

CSCI 1470/2470
Spring 2024

Ritambhara Singh

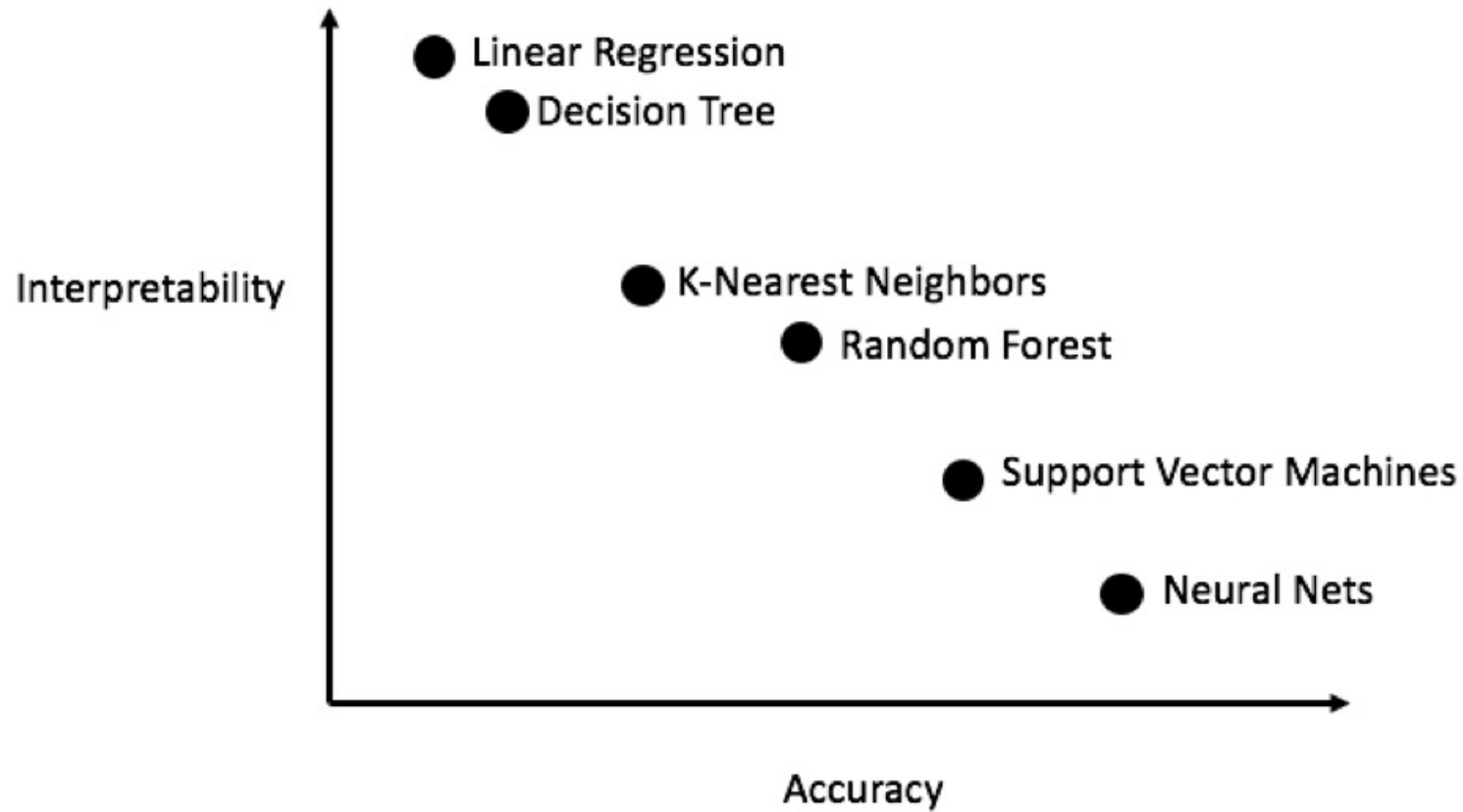
March 18, 2024
Monday

Interpretability

Deep Learning

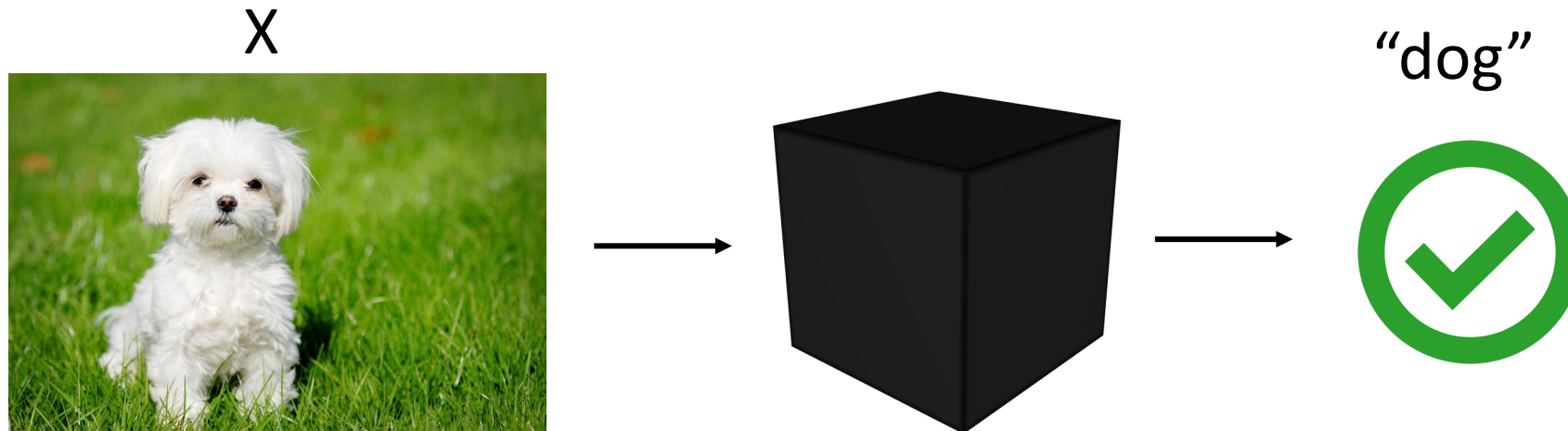


Machine Learning and Interpretability

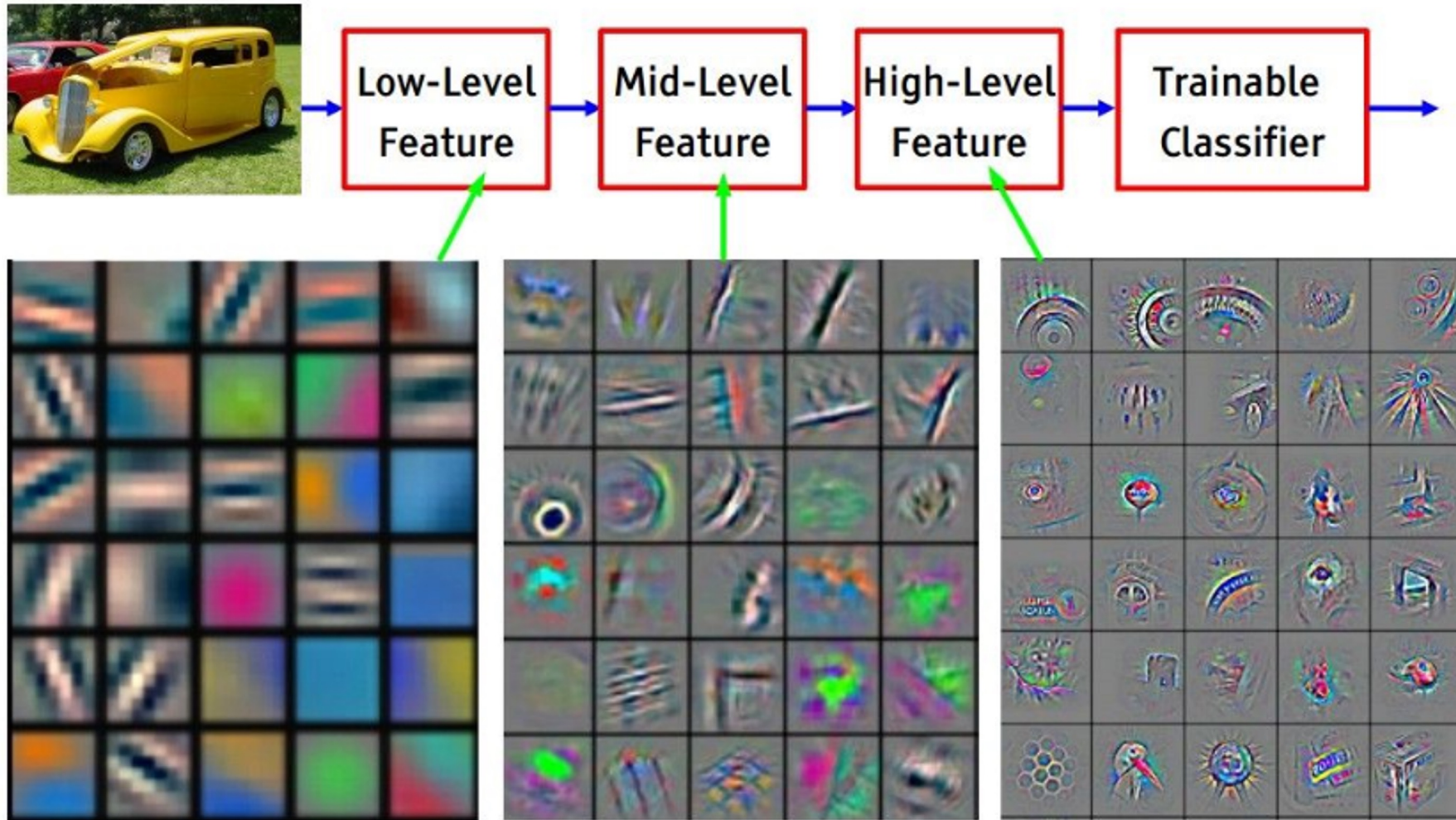


Deep Nets are “Black Box” Models

- What do the hidden layers of networks actually learn?
 - **Image recognition:** What do the many thousands of filters actually represent??
 - **Natural language processing:** What does the RNN hidden state actually store??



