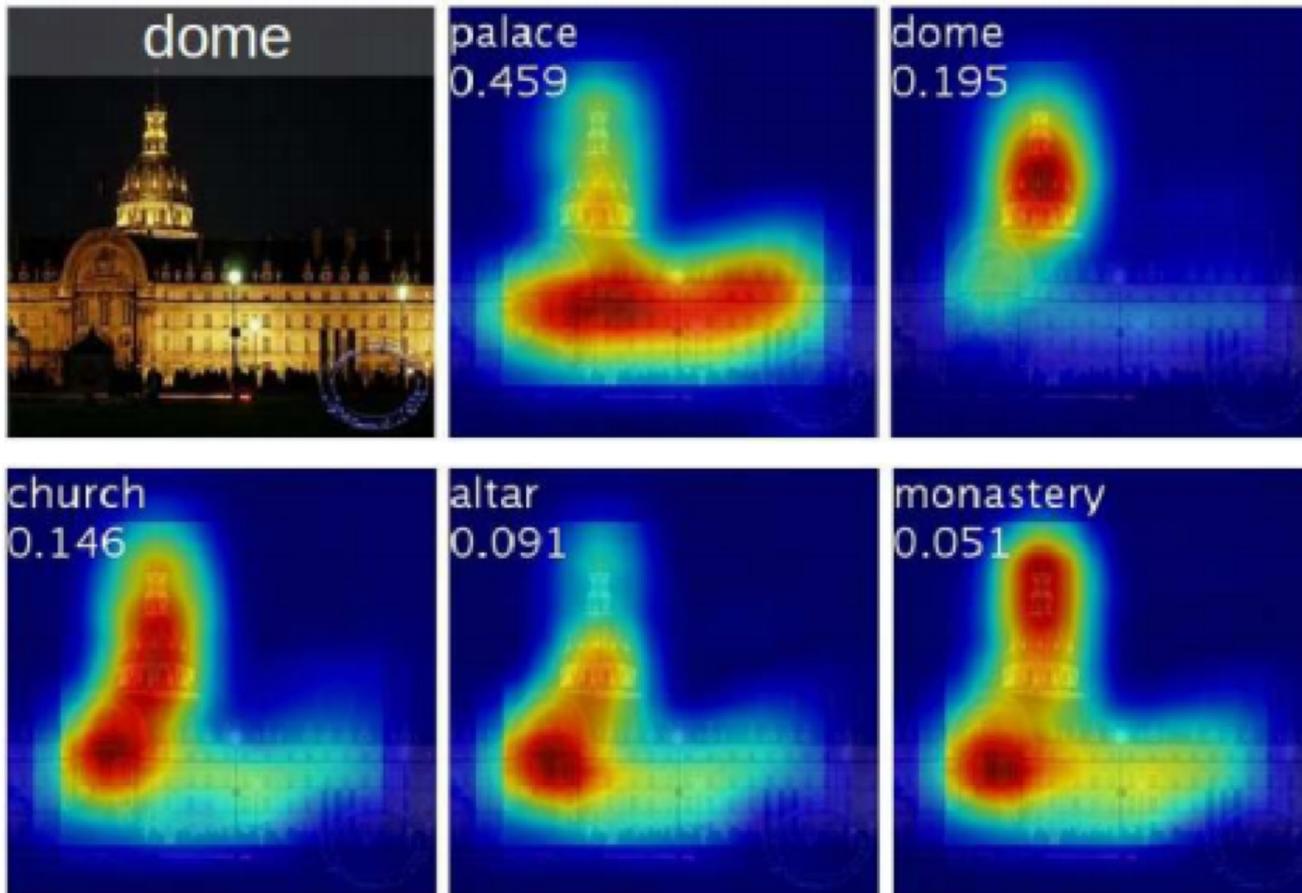
## Class Activation Mapping

$$W_1 * \text{[Heatmap 1]} + W_2 * \text{[Heatmap 2]} + \dots + W_n * \text{[Heatmap n]} = \text{[Final CAM Heatmap]}$$

