# Layered Explanations:
## Interpreting Neural Networks with Numerical Influence Measures
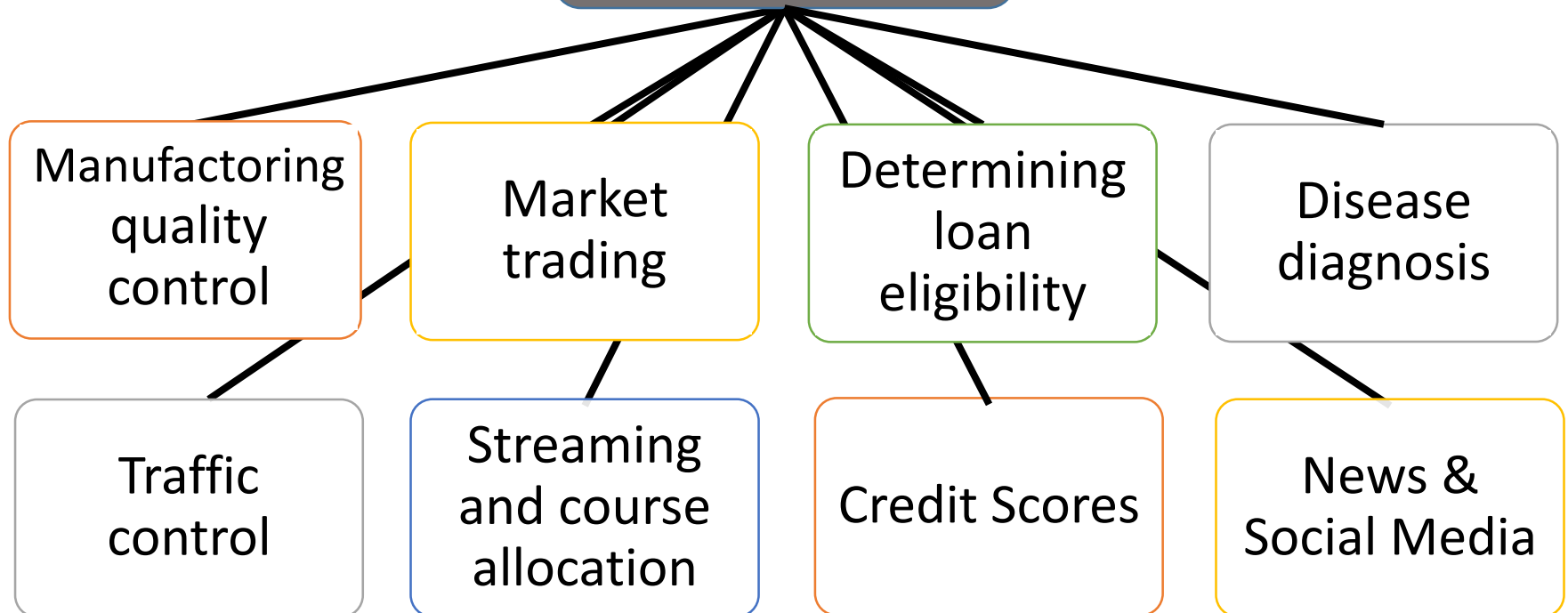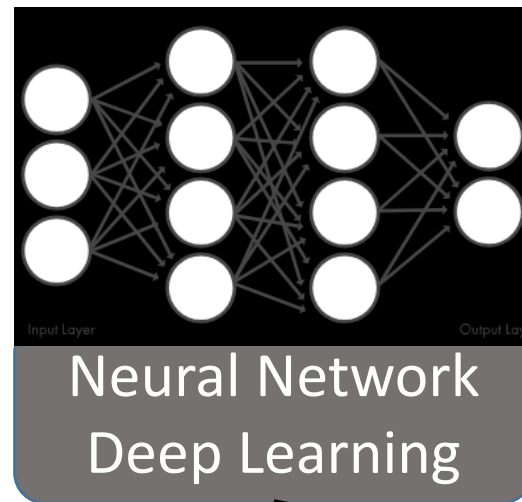
Ho Xuan Vinh
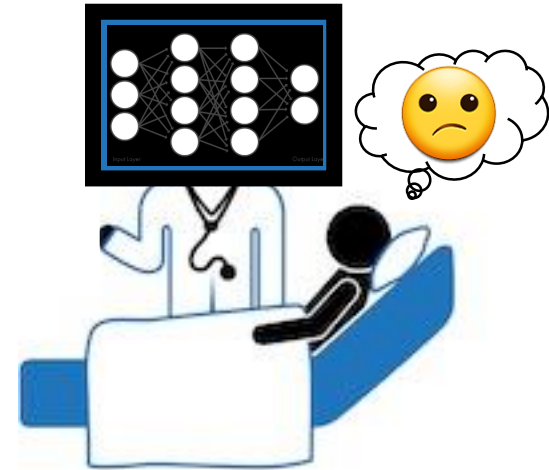
National University of Singapore

School of Computing

# Algorithmic Transparency: Motivation

Neural Network
Deep Learning

- Manufactoring quality control
- Market trading
- Determining loan eligibility
- Disease diagnosis
- Traffic control
- Streaming and course allocation
- Credit Scores
- News & Social Media

