# Explainable AI (XAI)

**Daisuke Sato**

2020/07/11

**Reference** (Survey on XAI)
- A. Adadi, M. Berrada, https://ieeexplore.ieee.org/document/8466590/
- Finale Doshi-Velez, Been Kim, arXiv:1702.08608

**Codes**
Notebook in Kaggle by Daisuke Sato
- "Random forest/XGBoost, and LIME/SHAP with Titanic"
- "Predict Housing price"

- **Daisuke Sato, Ph.D.**
  - **Graduate school**
    - **Kyoto university**
    - **Research topic: Theoretical Physics (Quark-Gluon Plasma)**
  - **Previous jobs**
    - **Postdoc (US., Italy, Germany)**
  - **Current position**
    - **CTO**

- **Concept/Motivation**

- **Recent trends on XAI**

- **Method 1: LIME/SHAP**
  - Example: Classification
  - Example: Regression
  - Example: Image classification

- **Method 2: ABN for image classification**

Generally speaking, AI is a blackbox.

We want AI to be explainable because⋯

1. **Users should trust AI to actually use it (prediction itself, or model)**
   Ex: diagnosis/medical check, credit screening

G. Tolomei, et. al.,  arXiv:1706.06691

People want to know why they were rejected by AI screening,
and what they should do in order to pass the screening.

## 2. It helps to choose a model from some candidates

Classifier of text to "Christianity" or "Atheism" (無神論)



Green: Christianity
Magenta: Atheism

**Both model give correct classification,
but it is apparent that model 1 is better than model 2.**

**3. It is useful to find overfitting, when train data is different from test data**

Cf: Famous example of "husky or wolf"

Training dataset contains pictures of wolfs with snowy background.



a) Husky classified as wolf

Then, the classifier trained on that dataset outputs "wolf" if the input image contains snow.

- Concept/Motivation

- **Recent trends on XAI**

- Method 1: LIME/SHAP
  - Example: Classification
  - Example: Regression
  - Example: Image classification

- Method 2: ABN for image classification

# of papers which includes
one of explanation-related words ("intelligible", "interpretable",···)
AND
one of AI-related words ("Machine learning", "deep learning",···)
FROM
7 repositories (arXiv, Google scholar, ···)



**Recently, researchers are studying XAI more and more.**

Contents

- Concept/Motivation

- Recent trends on XAI

- **Method 1: LIME/SHAP**
    - Example: Classification
    - Example: Regression
    - Example: Image classification

- Method 2: ABN for image classification

Objects: Classifier or Regressor

**Interpretable model**

$x$: specific input

**Basic idea: Approximate original ML model with interpretable model (linear model/DT), in the vicinity of specific features.**

Arithmer

*Original model*

$$f$$ ?

$$\mathbb{R}^d \longrightarrow \mathbb{R}$$

$$\cup \qquad\qquad \cup$$

$$x \qquad\qquad f(x)$$

*features*      *prediction*

*Interpretable model*

$$g$$

$$\{0,1\}^{d\prime} \longrightarrow \mathbb{R}$$

$$\cup \qquad\qquad \cup$$

$$z' \qquad\qquad g(z')$$

*Absence/presence of features*    *prediction*

Distance of predictions of original model (*f*) and its approximation

Complexity of model
(depth of threes for DT/
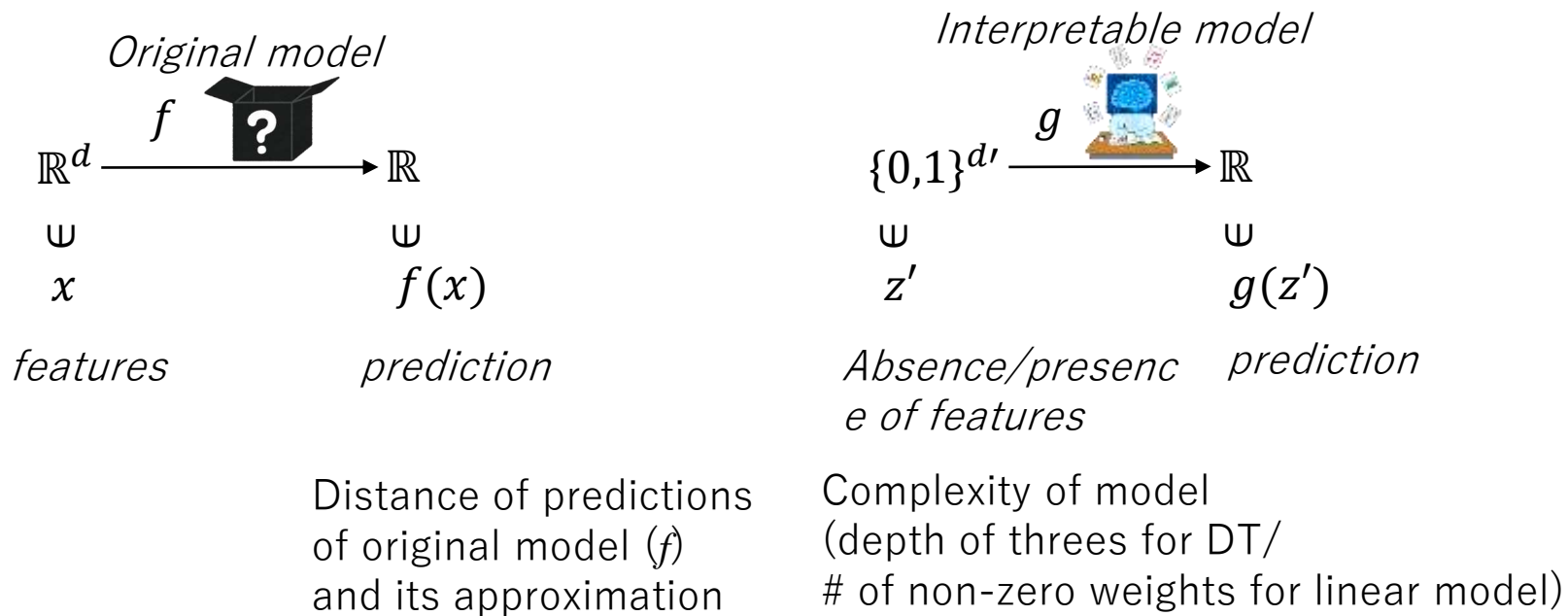# of non-zero weights for linear model)