Class Activation Map (Australian terrier)

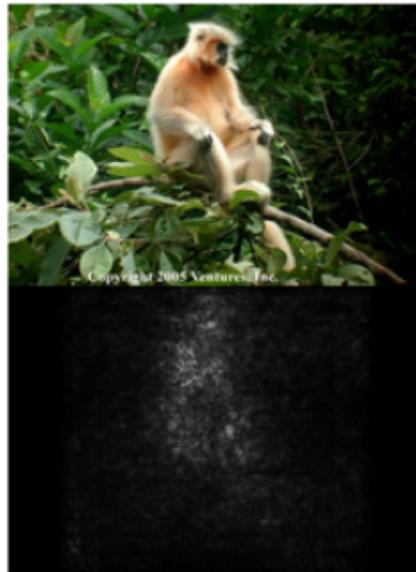
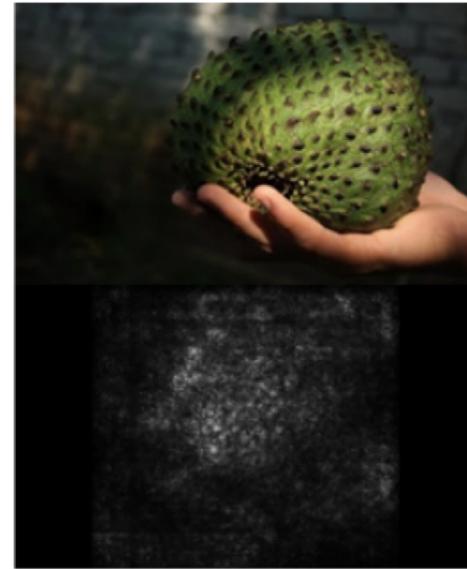
Use additional layer on top of the GAP (Global activation pooling) to learn **class specific** linear weights for each high level feature map and use them to weight the activations mapped back into input space.

# CAM: Class Activation Mapping



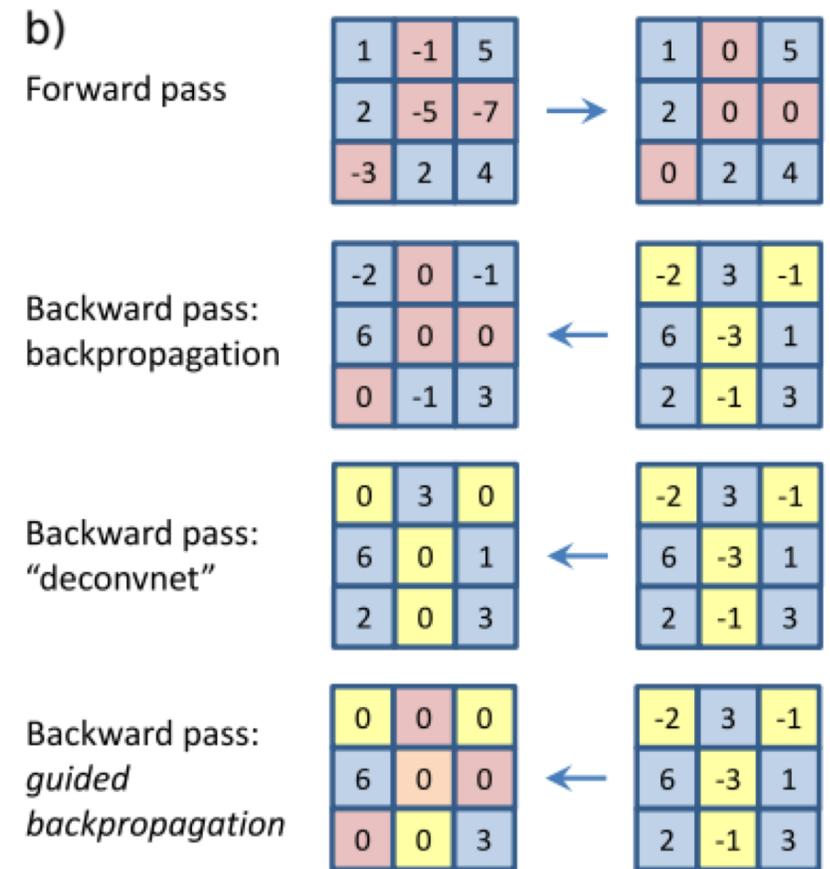
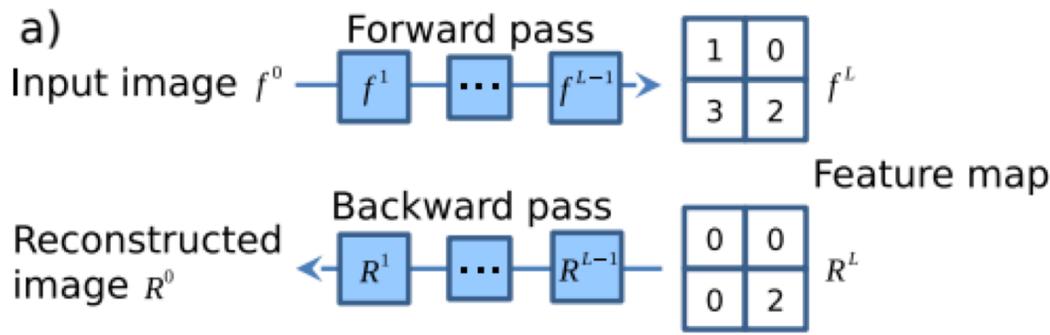
Use additional layer on top of the GAP (Global activation pooling) to learn **class specific** linear weights for each high level feature map and use them to weight the activations mapped back into input space.