## Trustworthy

- Doctor-patient relationship
- Doctor is a specialist
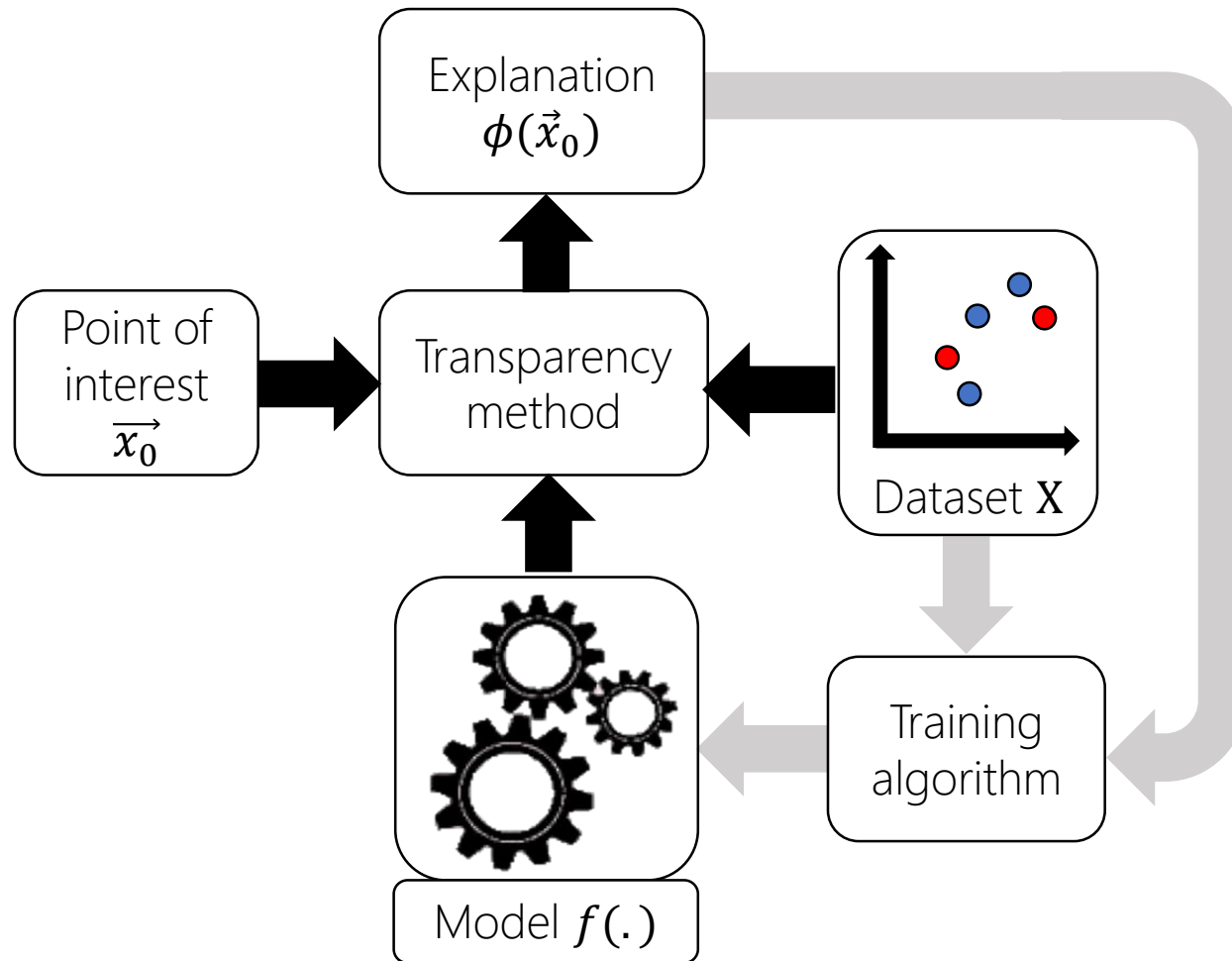- Doctor has a personal stake
- Reinforced reliability through time

## Untrustworthy

- How knowledgeable is NN?
- NN faces no consequence
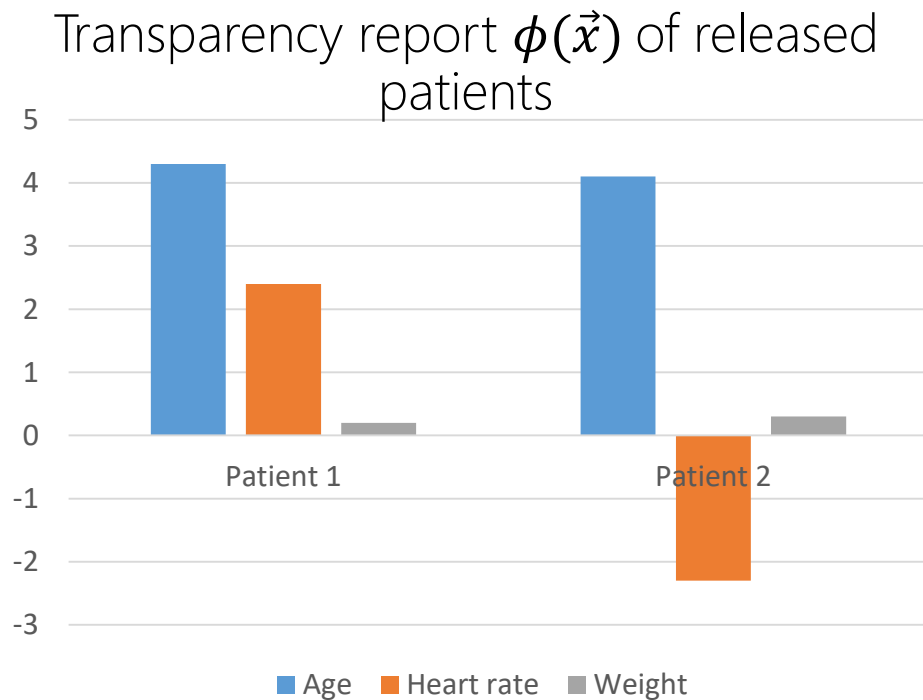- Unreliable historical records

Given $\vec{x} \in \mathbb{R}^n$, why was $\vec{x}$ labeled $f(\vec{x})$?

- $\phi(\vec{x}) \in \mathbb{R}^n$ - an explanation. May use the dataset $X$ and additional domain knowledge.