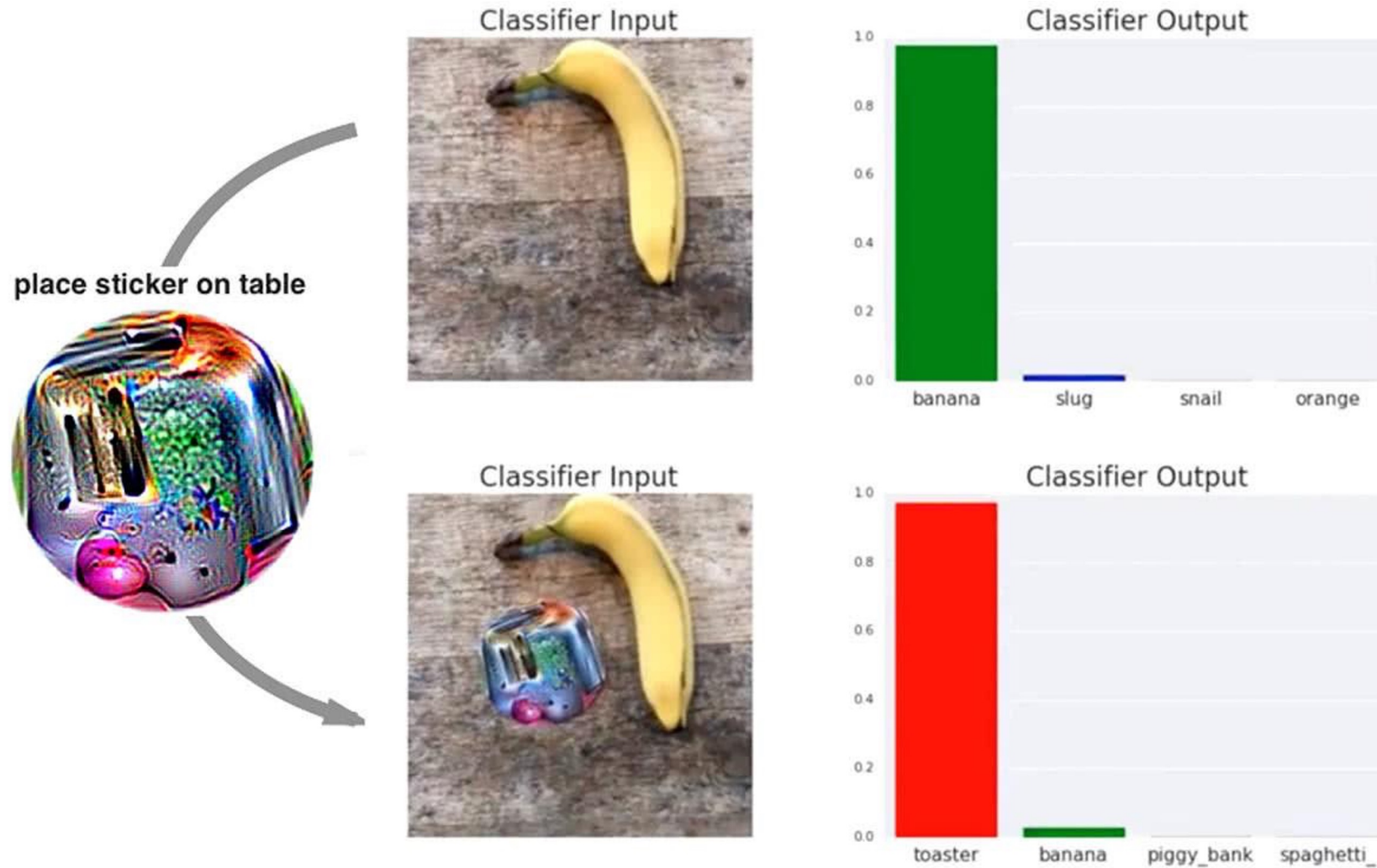
Example: What do CNN filters look like?



Example: What do CNN filters look like?



...which leads to situations like this:

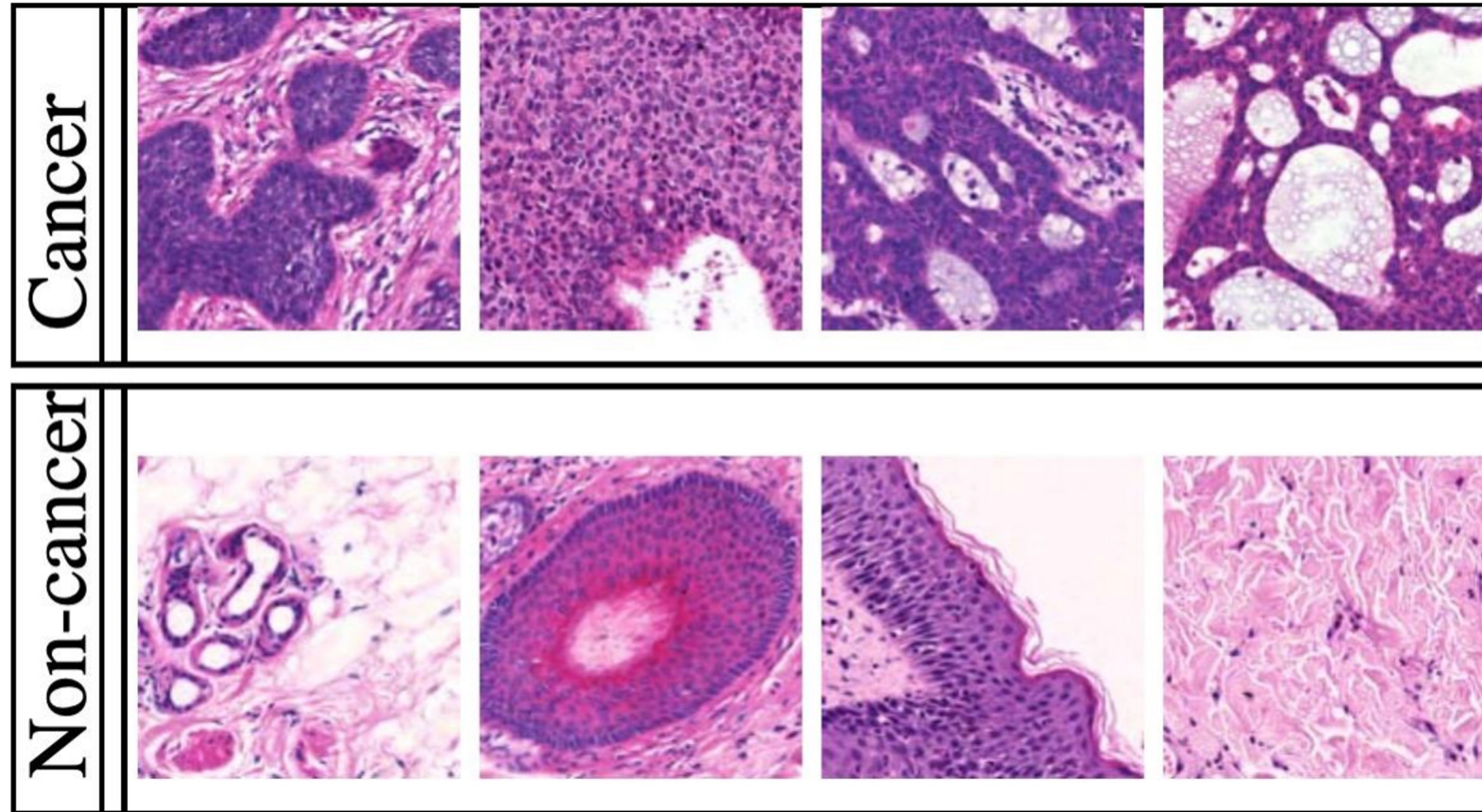


Deep Learning has an
interpretability problem

Maybe not such a big deal if we're
just classifying breakfast foods...

...but what if the decision is
really important?

What if the decision is *really* important?



From [Cruz-Roa et al., 2013](#)

What if the decision is *really* important?

- How can a human (e.g. a doctor) trust that a network is making a sensible decision (e.g. a patient has the flu, and *not* strep, mono, or pneumonia)?

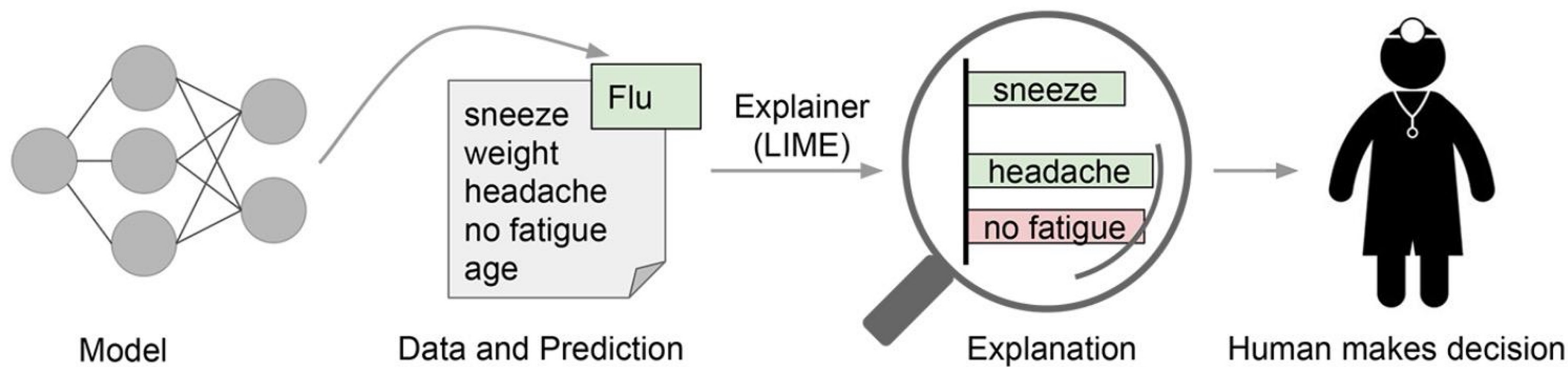


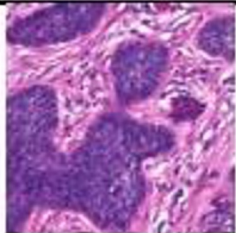
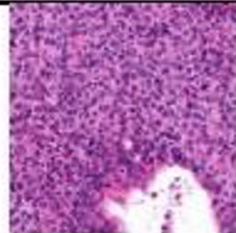
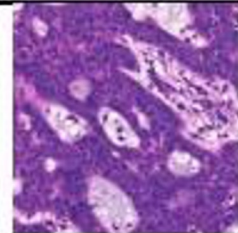
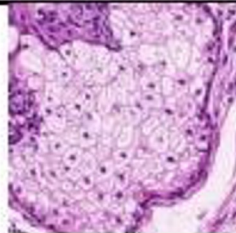
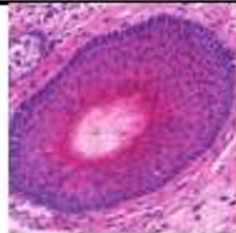
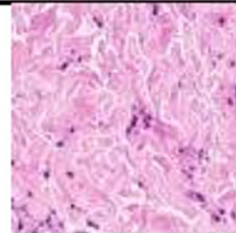
Figure 1. Explaining individual predictions to a human decision-maker. Source: Marco Tulio Ribeiro.