# Visualizing gradient: Saliency map



# Gradient variation 1: guided back-propagation

Only back propagate positive gradients



c)

activation:  $f_i^{l+1} = \text{relu}(f_i^l) = \max(f_i^l, 0)$

backpropagation:  $R_i^l = (f_i^l > 0) \cdot R_i^{l+1}$ , where  $R_i^{l+1} = \frac{\partial f^{out}}{\partial f_i^{l+1}}$

backward 'deconvnet':  $R_i^l = (R_i^{l+1} > 0) \cdot R_i^{l+1}$

guided backpropagation:  $R_i^l = (f_i^l > 0) \cdot (R_i^{l+1} > 0) \cdot R_i^{l+1}$

# Gradient variation 1: guided back-propagation

Only back propagate positive gradients

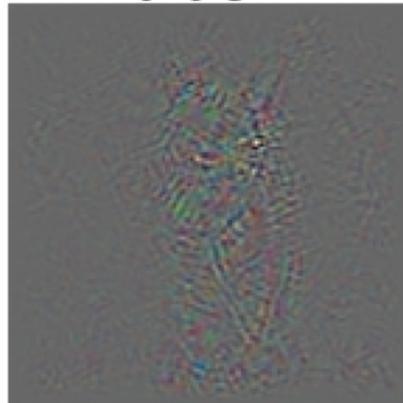
guided backpropagation



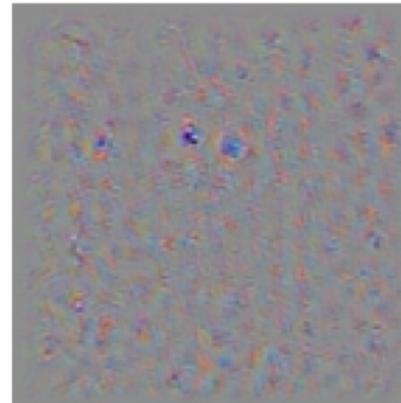
corresponding image crops



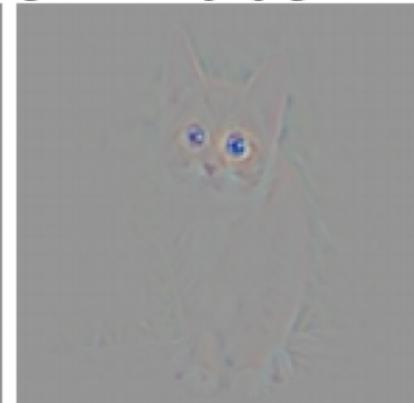
backpropagation



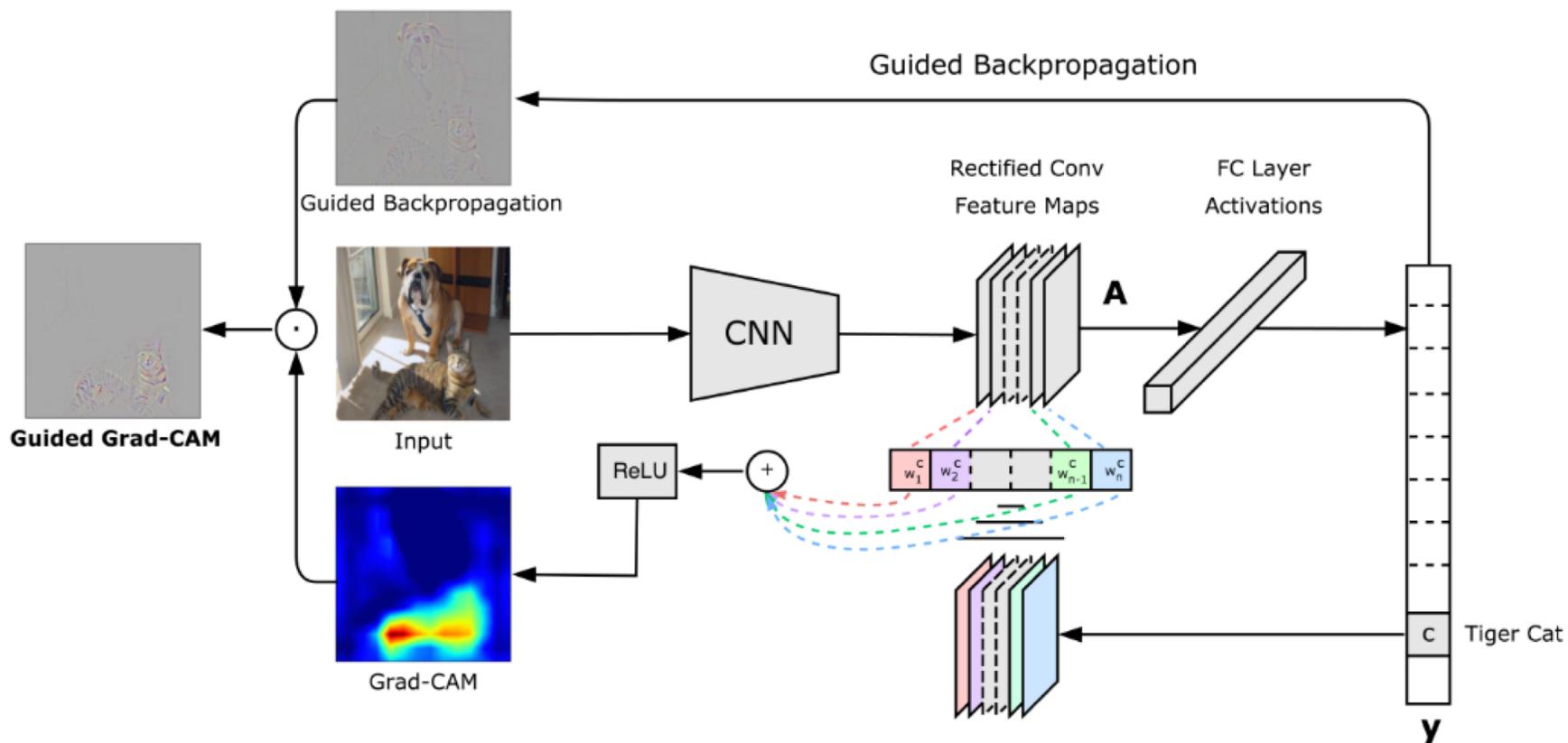
'deconvnet'



guided backpropagation



# Grad-CAM: combining CAM and guided backprop

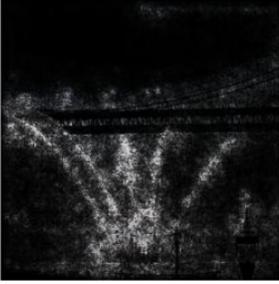
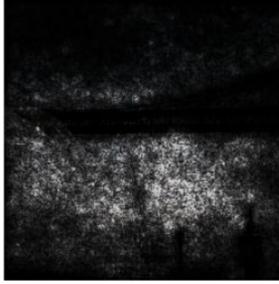