- *Linear explanation* is mainly focused

Transparency report $\phi(\vec{x})$ of released patients



Age   Heart rate   Weight

## Challenges

- Features often do not have intrinsic meaning
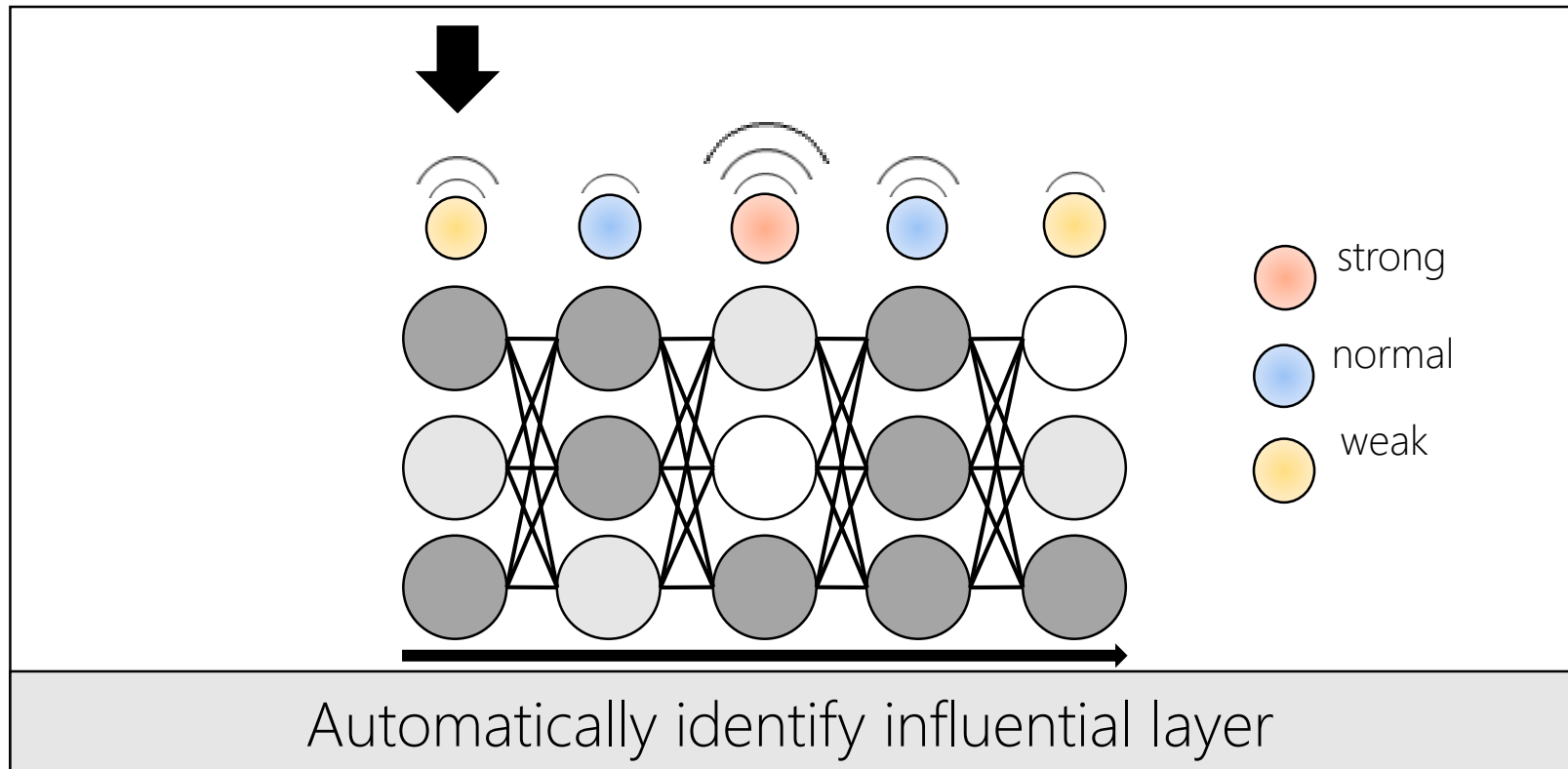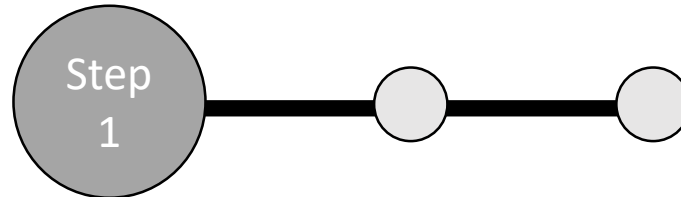- Individual measured effect is smalll