$$\underset{g \in G}{\mathrm{argmin}} \quad \mathcal{L}(f, g, \pi_x) + \Omega(g)$$

Linear model

$$\mathcal{L}(f, g, \pi_x) = \sum_{z, z' \in \mathcal{Z}} \pi_x(z) \left( f(z) - g(z') \right)^2$$

$$g(z') = w_g \cdot z'$$

Weights of sampling

$$\pi_x(z) = \exp(-D(x, x')^2 / \sigma^2)$$

1. Original model predicts from features (sneeze, headache,⋯) whether the patient is flu or not.

2. LIME approximates the model with linear model in the vicinity of the specific patient.

3. The weights of the linear model for each feature give "explanation" of the prediction

- **Interpretable**



- **Local fidelity**

- **Model-agnostic**
  （Original model is not
  affected by LIME at all）

- **Global perspective**
  （sample different inputs
  and its predictions）

S. Lundberg, S-I. Lee, arXiv:1705.07874

Generalization of methods for XAI,
*   LIME
*   DeepLIFT   A. Shrikumar et. al., arXiv:1605.01713
*   Layer-Wise Relevance Propagation  Sebastian Bach et al. In: PloS One 10.7 (2015), e0130140

Actually, they are utilizing concepts of cooperative game theory:
*   Shapley regression values
*   Shapley sampling values
*   Quantitative Input Influence

Label
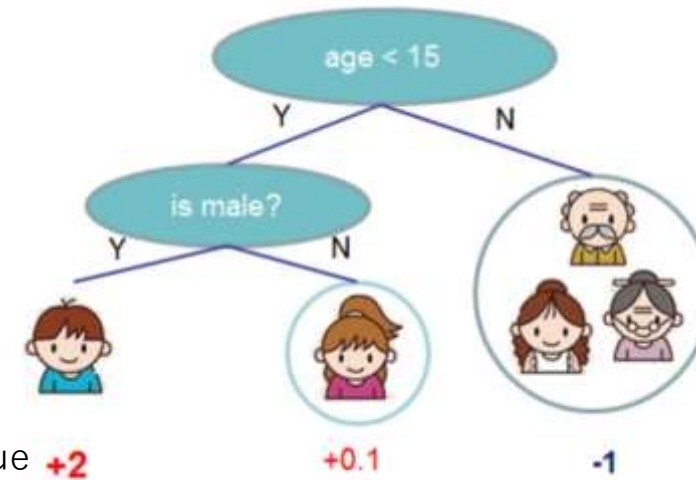
| Variable | Definition | Key |
|---|---|---|
| survival | Survival | 0 = No, 1 = Yes |
| pclass | Ticket class | 1 = 1st, 2 = 2nd, 3 = 3rd |
| sex | Sex | |
| Age | Age in years | |
| sibsp | # of siblings / spouses aboard the Titanic | |
| parch | # of parents / children aboard the Titanic | |
| ~~ticket~~ | ~~Ticket number~~ | |
| fare | Passenger fare | |
| ~~cabin~~ | ~~Cabin number~~ | |
| embarked | Port of Embarkation | C = Cherbourg, Q = Queenstown, S = Southampton |

Features

All values are unique

77% are NaN

\# of training samples = 891 →Train : Validation = 8 : 2
\# of test samples = 418

**Decision Tree (DT)**

**Selection of features and its threshold is done, so that the mutual information is maximized.**



Entropy in
parent node    Entropy in child nodes

"Boyness" value

Mutual information: $I_H(D_p) - I_H(D_{left}) - I_H(D_{right})$

Entropy at a node is maximized when the samples in that node are uniformly distributed to all the classes.