Model Interpretability

- Broadly, *interpretability* refers to ways of understanding/measuring how a model made a decision
- Can take the form of visualizations, summary statistics, metrics, ...
- A whole subfield of study: *Interpretable DL/AI*

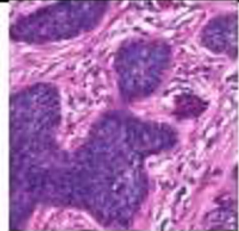
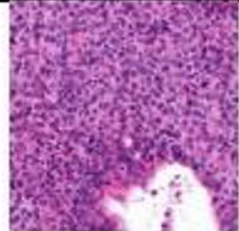
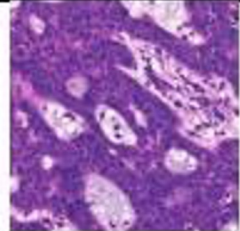
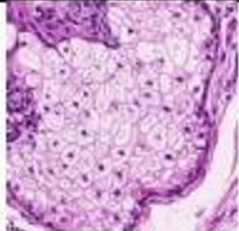
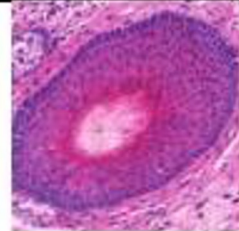
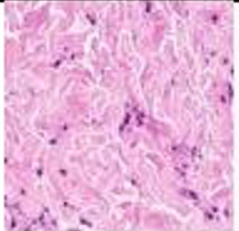
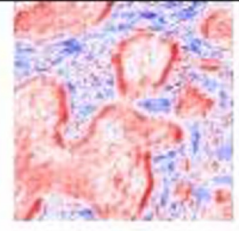
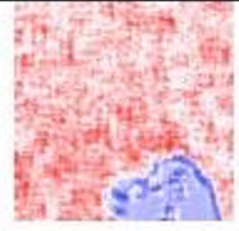
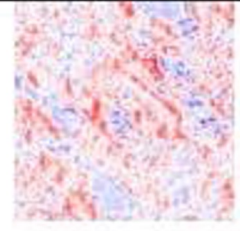
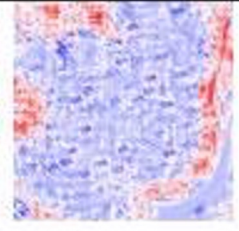
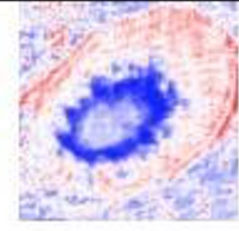
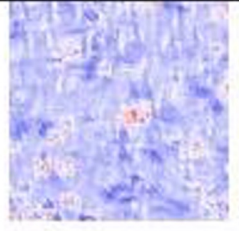
Making Cancer Predicting Interpretable

Simultaneously learn to predict which regions of the tissue are cancerous
(so a person can look and see if that makes sense)

True class	Cancer	Cancer	Cancer	Non-cancer	Non-cancer	Non-cancer
Input image						
Pred/Prob	Cancer (0.82)	Cancer (0.96)	Cancer (0.79)	Non-cancer (0.27)	Non-cancer (0.08)	Non-cancer (0.03)

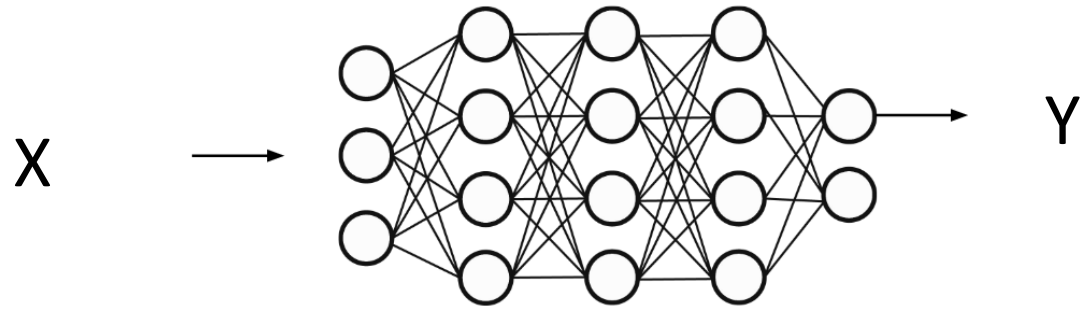
Making Cancer Predicting Interpretable

Simultaneously learn to predict which regions of the tissue are cancerous
(so a person can look and see if that makes sense)

True class	Cancer	Cancer	Cancer	Non-cancer	Non-cancer	Non-cancer
Input image						
Pred/Prob	Cancer (0.82)	Cancer (0.96)	Cancer (0.79)	Non-cancer (0.27)	Non-cancer (0.08)	Non-cancer (0.03)
Digital staining						

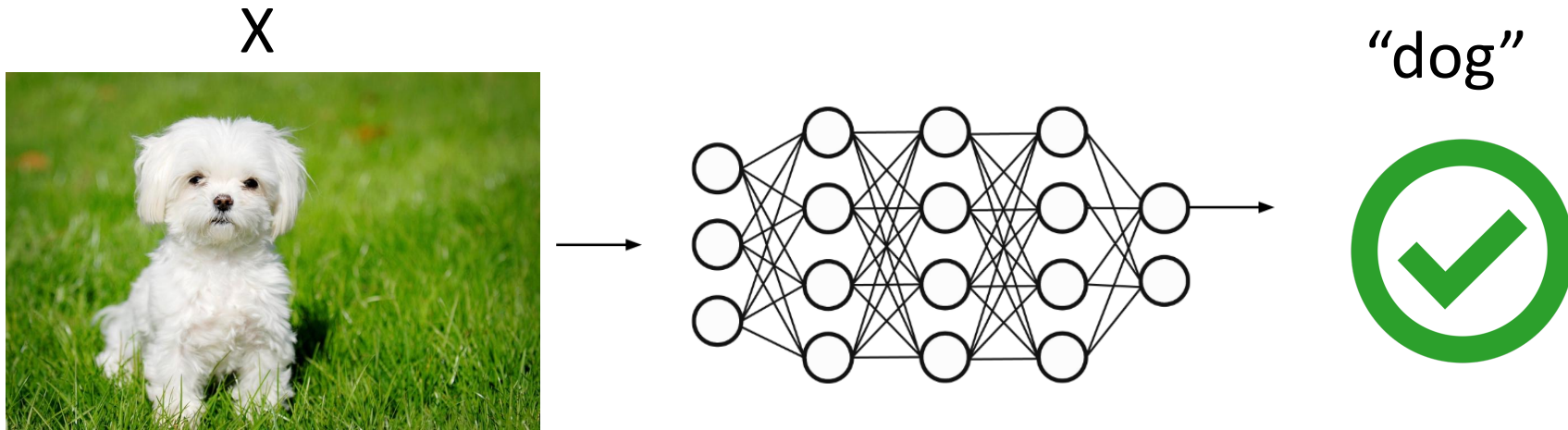
Red = cancerous tissue
Blue = healthy tissue

Task formulation



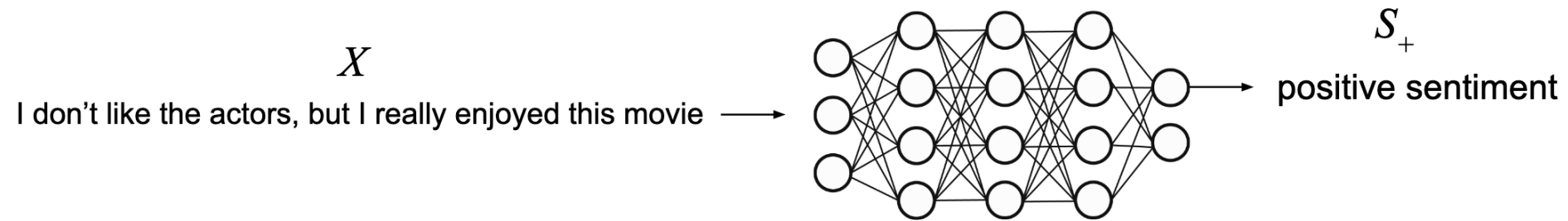
Which features in X are most important for the model prediction Y ?

Task formulation



Which pixels are most important for classification?

Task formulation



Which words/characters are most important for classification?

Today's goal – learn about interpretation in DL

- (1) Model architecture based methods (CNNs and RNNs)
- (2) Gradient-based methods
- (3) Model agnostic methods

Today's goal – learn about interpretation in DL

(1) Model architecture based methods (CNNs and RNNs)

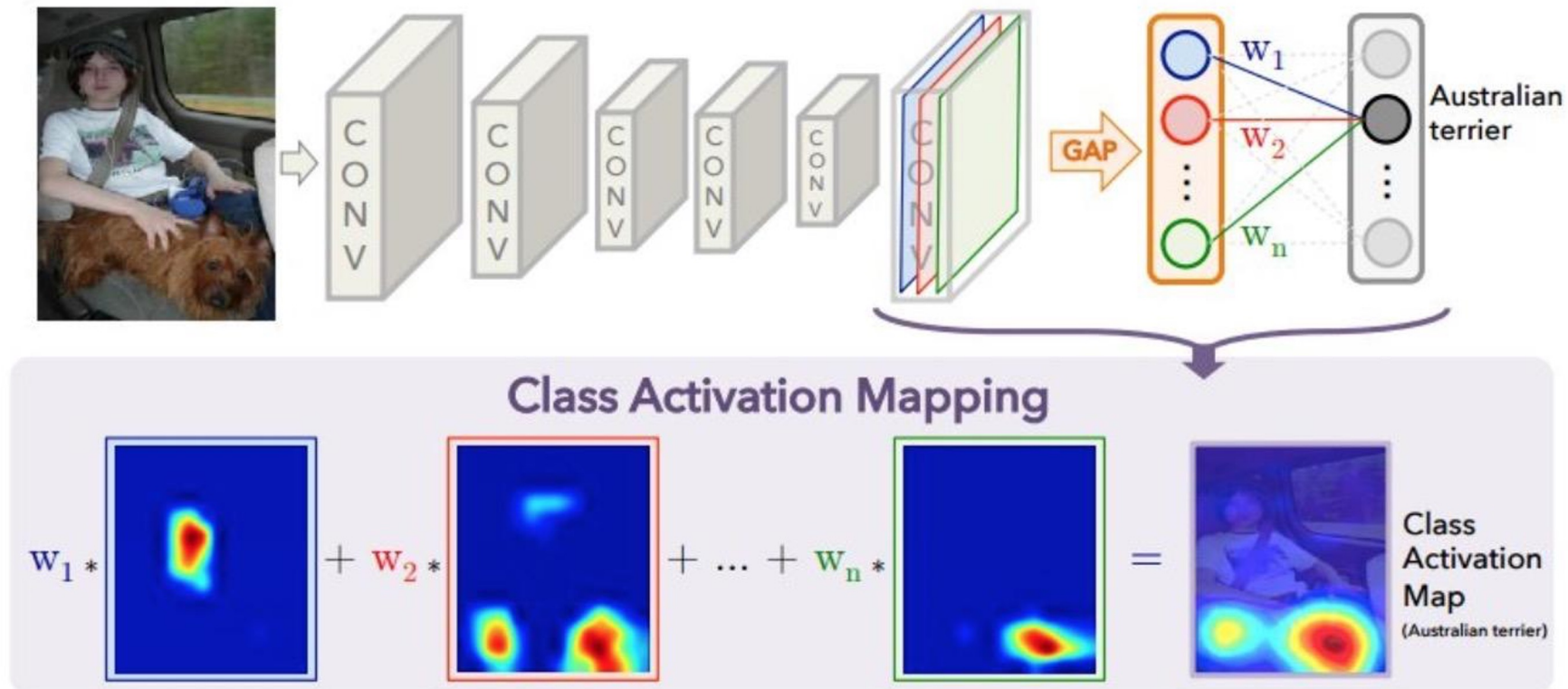
(2) Gradient-based methods

(3) Model agnostic methods

Identifying image regions that influence classification result

Global Average Pooling: Average all the pixels in the last feature map to produce a flat vector, then feed that through a linear layer to produce class logits

- A weighted sum of the last feature maps, according to the weights of the linear layer, localizes the region that leads to the classification