## Our approach
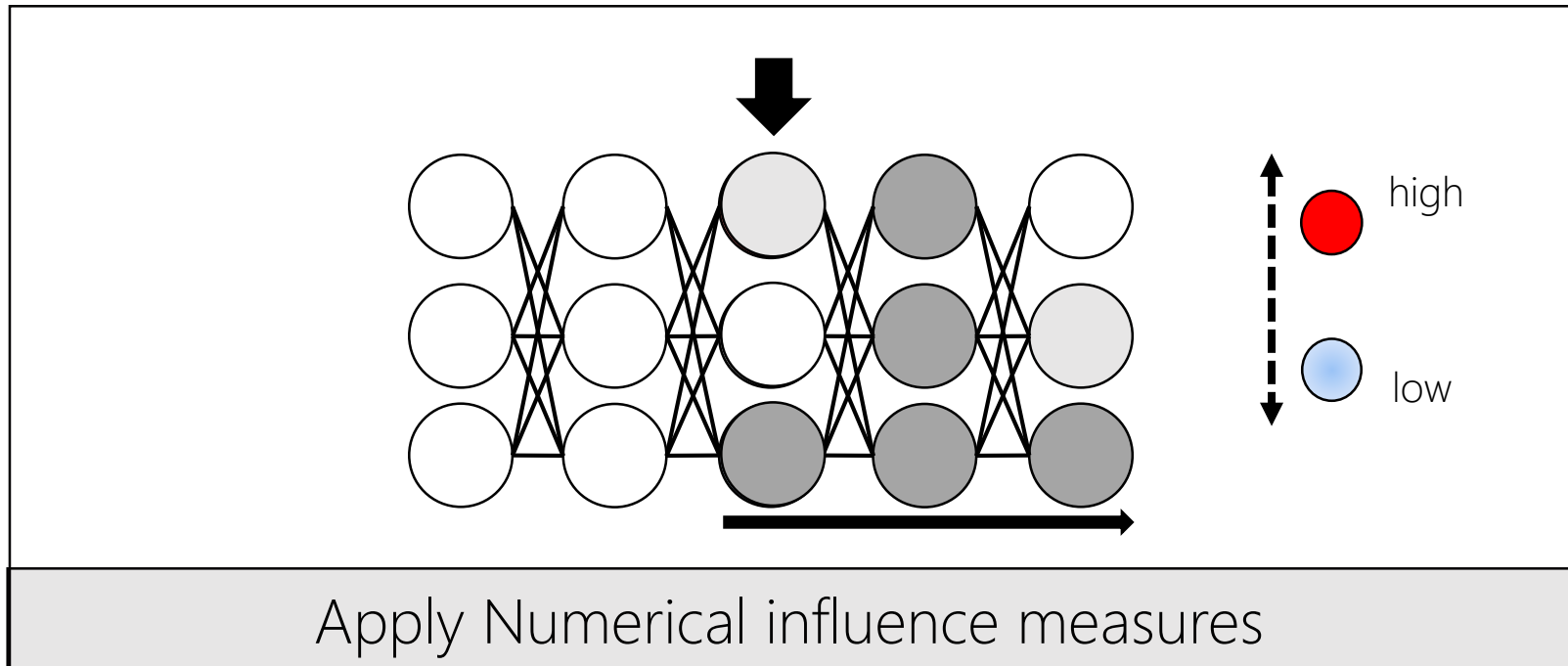
- Measure influence in Neural Network's internal layers

## Layered Explanations Framework

Provide explanation of learning model w.r.t. an instance



Automatically identify influential layer
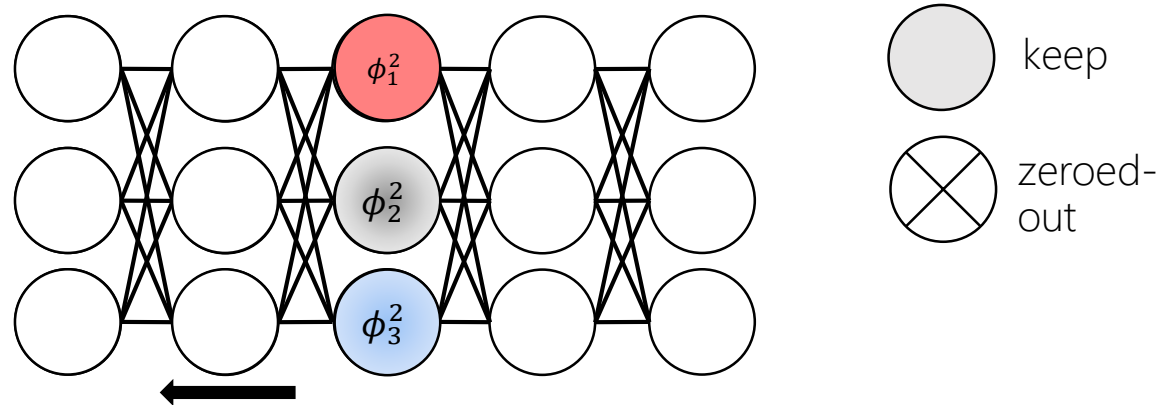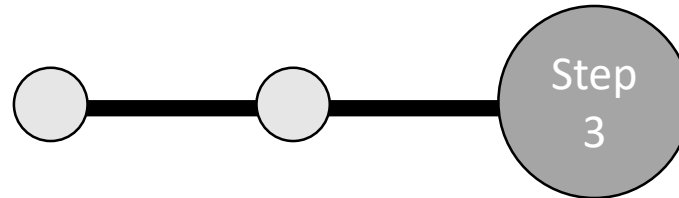
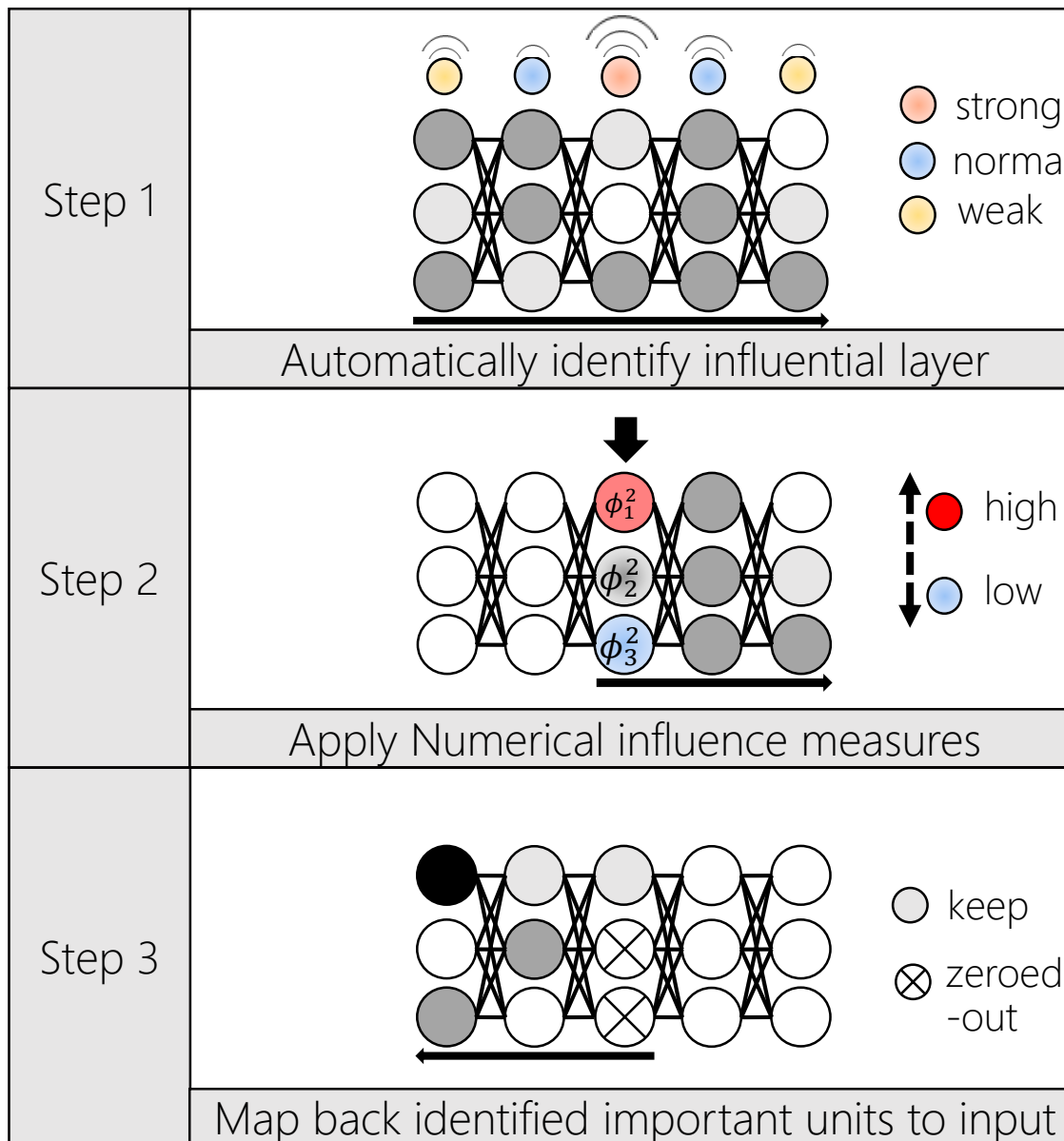Provide explanation of learning model w.r.t. an instance