→**Child nodes tends to be purely distributed**

☺

- Simple to understand and interpret
- Can represent interaction among features unlike linear model (at least up to # of depth_tree)
- No need to NaN preprocessing/standardization

☹

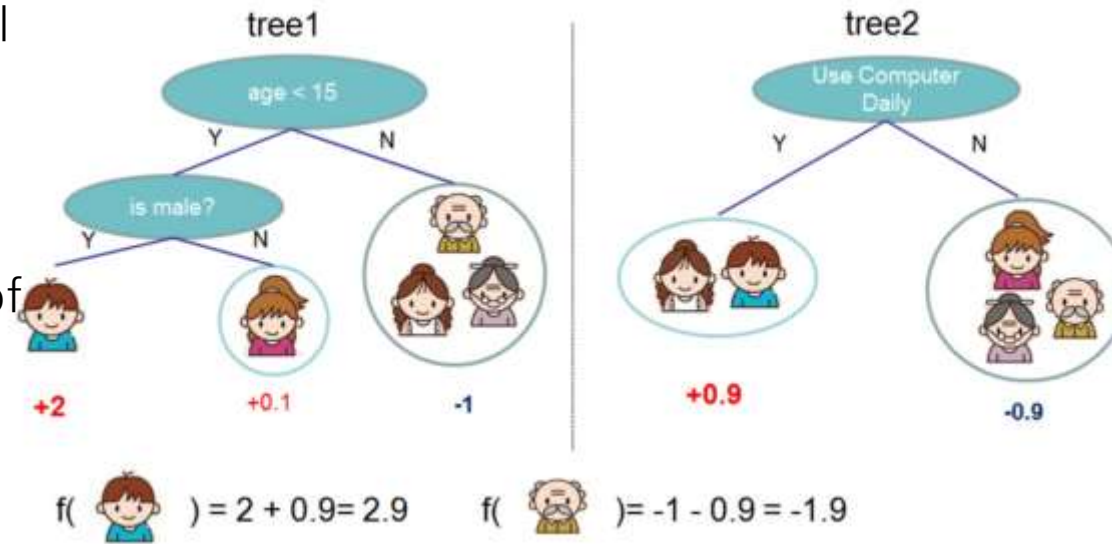- Tends to overlearn (if trees are deep)

## Gradient Boosting Decision Tree（GBDT）

In ensemble method, the final value is calculated from average of many trees.

How to decide the structure of tree ensemble (forest):
1. Start with one tree
2. Add another tree so that the loss function is minimized (gradient boost)



Cf: In random forest, all the trees are generated in parallel, at once.
Compared with it, GBDT can learn from the misprediction in the previous tree, so the result is expected to be better.

☺
- Overlearning problem is improved

☹
- Less interpretability
- Many hyperparameters to be tuned

Very popular method in Kaggle

I used one of GBDT (XGBoost) as a ML model, under the following conditions:

- No standardization of numerical features (as DT does not need it)
- No postprocessing of NaN (GBDT treats NaN as it is)
- No feature engineering
- Hyperparameter tuning for n_estimators and max_depth (optuna is excellent)

**Results**:

Best parameters: {'n_estimators': 20, 'max_depth': 12}
Validation score: 0.8659217877094972

Test score: 0.77033

Cf: Baseline (all women survive, all men die): 0.76555

**Not satisfying···**

Cf: Reported best score with ensemble method: 0.84210

Review of know-how on feature engineering by experts:
Kaggle notebook "How am I doing with my score?"

## Cf: Test score at Kaggle competition

| # | Team Name | Notebook | Team Members | Score |
|---|-----------|----------|--------------|-------|
| 1 | hongjungu | | | 1.00000 |
| 2 | qianbi | | | 1.00000 |
| 3 | Bayo Adekanmbi | | | 1.00000 |
| 4 | Yudistira Ashadi | | | 1.00000 |
| 5 | Alex Stone | | | 1.00000 |
| 6 | SteveKane | | | 1.00000 |
| 7 | Rum Yue | | | 1.00000 |
| 8 | umang aditya | | | 1.00000 |
| 9 | Keewon Shin | | | 1.00000 |
| 10 | jatin grover | | | 1.00000 |

Cheaters… (Memorizing all the names of the survivors? Repeat trying with GA?)

**Linear correlation between features and label**

（explanation of data itself, not model）

The method which you try first, to select important features.
You can work when the number of features is small.