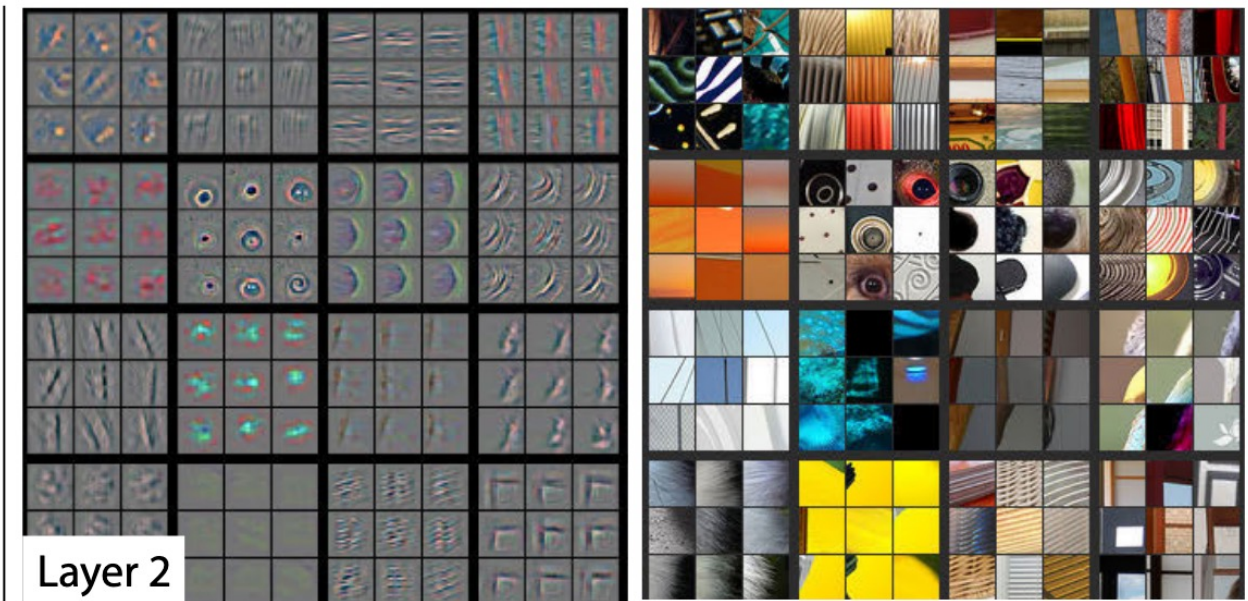
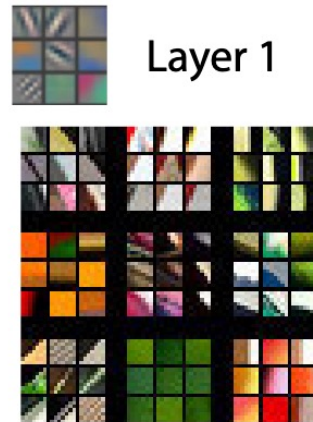
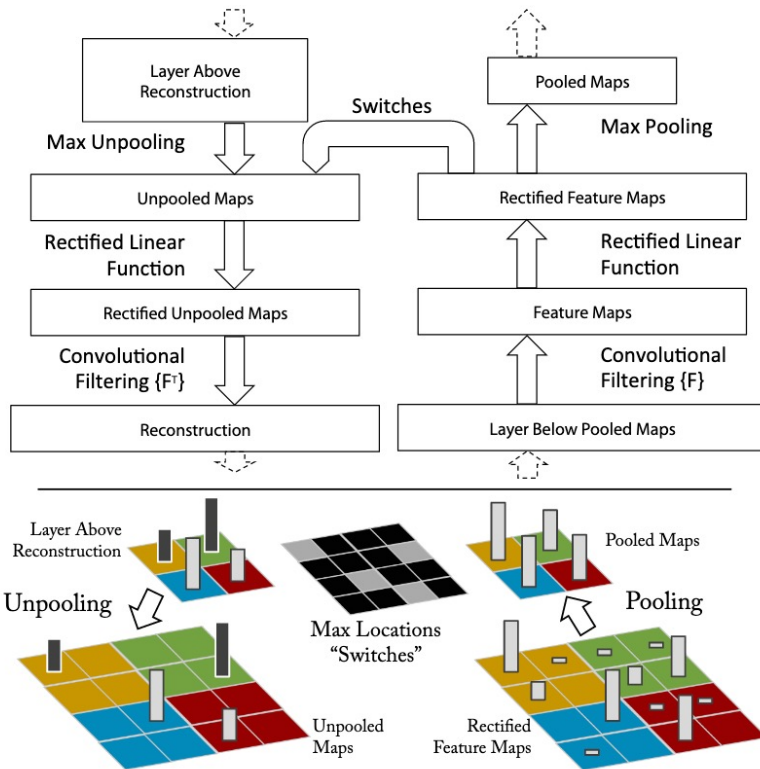


Deconvolution

Map filter activations back to the input pixel space, showing what input pattern originally caused a given activation in the feature maps.

Perform this mapping with a [Deconvolutional Network](#) (deconvnet).

Decovnet: a convnet model that uses the same components (filtering, pooling) but in reverse, so instead of mapping pixels to features it does the opposite.



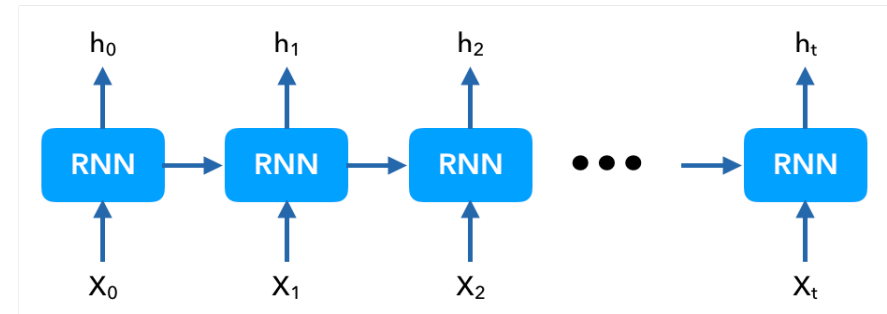
Interpreting RNNs

Pick one entry (cell) of the hidden state, highlight characters that cause that cell to take on a high value

- This is a ***character-level language model***, not a word-level one

Cell sensitive to position in line:

The sole importance of the crossing of the Berezina lies in the fact that it plainly and indubitably proved the fallacy of all the plans for cutting off the enemy's retreat and the soundness of the only possible line of action--the one Kutuzov and the general mass of the army demanded--namely, simply to follow the enemy up. The French crowd fled at a continually increasing speed and all its energy was directed to reaching its goal. It fled like a wounded animal and it was impossible to block its path. This was shown not so much by the arrangements it made for crossing as by what took place at the bridges. When the bridges broke down, unarmed soldiers, people from Moscow and women with children who were with the French transport, all--carried on by vis inertiae--pressed forward into boats and into the ice-covered water and did not, surrender.



Interpreting RNNs

Pick one entry (cell) of the hidden state, highlight characters that cause that cell to take on a high value

- This is a ***character-level language model***, not a word-level one

Cell that turns on inside quotes:

"You mean to imply that I have nothing to eat out of.... On the contrary, I can supply you with everything even if you want to give dinner parties," warmly replied Chichagov, who tried by every word he spoke to prove his own rectitude and therefore imagined Kutuzov to be animated by the same desire.

Kutuzov, shrugging his shoulders, replied with his subtle penetrating smile: "I meant merely to say what I said."

Interpreting RNNs

Pick one entry (cell) of the hidden state, highlight characters that cause that cell to take on a high value

- This is a ***character-level language model***, not a word-level one

Cell that robustly activates inside if statements:

```
static int __dequeue_signal(struct sigpending *pending, sigset_t *mask,
                           siginfo_t *info)
{
    int sig = next_signal(pending, mask);
    if (sig) {
        if (current->notifier) {
            if (sigismember(current->notifier_mask, sig)) {
                if (!(current->notifier)(current->notifier_data)) {
                    clear_thread_flag(TIF_SIGPENDING);
                    return 0;
                }
            }
        }
        collect_signal(sig, pending, info);
    }
    return sig;
}
```

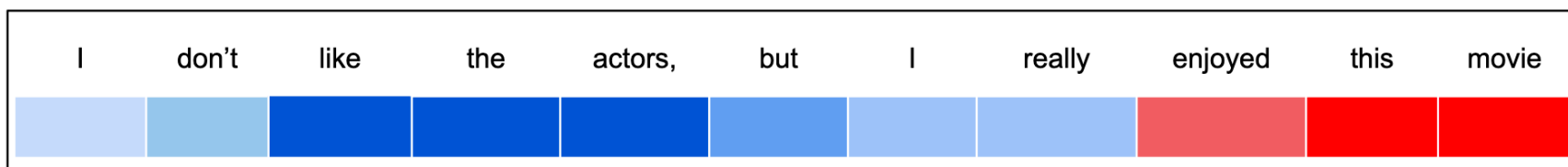
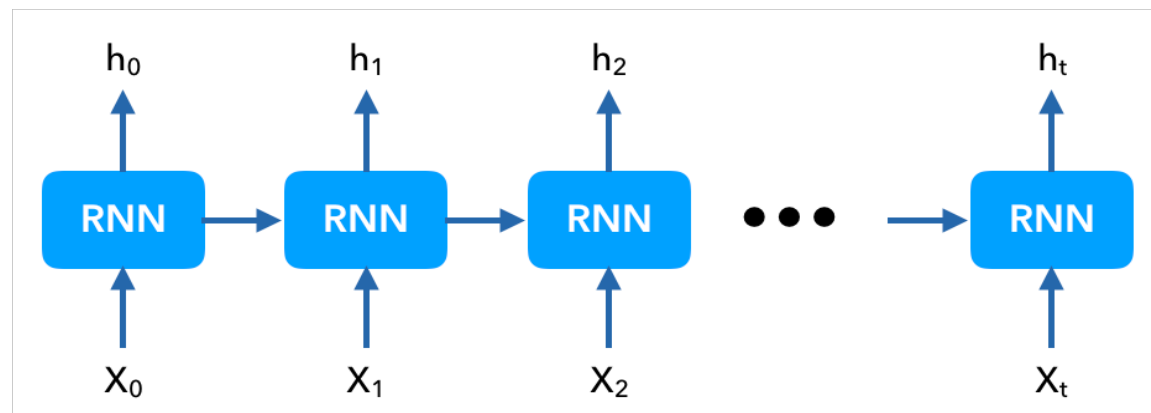
Note that in most cases, very few cells (~5%) respond to interpretable patterns like this

Any questions?