Apply Numerical influence measures

Provide explanation of learning model w.r.t. an instance



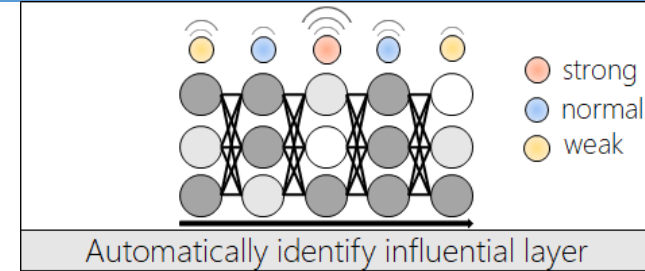Map back identified important units to input

Provide explanation of learning model w.r.t. an instance

## Motivation:

- Potentially capture high-level feature
- Receptive field are bigger
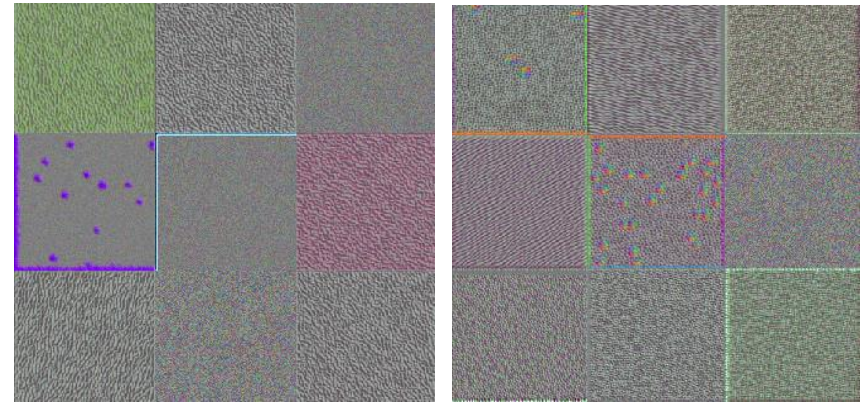


Automatically identify influential layer

## Input:

- Model $f$

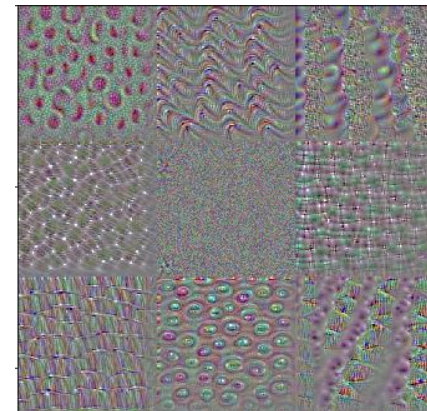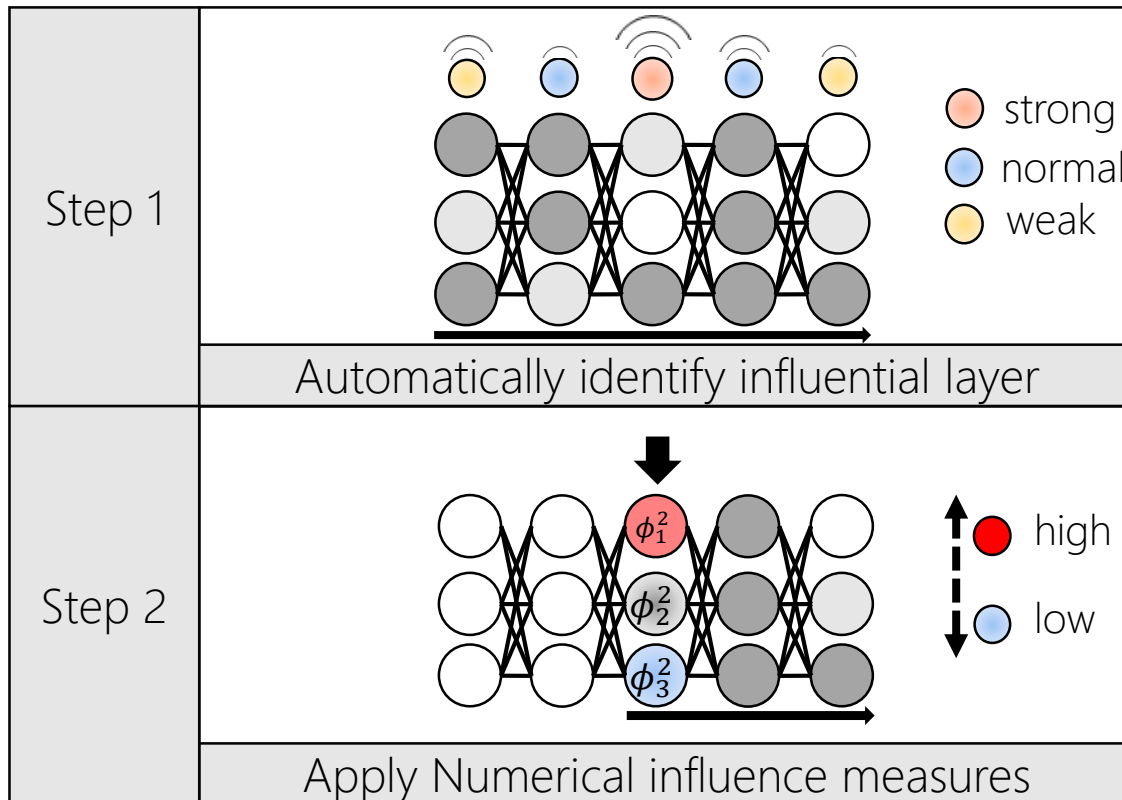## Output