Correlation between Survived and features

**Sex_female     0.543351**
**Pclass_3       0.322308**
Pclass_1       0.285904
Fare           0.257307
Embarked_C     0.168240
Embarked_S     0.155660

Fare has large correlation with Pclass, as you can easily expect (~0.4-0.6).
Also with Embarked_C (~0.27), which I do not know why.

**It is likely that the model trained with this data put large importance on the features with large correlation.**

## Importance of features in GBDT

（explanation of whole model, not for specific sample）

Feature importance

Women had high priority for getting on boats

Passengers in 3rd class are not likely to survive.
(Sign of feature's contribution can not be read, only the magnitude)

The average gain (for mutual information, or inpurity) of splits which use the feature

**Top 3 agrees with that in linear correlation.**

However, this method has a problem of "*inconsistency*" (when a model is changed such that a feature has a higher impact on the output, the importance of that feature can decrease)

**This problem is overcome by LIME/SHAP.**

## Explanation by LIME

（explanation of model prediction for specific sample）

My code did not work on Kaggle kernel, because of a bug in LIME package…
So, here I quote results from other person.

https://qiita.com/fufufukakaka/items/d0081cd38251d22ffebf



**As LIME approximates the model with linear function locally, the weights of the features are different depending on sample.**

In this sample, the top 3 features are Sex, Age, and Embarked.

**Results of SHAP**

Red points（Sex_female=1:female）have positive SHAP values, while the blue points（Sex_ｆemale=0:male） negative.
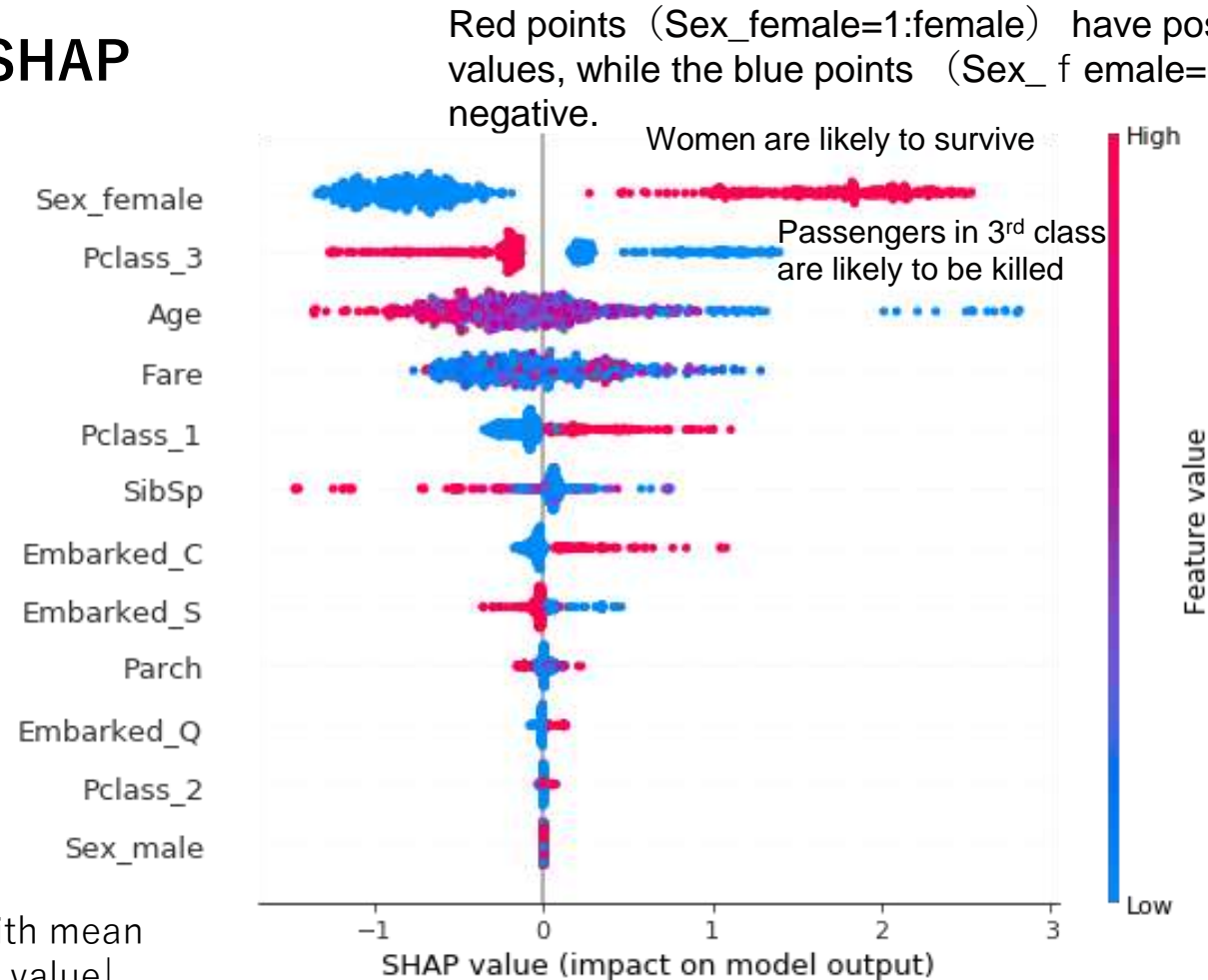