Temporal output (RNNs)

Track the prediction of the RNN for one hidden unit at a time



 = negative sentiment

 = positive sentiment

Today's goal – learn about interpretation in DL

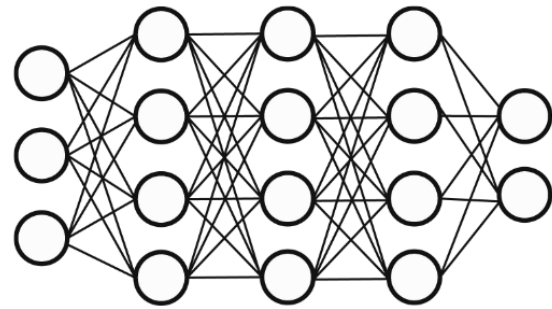
(1) Model architecture based methods

(2) Gradient-based methods

(3) Model agnostic methods

Saliency maps

X

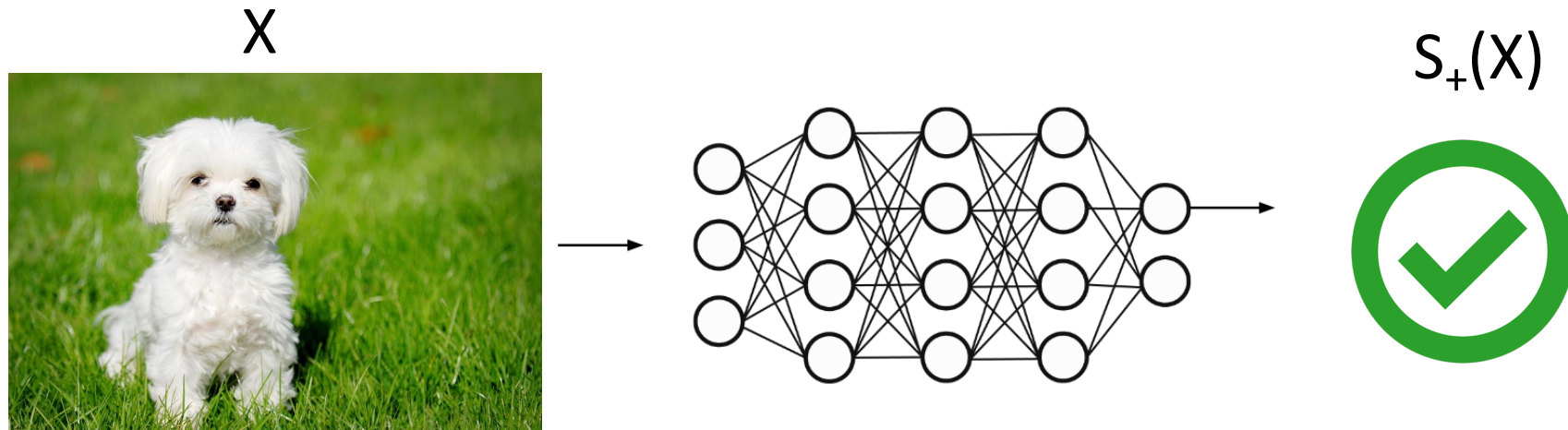


“dog”



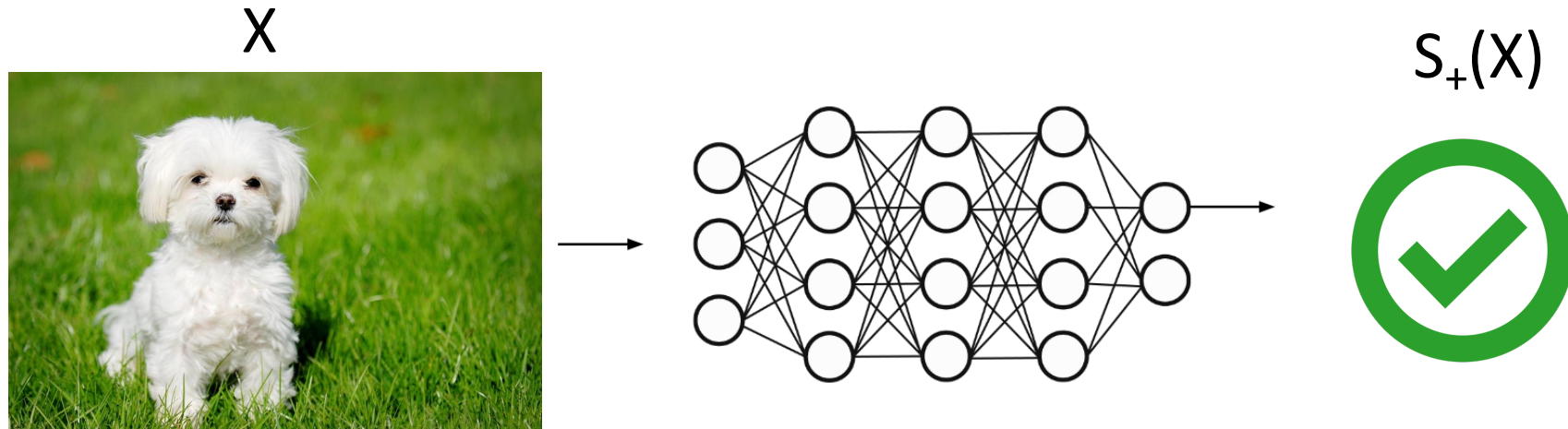
Which pixels are most important for classification?

Saliency maps



$$S_+(X) \approx w^T X + b = \sum_{i=1}^{|X|} w_i x_i$$

Saliency maps



$$S_+(X) \approx w^T X + b = \sum_{i=1}^{|X|} w_i x_i$$

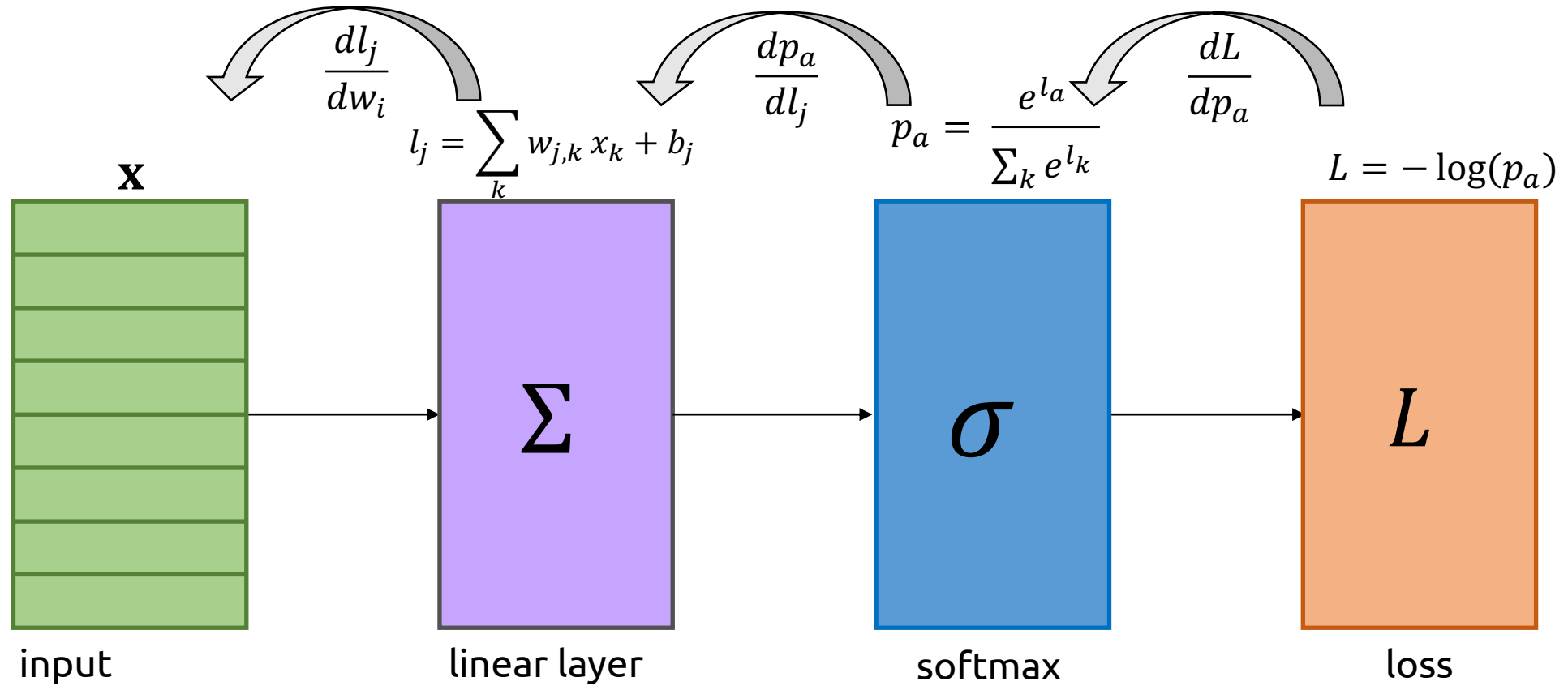
$$w = \left. \frac{\partial S_+}{\partial X} \right|_{X_0} = \text{“saliency map”}$$

How can we calculate this gradient?

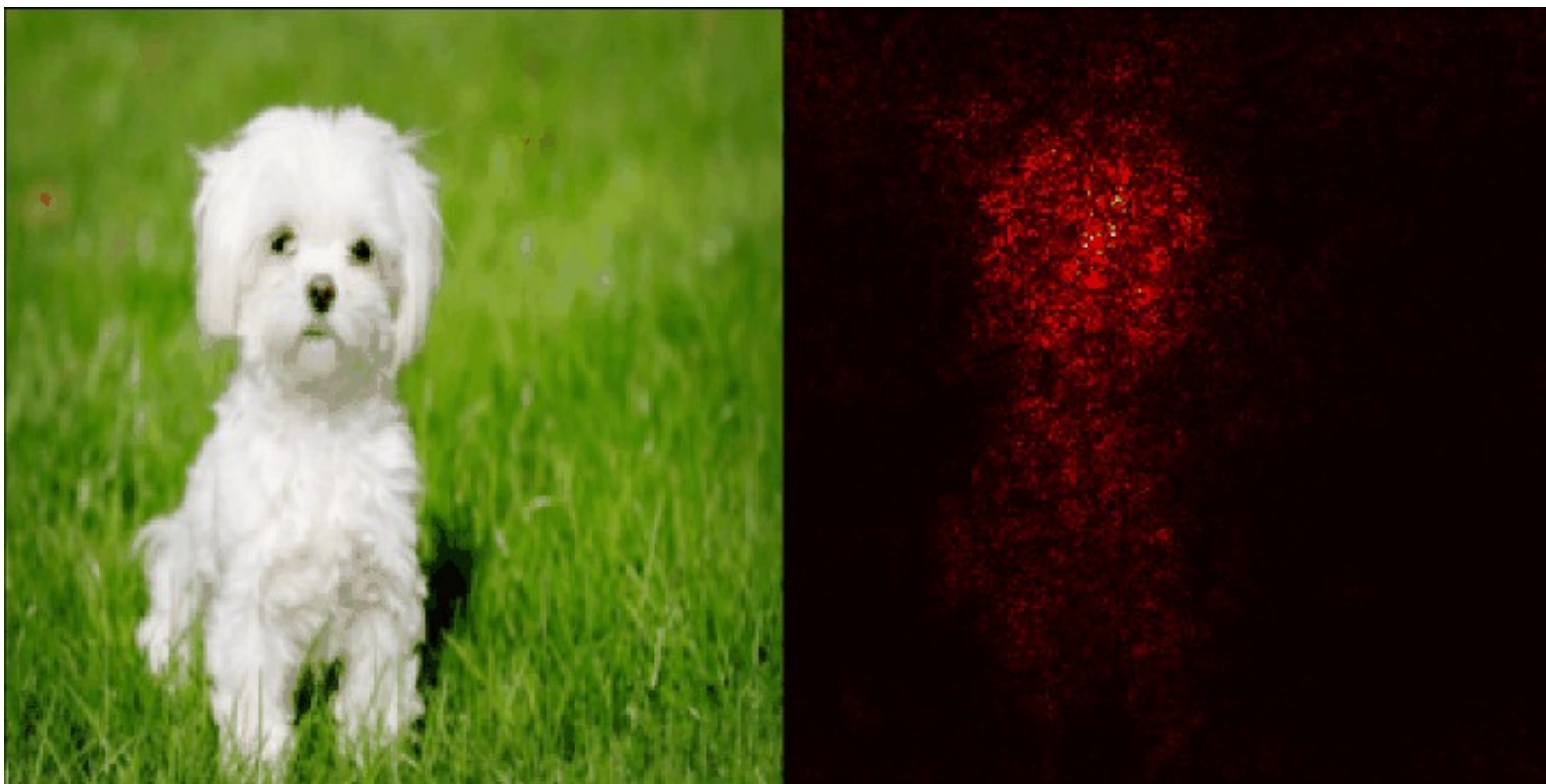
Backpropagation is back!