# Integrated Gradient

Given an input image  $x_i$  and a **baseline input**  $x'_i$  :

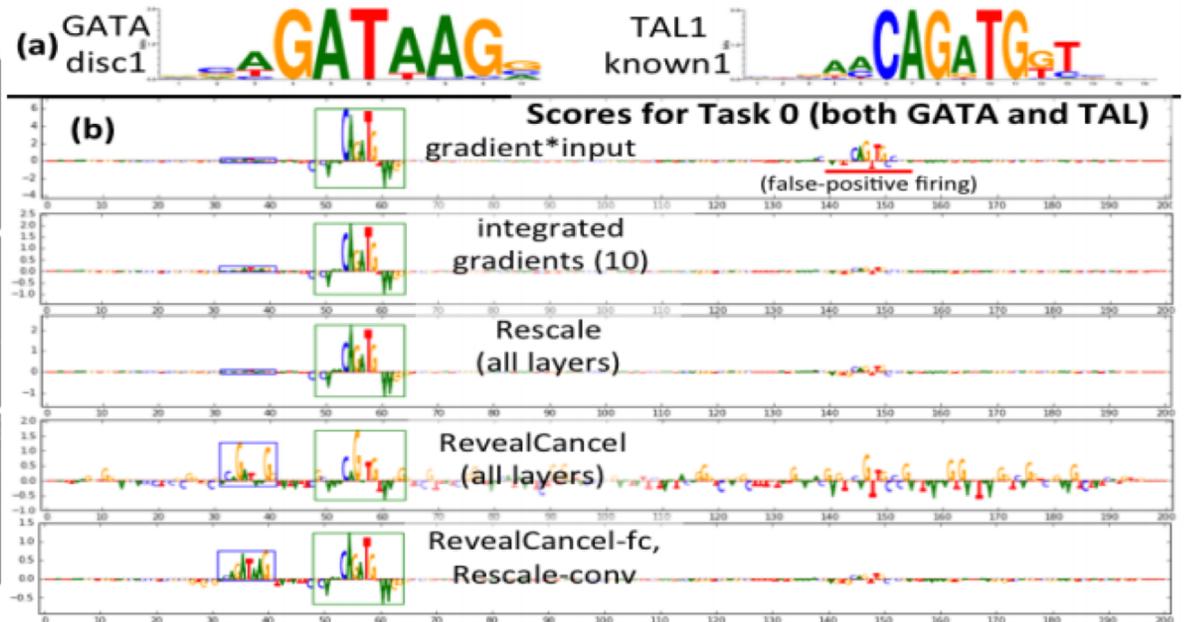
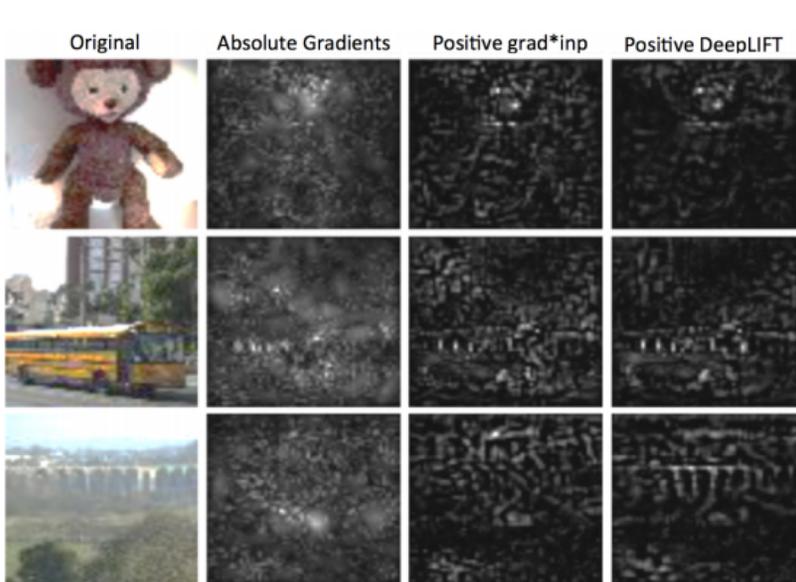
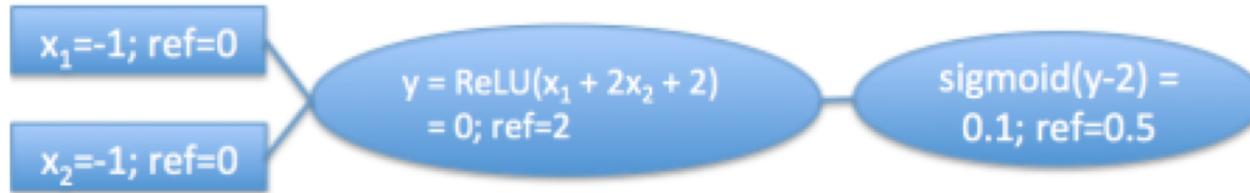
$$\text{IntegratedGrads}_i(x) ::= (x_i - x'_i) \times \int_{\alpha=0}^1 \frac{\partial F(x' + \alpha \times (x - x'))}{\partial x_i} d\alpha$$