Women are likely to survive

Passengers in 3rd class are likely to be killed

Sorted with mean of |SHAP value|

**Top 3 does not agree with that in linear correlation/XGBoost score（Age enters）.**

**SHAP is consistent（unlike feature importance of XGBoost）and has local fidelity（unlike linear correlation）, I would trust SHAP result than the other two.**

| | Variable | Definition | Key |
|---|---|---|---|
| **Label** | SalePrice | the property's sale price in dollars | - |
| **Features** | OverallQual | Rates the overall material and finish of the house | 10:Very Excellent<br>9:Excellent<br>…<br>1:Very Poor |
| | GarageCars | Size of garage in car capacity | - |
| | KitchenAbvGr | Kitchens above grade | Ex:Excellent<br>Gd:Good<br>TA:Typical/Average<br>Fa:Fair<br>Po:Poor |
| | …(80 features in total) | … | … |

Dataset describing the sale of individual residential property in Ames, Iowa, from 2006 to 2010.

\# of training samples = 1460　→Train : Validation =75 : 25
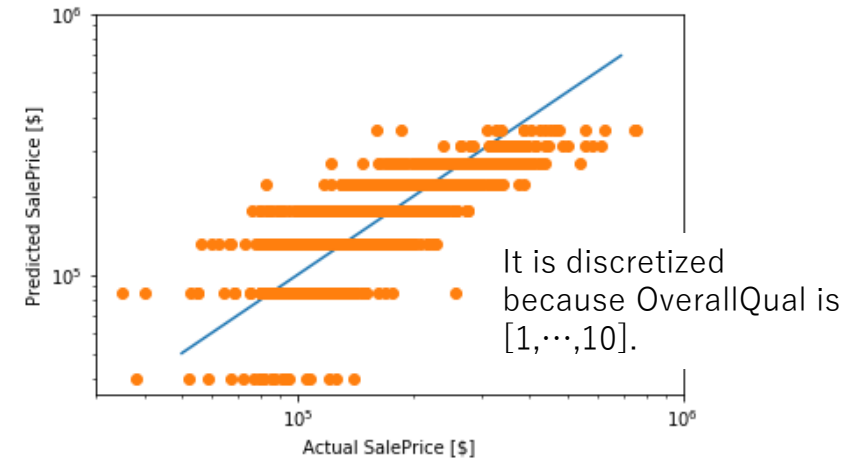\# of test samples =1459

Model: XGBoost

**Results** (metric is Root Mean Squared Error of Log (RMSEL) )

Validation : 0.13830

Test: 0.16332

Cf: Baseline (linear regression in terms of OverallQual): 1.1608



Prediction/actual value for validation data



It is discretized because OverallQual is [1,···,10].

**This time, it is much better than simple baseline.**

Cf: Score with the same method (XGBoost), but with feature engineering: 0.11420

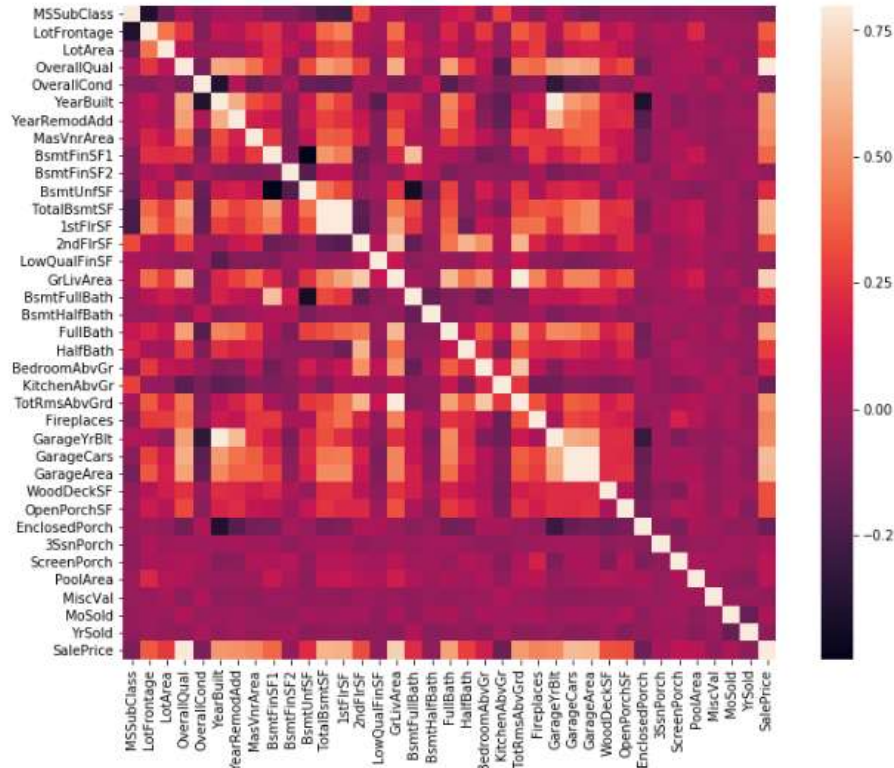Kaggle notebook "Stacked Regressions : Top 4% on LeaderBoard"