How do we change this?

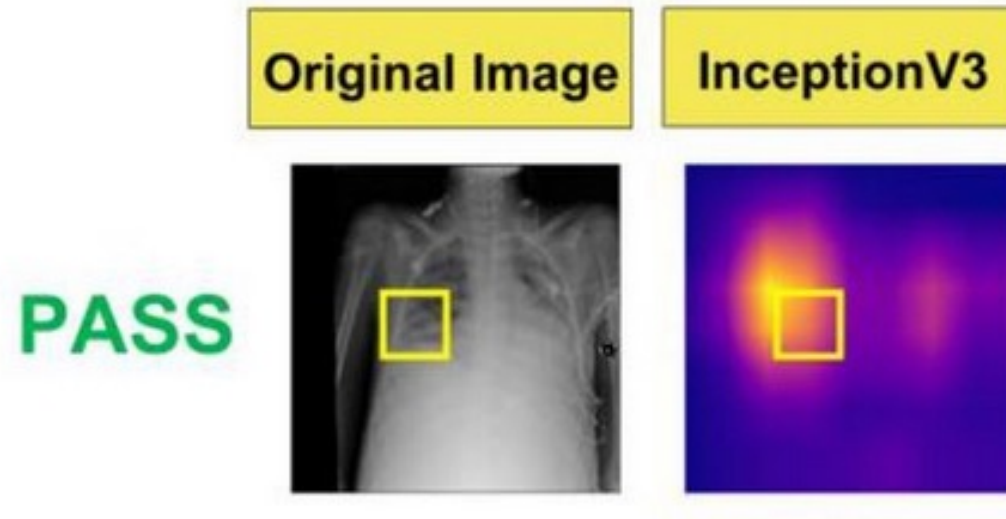
$$\Delta w_{j,i} = -\alpha \frac{\partial L}{\partial w_{j,i}} = -\alpha \cdot \frac{\partial L}{\partial p_a} \cdot \frac{\partial p_a}{\partial l_j} \cdot \frac{\partial l_j}{\partial w_{j,i}}$$



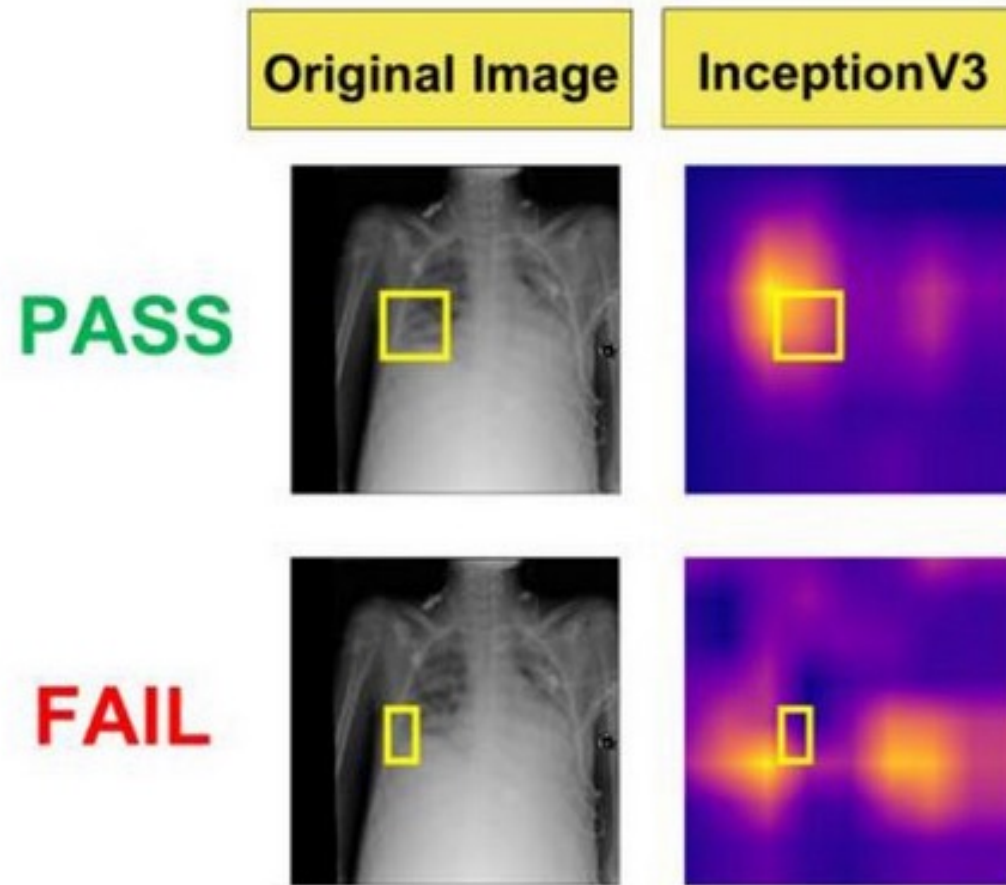
Saliency maps work well



Saliency maps can also fail



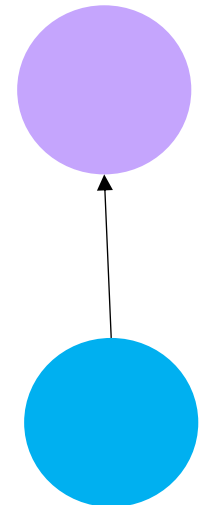
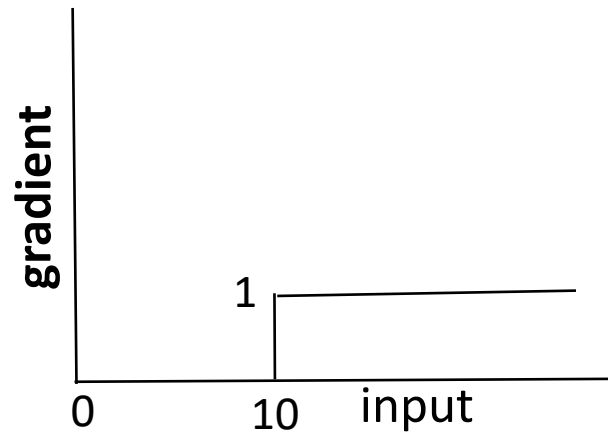
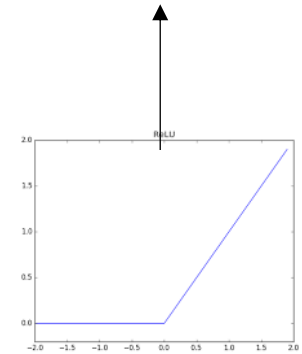
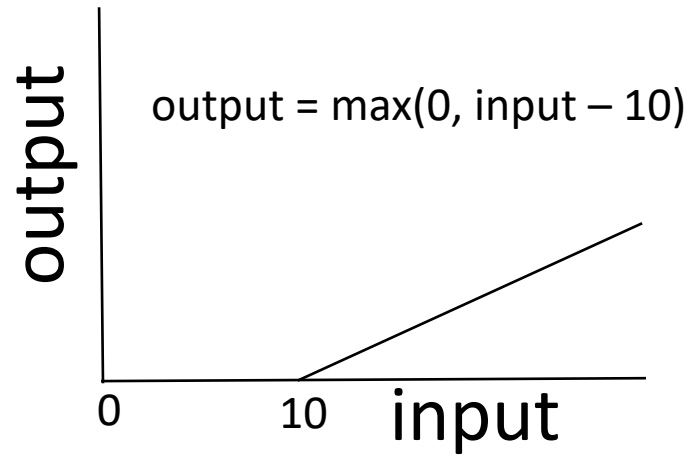
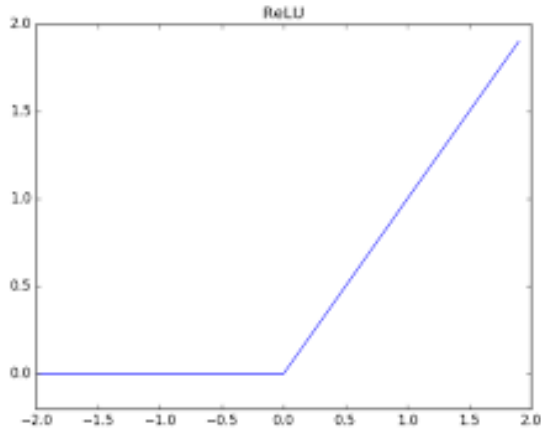
Saliency maps can also fail



What could be going wrong?

Backpropagation through activation functions

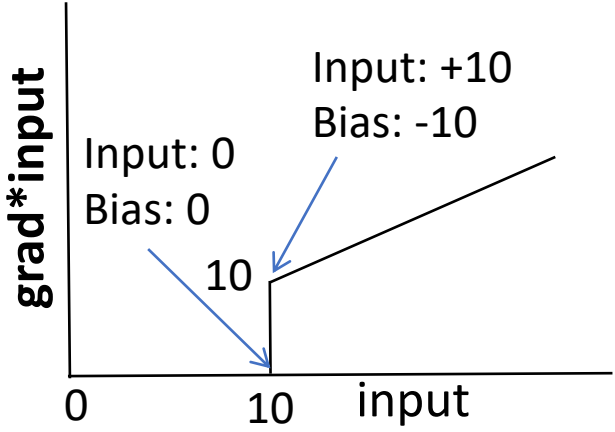
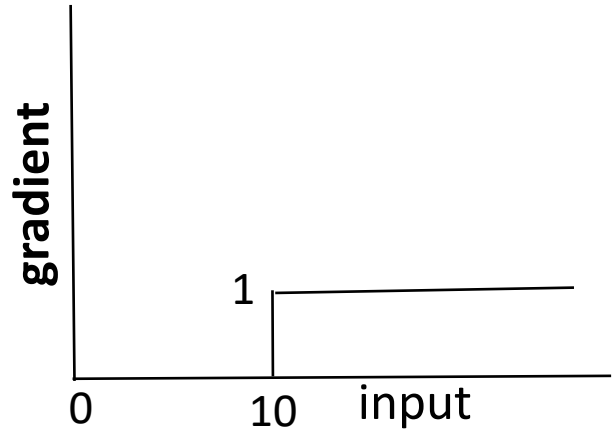
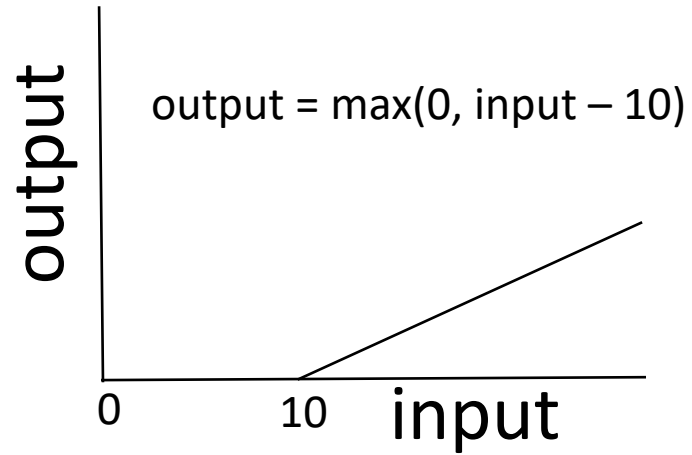
Thresholding problem