- Chosen layer $L_k$



## Method:

- Influence measures
- Linear probe
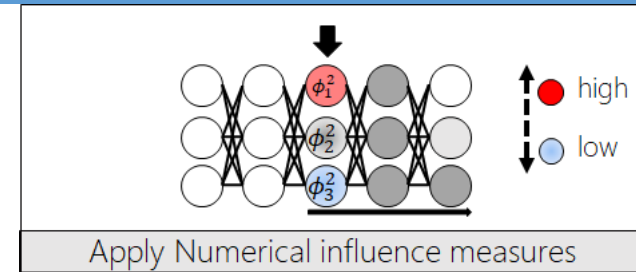  - Most robust layer
  - Most discriminative layer

Provides explanation of learning model w.r.t. an instance

Apply Numerical influence measures

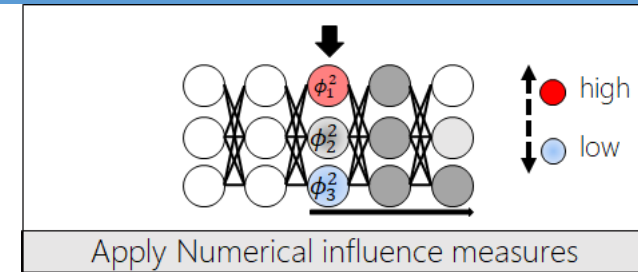Given $\vec{x}$, why was $\vec{x}$ labeled $f(\vec{x})$?

- $\phi(\vec{x})$ - an explanation. May use the dataset $\mathbf{X}$ and additional domain knowledge.

- <u>*This talk*</u>: numerical explanations
  - $\phi(x) \in \mathbb{R}^n$, where $\phi_i(\vec{x})$ is how important was feature $i$ in determining $f(\vec{x})$?
  - MIM, QII, LIME, Parzen, DeepLIFT...

- Our method can use any numerical influence measure. We focus on MIM
  - Computationaly efficient
  - Axiomatically justified

## Input:

- Datapoint $\vec{x} = \{x_1, x_2, \ldots, x_n\}$
- Dataset $\mathrm{X}$
- Model $f$



Apply Numerical influence measures

## Output:

- Influence score $\phi^i(\vec{x})$ of each feature/unit $i$ in $\vec{x}$

## Method:

Quantitative Input Influence | Monotone Influence Measure

(QII)             (MIM)

## Quantitative Input Influence | Monotone Influence Measure

### Intuition

The obtained influence measure $\phi_{QII}(\vec{x}, X)$ tells how much each feature _marginally contributes_ to model's outcome $Q(\vec{x}, N)$ w.r.t. instance $\vec{x}$.

### Marginal contribution $\phi_{QII}^i(\vec{x}, X)$

$$\phi_{QII}^i(\vec{x}, X) = \frac{1}{n!} \sum_{\pi \in \Pi} [Q(\vec{x}, S_\pi(i) \cup \{i\}) - Q(\vec{x}\, S_\pi(i))]$$