## Linear correlation between features and label

(explanation of data itself, not model)
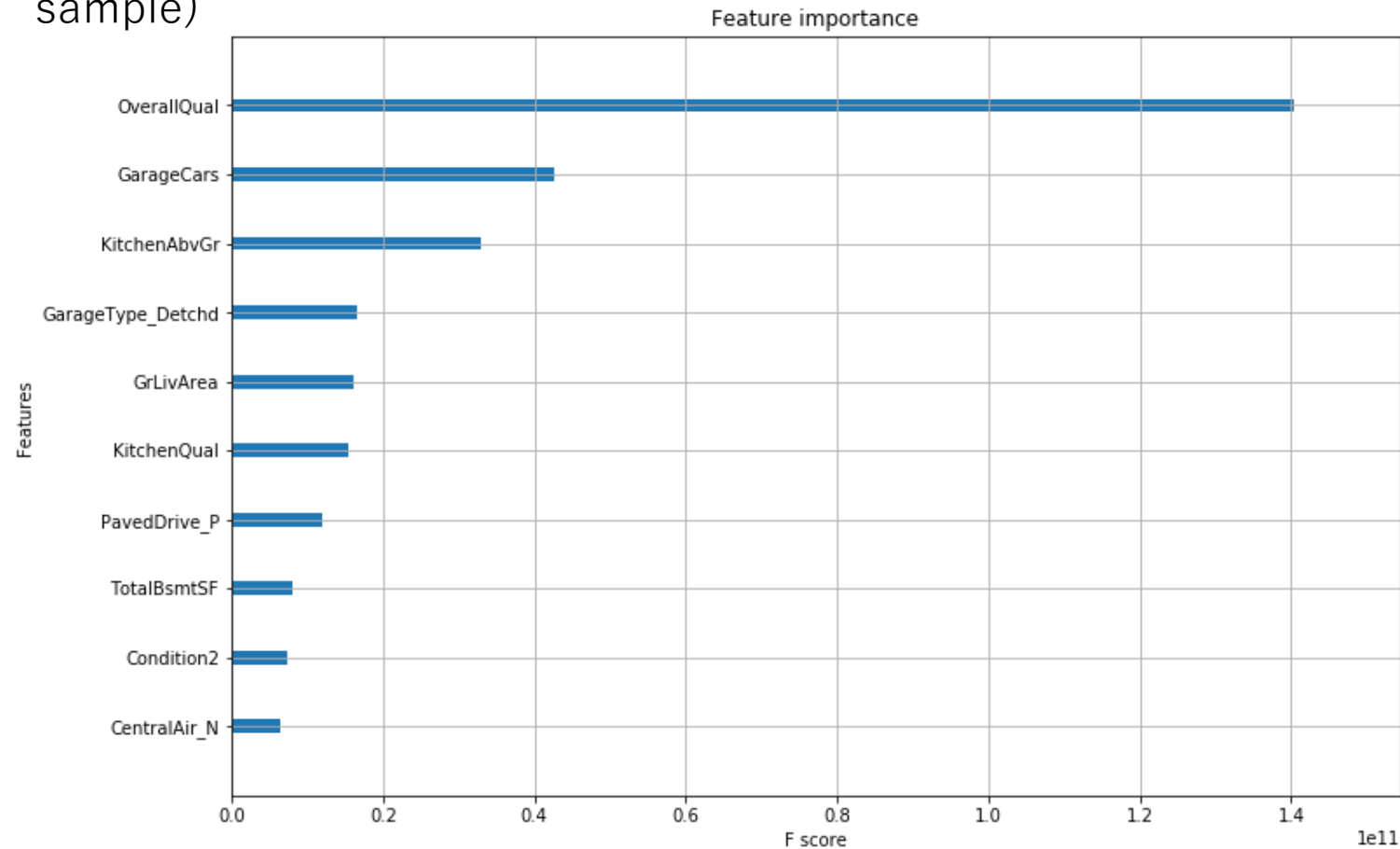
Correlation between SalePrice and features