$$\text{IntegratedGrads}_i^{\text{approx}}(x) ::= (x_i - x'_i) \times \sum_{k=1}^m \frac{\partial F(x' + \frac{k}{m} \times (x - x'))}{\partial x_i} \times \frac{1}{m}$$

Original image	Top label and score	Integrated gradients	Gradients at image
	Top label: reflex camera Score: 0.993755		
	Top label: fireboat Score: 0.999961		
	Top label: school bus Score: 0.997033		

# DeepLIFT

compares the activation of each neuron to its reference activation and assigns contribution scores according to the difference



# Other input dependent attribution score approaches:

- LIME (Local Interpretable Model-agnostic Explanations)
  - identify an interpretable model over the interpretable representation that is locally faithful to the classifier by approximating the original function with interpretable models locally.
- SHAP(SHapley Additive explanation)
  - Unified several additive attribution score methods by using definition of Sharpley value from game theory
- maximum entropy
  - Locally sample inputs that maximum the entropy of predicted score

# Input independent visualization: gradient ascent

Generate input that maximum activation of certain neuron or final activation of the class

$$\arg \max_I S_c(I) - \lambda \|I\|_2^2$$

Simple regularizer: Penalize L2 norm of generated image



**dumbbell**



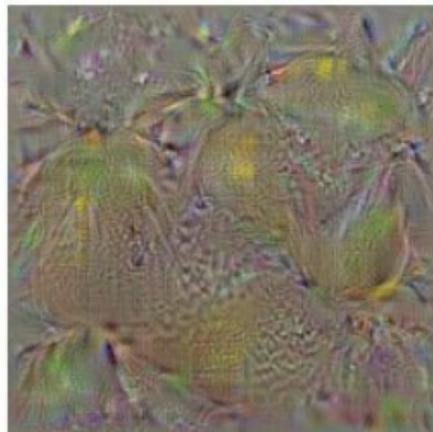
**cup**



**dalmatian**



**bell pepper**



**lemon**



**husky**

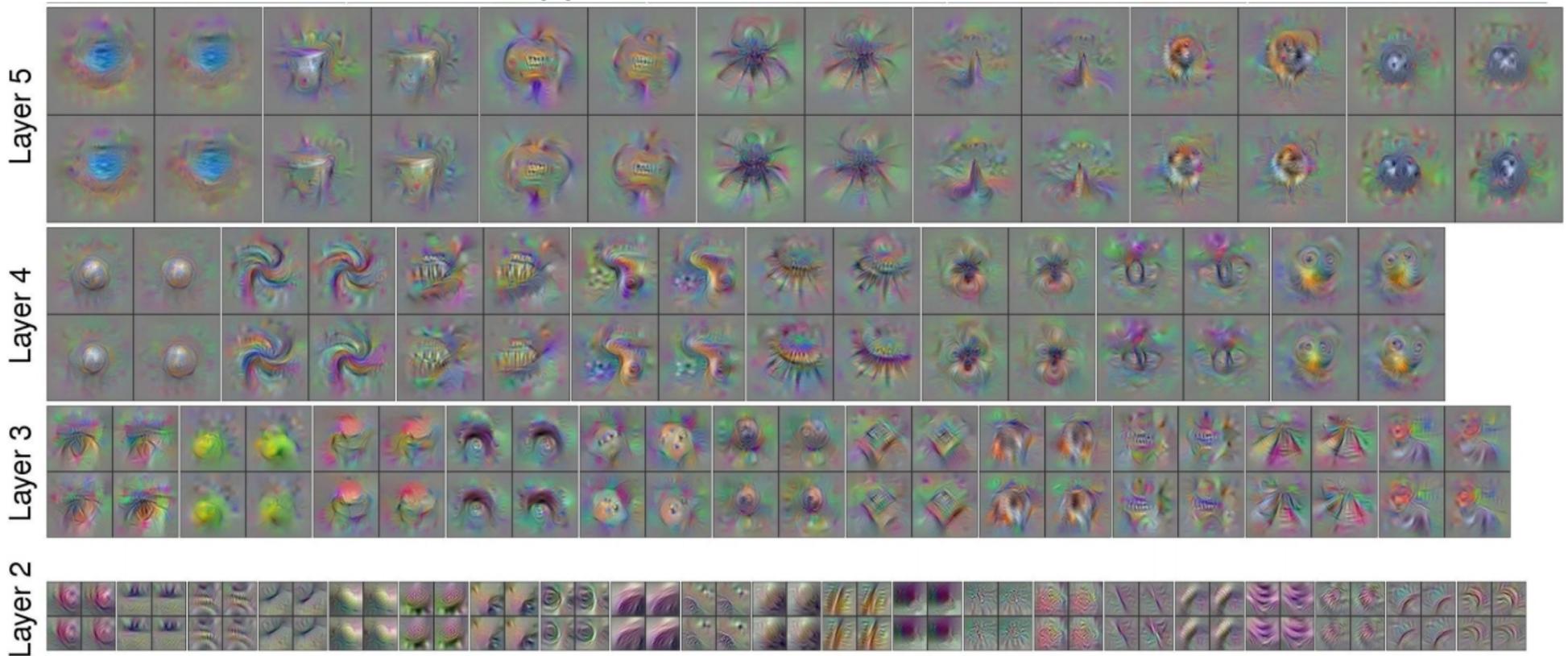
Simonyan et al., *Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps*

# Input independent visualization: gradient ascent

Generate input that maximum activation of certain neuron or final activation of the class

$$\arg \max_I S_c(I) - \lambda \|I\|_2^2$$

Simple regularizer: Penalize L2 norm of generated image



# DeepMotif uses gradient ascent

## NFYB

JASPAR Motifs	Forward: 	Backward: 
CNN Positive Class Maximization		
RNN Positive Class Maximization		

Lanchantin et al., *Deep Motif Dashboard: Visualizing and Understanding Genomic Sequences Using Deep Neural Networks*

**FIN - Thank You**

# SIS Resources

**Full paper in arXiv:**

**<https://arxiv.org/abs/1810.03805>**

**Code for paper and analysis:**

**<https://github.com/b-carter/SufficientInputSubsets>**

**Code for open-source SIS library and tutorial:**

**[https://github.com/google-research/google-research/tree/master/sufficient\\_input\\_subsets](https://github.com/google-research/google-research/tree/master/sufficient_input_subsets)**