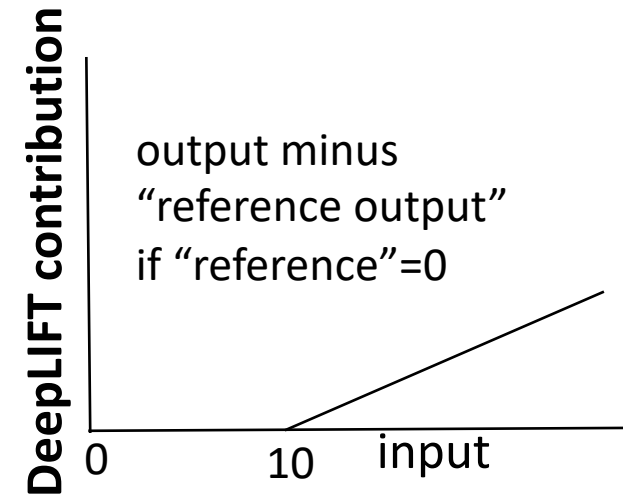
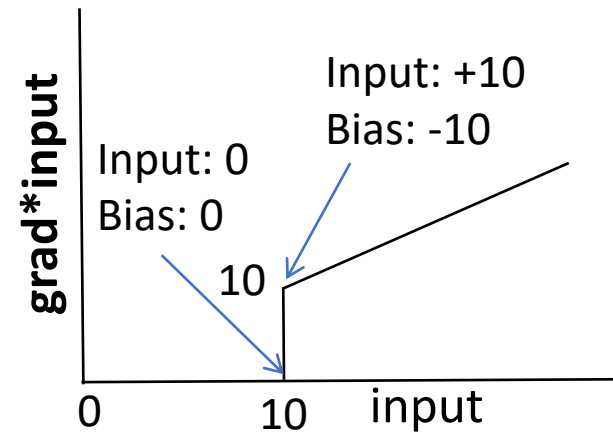
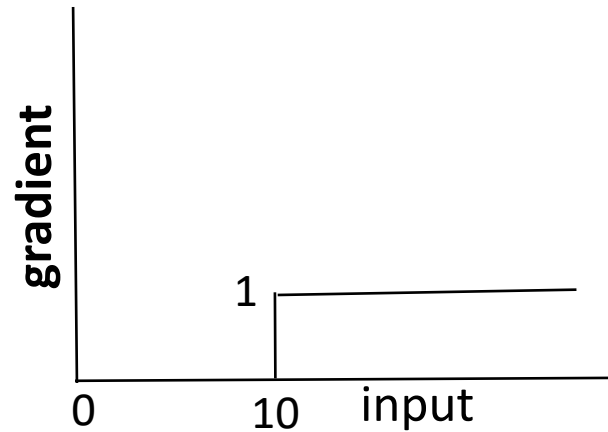
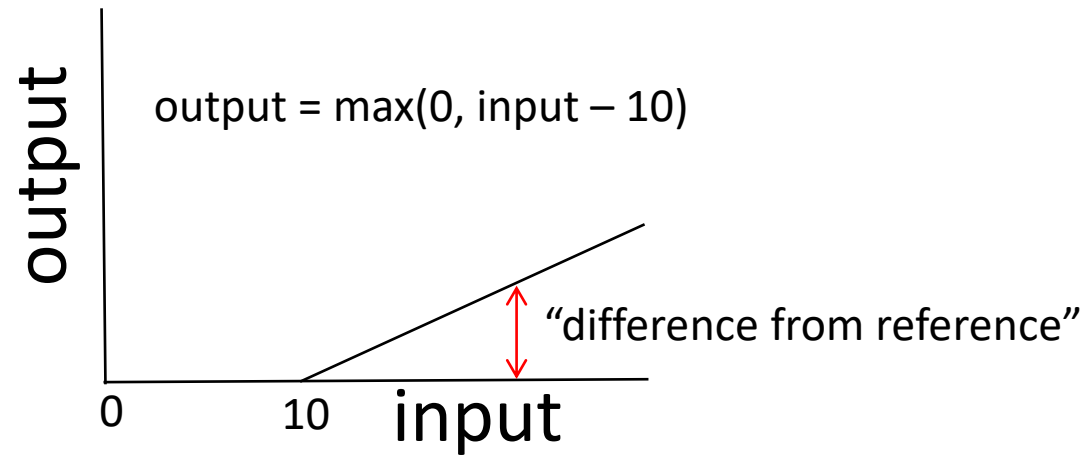
Any questions?



Gradient * Input



Solves thresholding problem



Integrated gradients

Input image



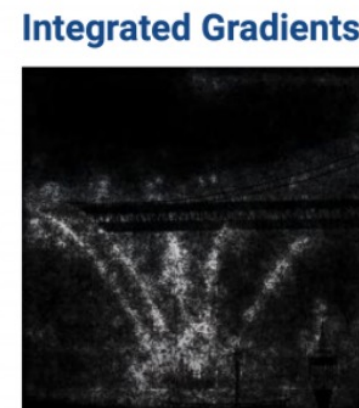
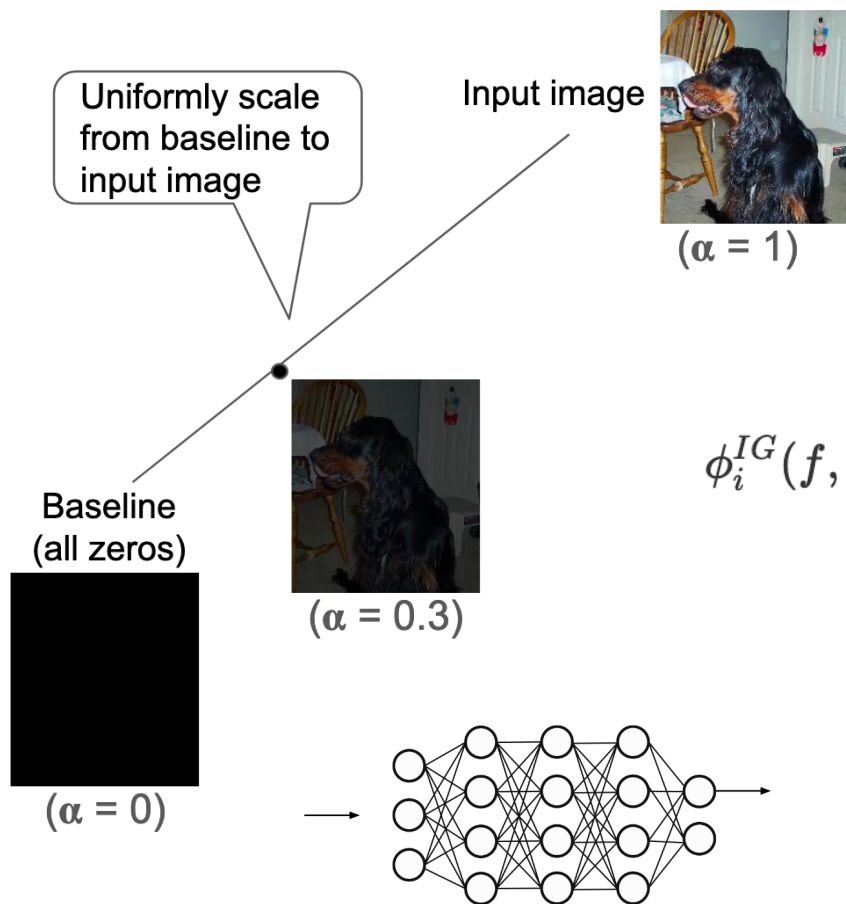
($\alpha = 1$)

Baseline
(all zeros)



($\alpha = 0$)

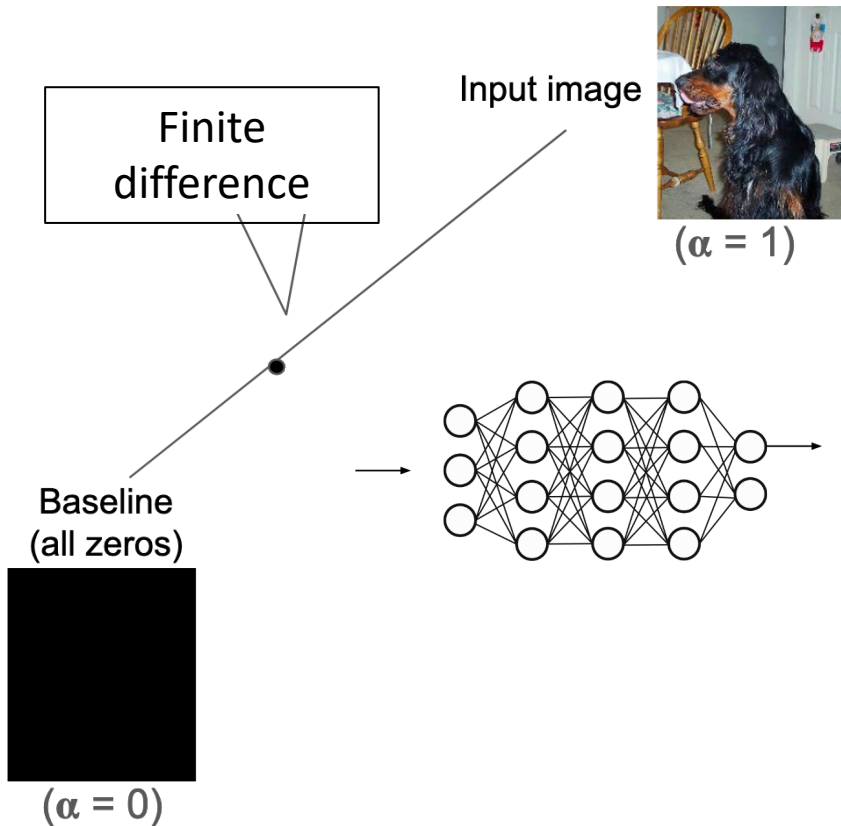
Integrated gradients



$$\phi_i^{IG}(f, x, x') = \underbrace{(x_i - x'_i)}_{\text{Difference from baseline}} \times \underbrace{\int_{\alpha=0}^1}_{\text{From baseline to input...}} \underbrace{\frac{\delta f(x' + \alpha(x - x'))}{\delta x_i}}_{\text{...accumulate local gradients}} d\alpha$$