Rates the overall material and finish of the
house
Above grade (ground) living area square feet

Size of garage in car capacity



| | |
|---|---|
| **OverallQual** | **0.790982** |
| **GrLivArea** | **0.708624** |
| **GarageCars** | **0.640409** |
| GarageArea | 0.623431 |
| TotalBsmtSF | 0.613581 |
| 1stFlrSF | 0.605852 |
| FullBath | 0.560664 |
| TotRmsAbvGrd | 0.533723 |
| YearBuilt | 0.522897 |
| YearRemodAdd | 0.50710 |
| … | |

**Importance of features in GBDT**

(explanation of whole model, not for specific sample)
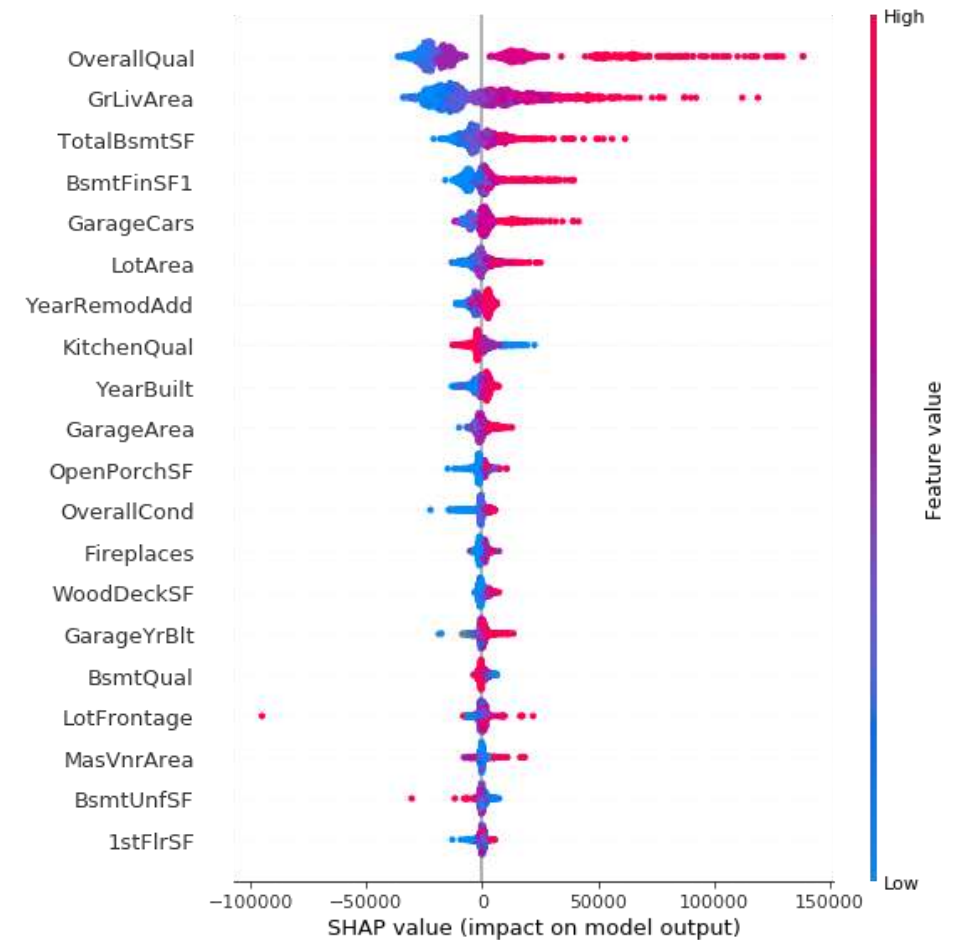