Over the permutation set $\Pi$ of $n$ features

Quantitative Input Influence | Monotone Influence Measure

## Intuition

The obtained influence measure $\phi_{MIM}(\vec{x}, \mathrm{X})$ indicates global direction that <u>*strengthens*</u> the current label of instance $\vec{x}$.
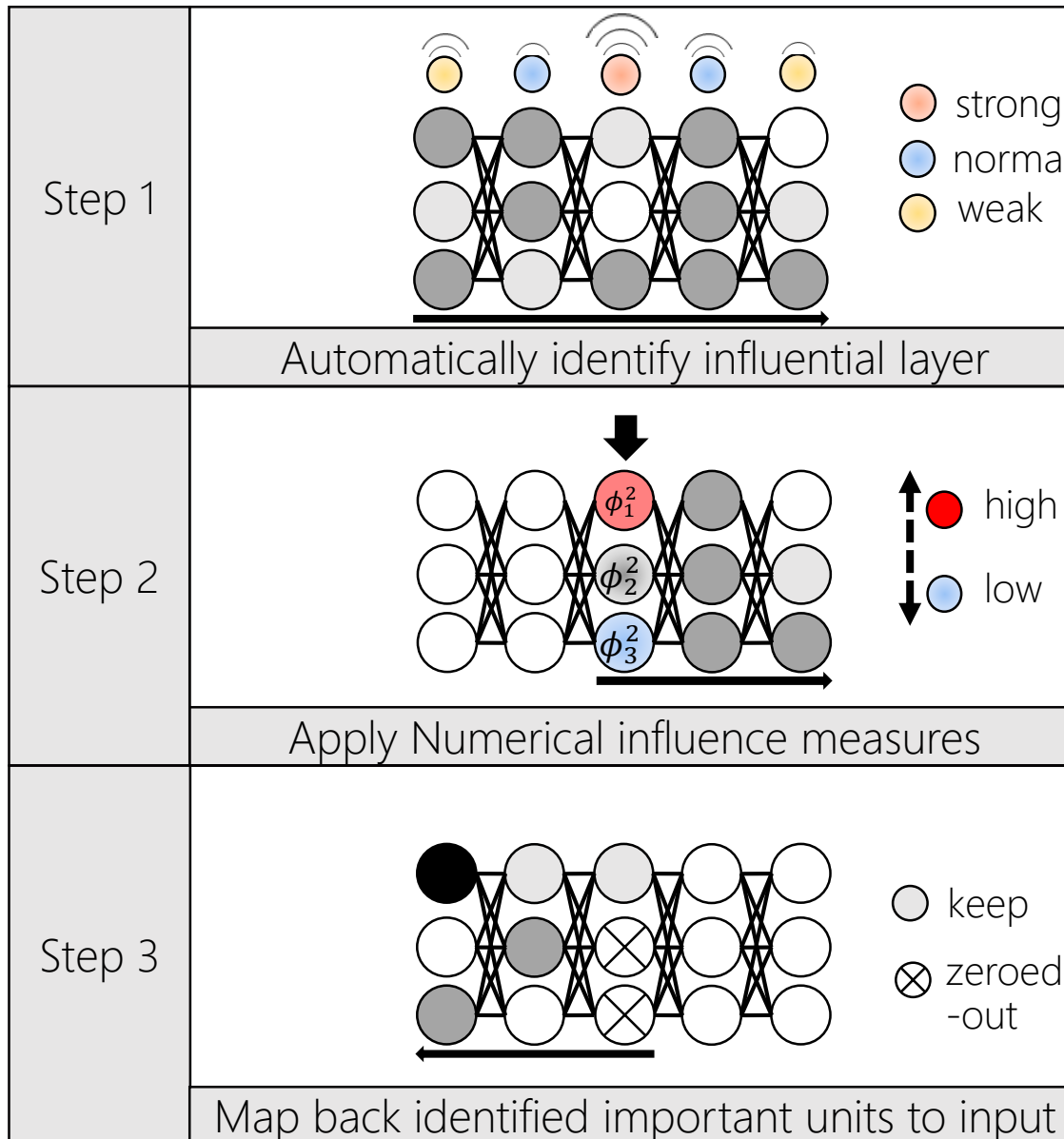
## Influence measure $\phi^i_{MIM}(\vec{x}, \mathrm{X})$

Function of Distance

$$\phi^i_{MIM}(\vec{x}, \mathrm{X}) = \sum_{\vec{y} \in \mathrm{X} \setminus \vec{\mathrm{x}}} (\vec{y} - \vec{x})\, \alpha(\| \vec{y}, \vec{x} \|)\, \mathbb{I}(f(\vec{x}) = f(\vec{y}))$$

{−1,1}-valued Indicator function

Over all data points except for $x$

Provides explanation of learning model w.r.t. an instance

Map back identified important units to input

○ keep

⊗ zeroed -out

## Input:
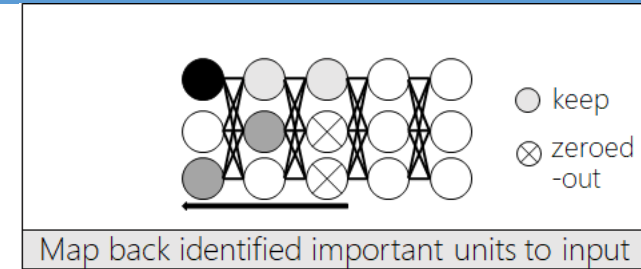
- Datapoint $\vec{x}$ and its value in each layer $L_i$, $\forall i \in \{1, \dots, K\}$
- Influence scores at chosen layer $L_k$
- Neural network $f$
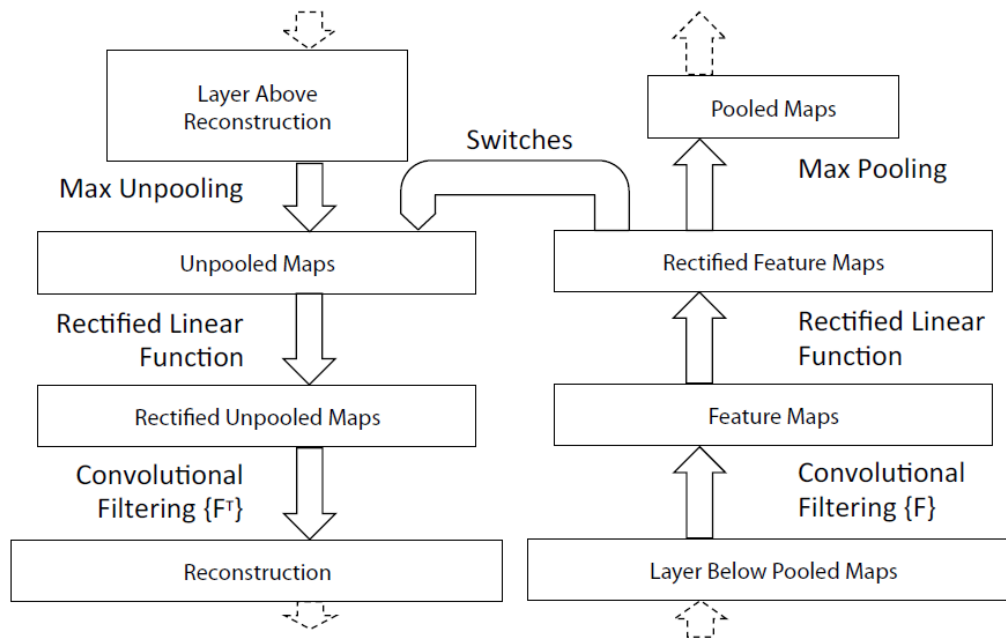
## Output:

- Back-propagated value at input space $L_1$

## Method:
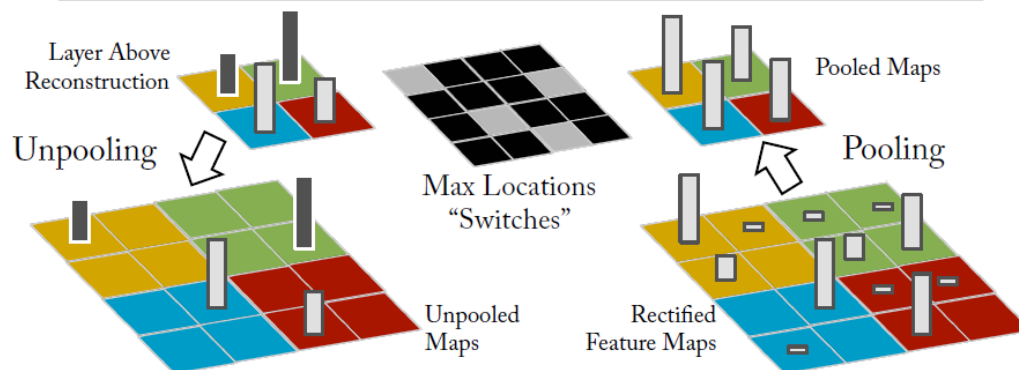
Deconvolutional Neural Net | Guided Backpropagation

(DeconvNet)　　　　　　　　(GuidedBackprop)

## Deconvolutional Neural Net

Guided Backpropagation



- Unpooling: switch
- Deconvolutional layer
- ReLU in backward phase

Zeiler, M. D. and Fergus, R. (2014). Visualizing and understanding convolutional networks. In European conference on computer vision, pages 818–833. Springer.

## Deconvolutional Neural Net

## Guided Backpropagation

Forward pass



Backward pass



$$\mathrm{R}^i = (f^i > 0).(R^{i+1} > 0).R^{i+1}$$

$$\text{where } R^{i+1} = \frac{\partial f^{out}}{\partial f^{i+1}}$$

# Experiment

- Dataset & classifier:
  - MNIST (digit pair 6-9) + 3-convolutional-layer CNN
  - Dog-Fish (extracted from ImageNet) + VGG16 net

- Comparison methods:
  - MIM on input layer
  - Guided-Activation:
    - Masking $M = \mathbb{I}\left(|\vec{x}_{L_i}| \geq \delta\right)$
  - Guided-MIM: MIM + GuidedBackprop
    - Masking $M = \mathbb{I}\left(|\phi_{MIM}\left(\vec{x}_{L_i}\right)| \geq \delta\right)$

- Parameters:
  - We set $\delta$ such that only top 1% influential hidden units are chosen.

Normalized value of $\left\|\phi_{MIM}^{i}(\vec{x}, X)\right\|$ of four samples.

$$L^* = \underset{L_i \in \{L_2, \ldots, L_{K-1}\}}{\arg\max} \frac{\left\|\phi_{MIM}^{i}(\vec{x}, X)\right\|_2}{\max\limits_{\vec{y} \in X \backslash \vec{x}} \left\| \vec{y}_{L_i} - \vec{x}_{L_i} \right\|_2}$$

# Experiment: MNIST



Heatmaps of three methods: MIM, Guided-Activation, and Guided-MIM on four samples of the MNIST dataset, only 1% of total units in each layer backpropagated.

Normalized value of $\left\| \phi_{MIM}^{i}(\vec{x}, X) \right\|$ of four samples.

$$L^{*} = \underset{L_i \in \{L_2, \dots, L_{K-1}\}}{\operatorname{argmax}} \frac{\left\| \phi_{MIM}^{i}(\vec{x}, X) \right\|_2}{\max_{\vec{y} \in X \setminus \vec{x}} \left\| \vec{y}_{L_i} - \vec{x}_{L_i} \right\|_2}$$

a) Low predicted probability

b) Incorrect prediction

c) Clear background

d) Noisy background

| Dataset | Label $t$ | $f(\vec{x}) = t$ | $f(\vec{x} + \phi_{MIM}(\vec{x})) = t$ | $f_P(\vec{x} + \phi_{MIM}(\vec{x})) > f_P(\vec{x})$ |
|---------|-----------|------------------|----------------------------------------|------------------------------------------------------|
| MNIST | Digit 6 | 1008 | 1008 | 262 |
|  | Digit 9 | 959 | 959 | 875 |
| Dog-Fish | Dog | 298 | 298 | 138 |
|  | Fish | 302 | 302 | 0 |

Number of samples increases their predicted probabilities $f_P(.)$ by shifting toward the vector $\phi_{MIM}(\vec{x})$ on MNIST and Dog-fish dataset.

- Guided-MIM is a preliminary step toward a better model.
  - Guided-Activation and Guided-MIM can outline the full shape of target object.

- Possible extensions:
  - Step 1: A NN slice marks transition from general to specific-class feature.
  - Step 2: Other numerical influence methods.
  - Step 3: Backpropate influence scores and redistribute to input layer.
  - Combination of different influence measure – backpropagation pair requires different interpretation.
  - Quantitative evaluation

- Assumption
  - pretrained model has acceptable accuracy

# Conclusion

- Proposes Layered Explanations Framework
    1. Identify the greatest-explanatory layer and
    2. Influential units using numerical influence measures, then we
    3. Reconstruct relevant input regions responsible for activating these influential units.

- Experiment on MNIST & Dog-fish dataset:
    - combination of MIM and GuidedBackprop
    - able to outline target object, but no clear signal of identifying influential components.

- Possible extensions
    - applicable methods on step 2 and 3
    - possible intepretation of selected method pairs
    - quantitative evaluation

# References

- **DeconvNet**: Zeiler, M. D. and Fergus, R. (2014). Visualizing and understanding convolutional networks. In European conference on computer vision, pages 818–833. Springer.

- **Guided Backpropagation**: Springenberg, J. T., Dosovitskiy, A., Brox, T., and Riedmiller, M. (2014). Striving for simplicity: The all convolutional net. arXiv preprint arXiv:1412.6806

- **QII**: Datta, A., Sen, S., and Zick, Y. (2016). Algorithmic transparency via quantitative input influence: Theory and experiments with learning systems. In Security and Privacy (SP), 2016 IEEE Symposium on, pages 598–617. IEEE.

- **MIM**: Sliwinski, J., Strobel, M., & Zick, Y. (2017). A Characterization of Monotone Influence Measures for Data Classification. arXiv preprint arXiv:1708.02153.