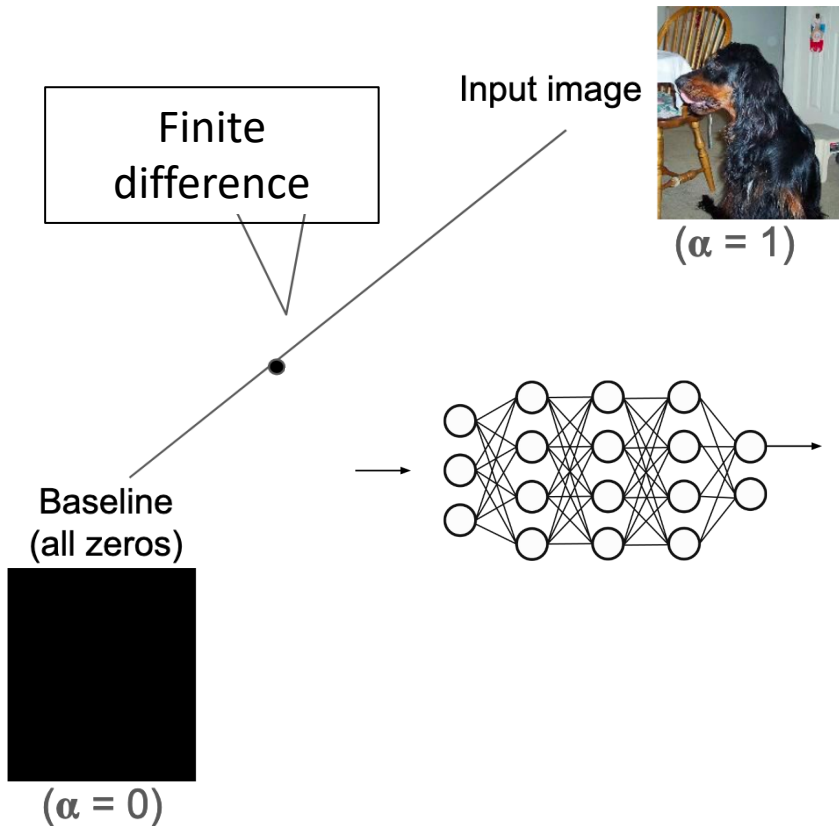
Can we think of disadvantages of this method?

DeepLIFT



- Explain “difference from reference value” of output in terms of “difference from reference value” of inputs

DeepLIFT

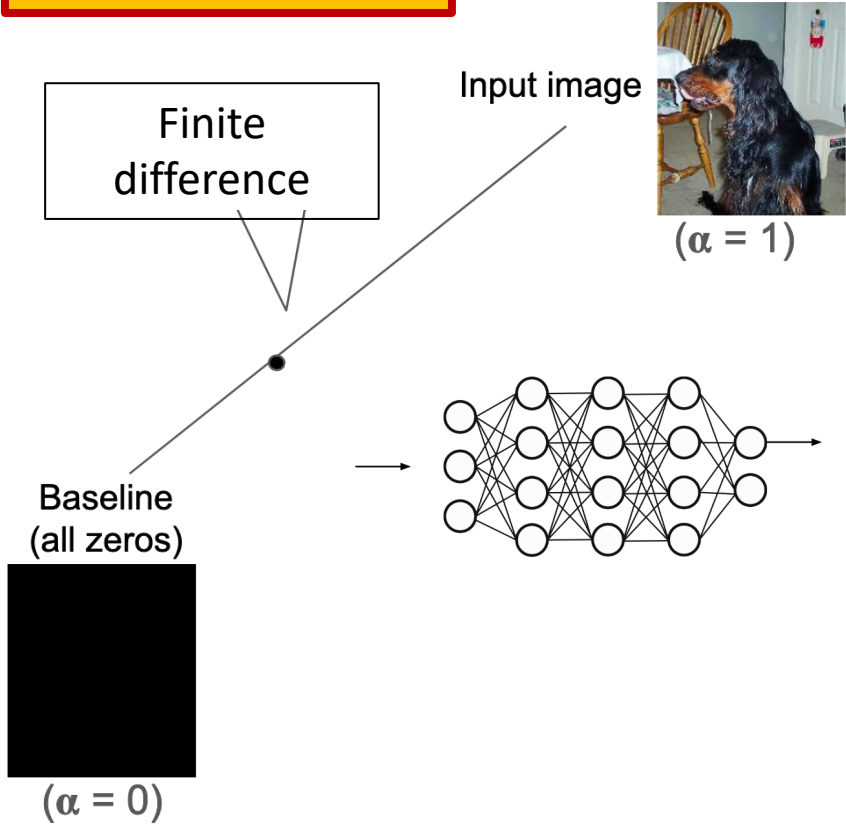


- Explain “difference from reference value” of output in terms of “difference from reference value” of inputs
- Target neuron t with diff-from-ref Δt



DeepLIFT

Can you think of any disadvantage of this method?



shap.DeepExplainer

```
class shap.DeepExplainer(model, data, session=None, learning_phase_flags=None)
```

- Explain “difference from reference value” of output in terms of “difference from reference value” of inputs
- Target neuron t with diff-from-ref Δt
- “Blame” Δt on $\Delta x_1 \dots \Delta x_n$
- Assign contributions $C_{\Delta x_i \Delta t}$ such that:

$$\sum_{i=1}^n C_{\Delta x_i \Delta t} = \Delta t$$

Any questions?

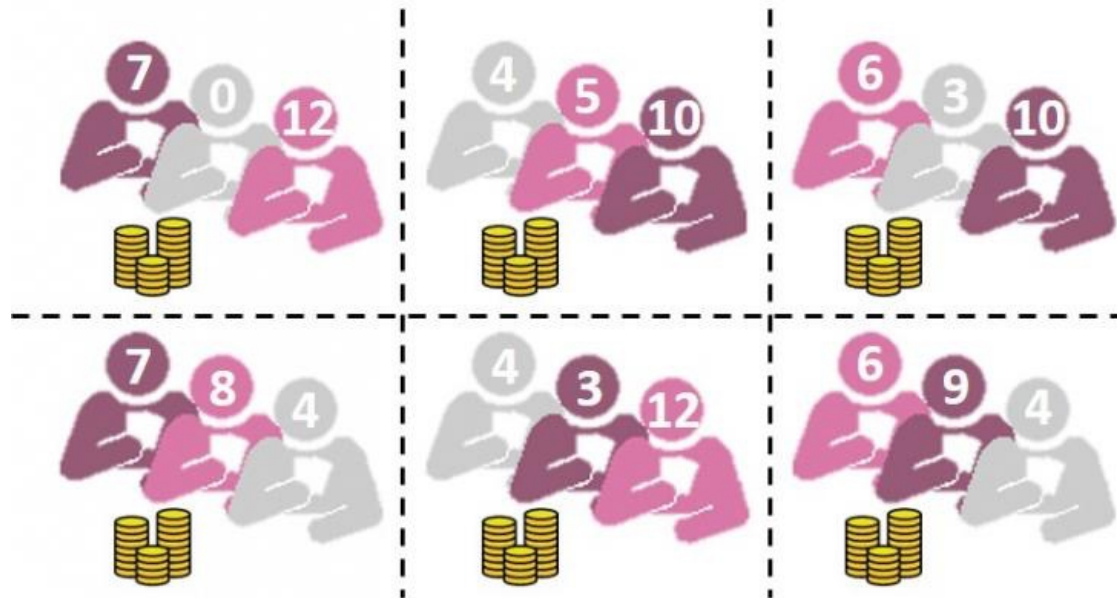


DeepLIFT (and Shapley values)

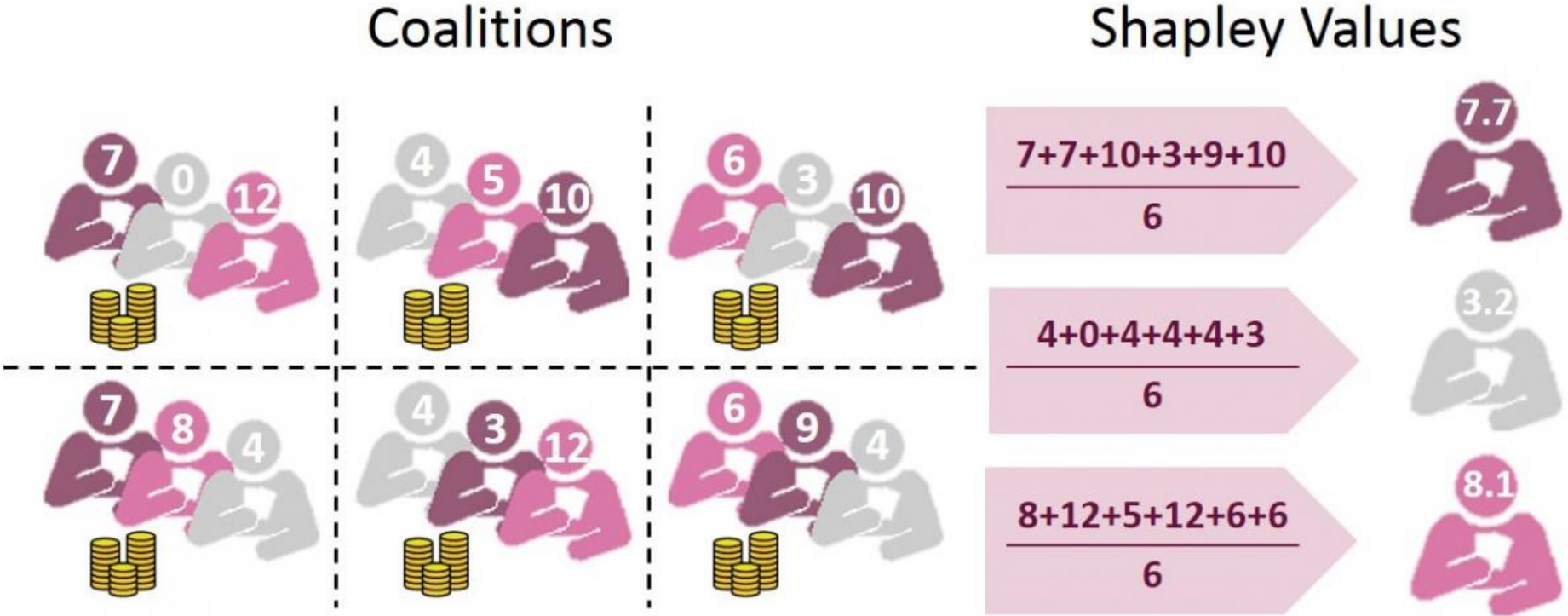


DeepLIFT (and Shapley values)

Coalitions



DeepLIFT (and Shapley values)



Shapley value not only considers the ability of each member, but also takes into account the cooperation among the members.

Today's goal – learn about interpretation in DL

(1) Model architecture based methods

(2) Gradient-based methods

(3) Model agnostic methods

Next time!

Recap

Model architecture based

CNNs (Global Avg Pooling, Deconv)

RNNs (Cell mapping, temporal output)

Gradient based and Shapley based

Saliency map (Gradient * Input)

Integrated gradients

DeepLIFT

Cell sensitive to position in line:
The sole importance of the crossing of the Berezina lies in the fact that it plainly and indubitably proved the fallacy of all the plans for cutting off the enemy's retreat and the soundness of the only possible line of action--the one Kutuzov and the general mass of the army demanded--namely, simply to follow the enemy up. The French crowd fled at a continually increasing speed and all its energy was directed to reaching its goal. It fled like a wounded animal and it was impossible to block its path. This was shown not so much by the arrangements it made for crossing as by what took place at the bridges. When the bridges broke down, unarmed soldiers, people from Moscow and women with children who were with the French transport, all--carried on by vis inertiae--pressed forward into boats and into the ice-covered water and did not surrender.