Feature importance

**Top 3 is different from linear correlation (KitchenAbvGr).**

## Explanation by SHAP

（explanation of model prediction for specific sample）
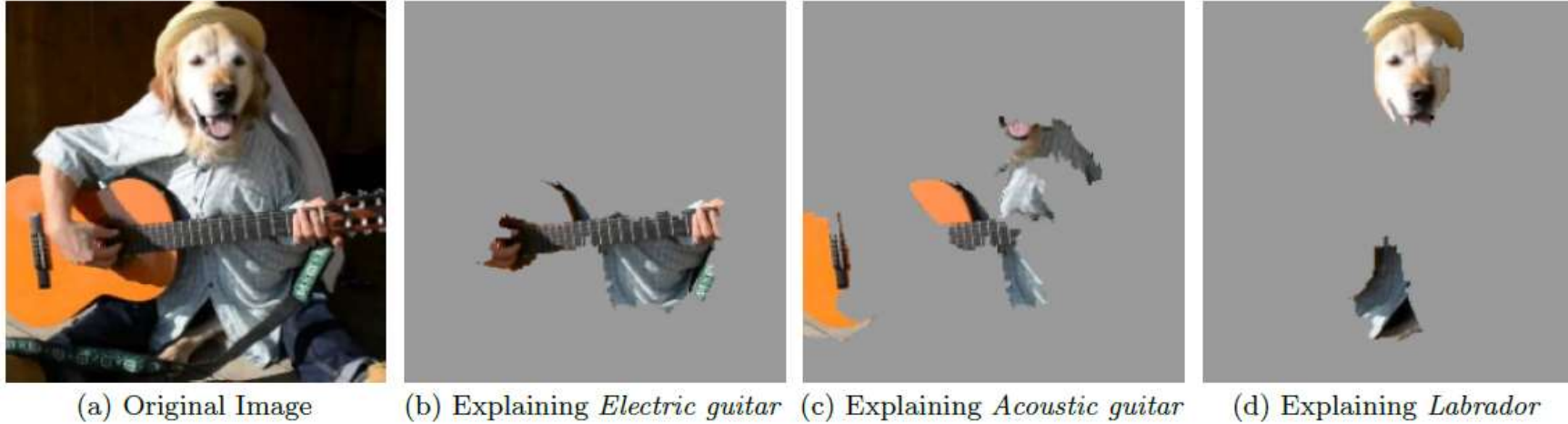
GrLivArea: Above grade (ground) living area square feet

TotalBsmtSF: Total square feet of basement area



**Top 3 is different from linear correlation/LIME.**

**Results of LIME**



(a) Original Image  (b) Explaining *Electric guitar*  (c) Explaining *Acoustic guitar*  (d) Explaining *Labrador*
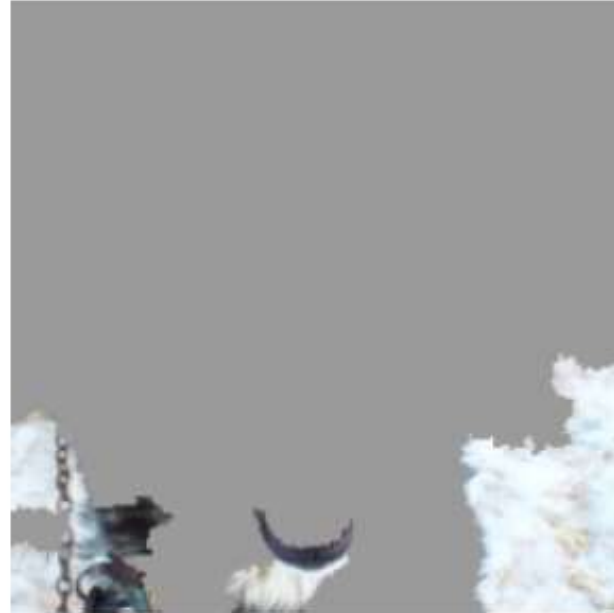
**The model seems to focus on the right places.**

However, there are models which can not be approximated with LIME.
Ex: Classifier whether the image is "retro" or not considering the
values of the whole pixels (sepia?)

**Husky or wolf example**
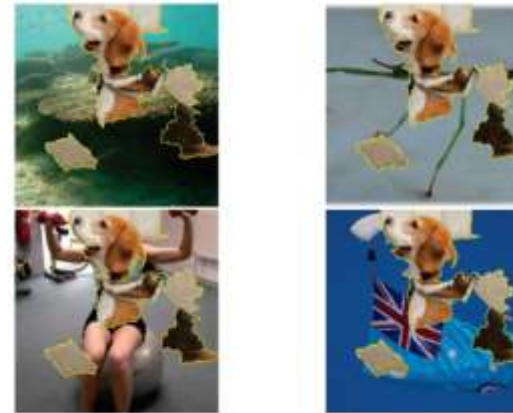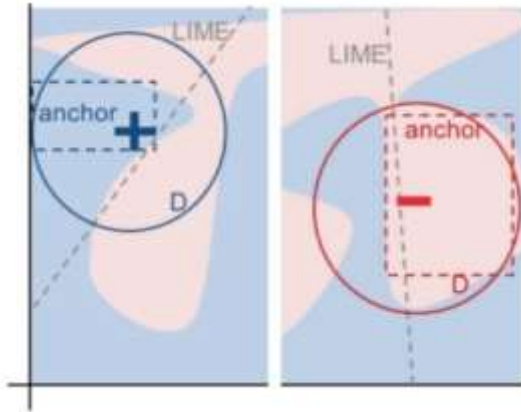


a) Husky classified as wolf    (b) Explanation

**By looking at this explanation,
it is easy to find that the model is focusing on snow.**

## Other approaches

- Anchor M. Ribeiro, et. al., *Thirty-Second AAAI Conference on Artificial Intelligence.* 2018.
  Gives range of features which does not change the prediction



(c) Images where Inception predicts $P(\text{beagle}) > 90\%$

- Influence P. Koh, P. Liang, arXiv:1703.04730
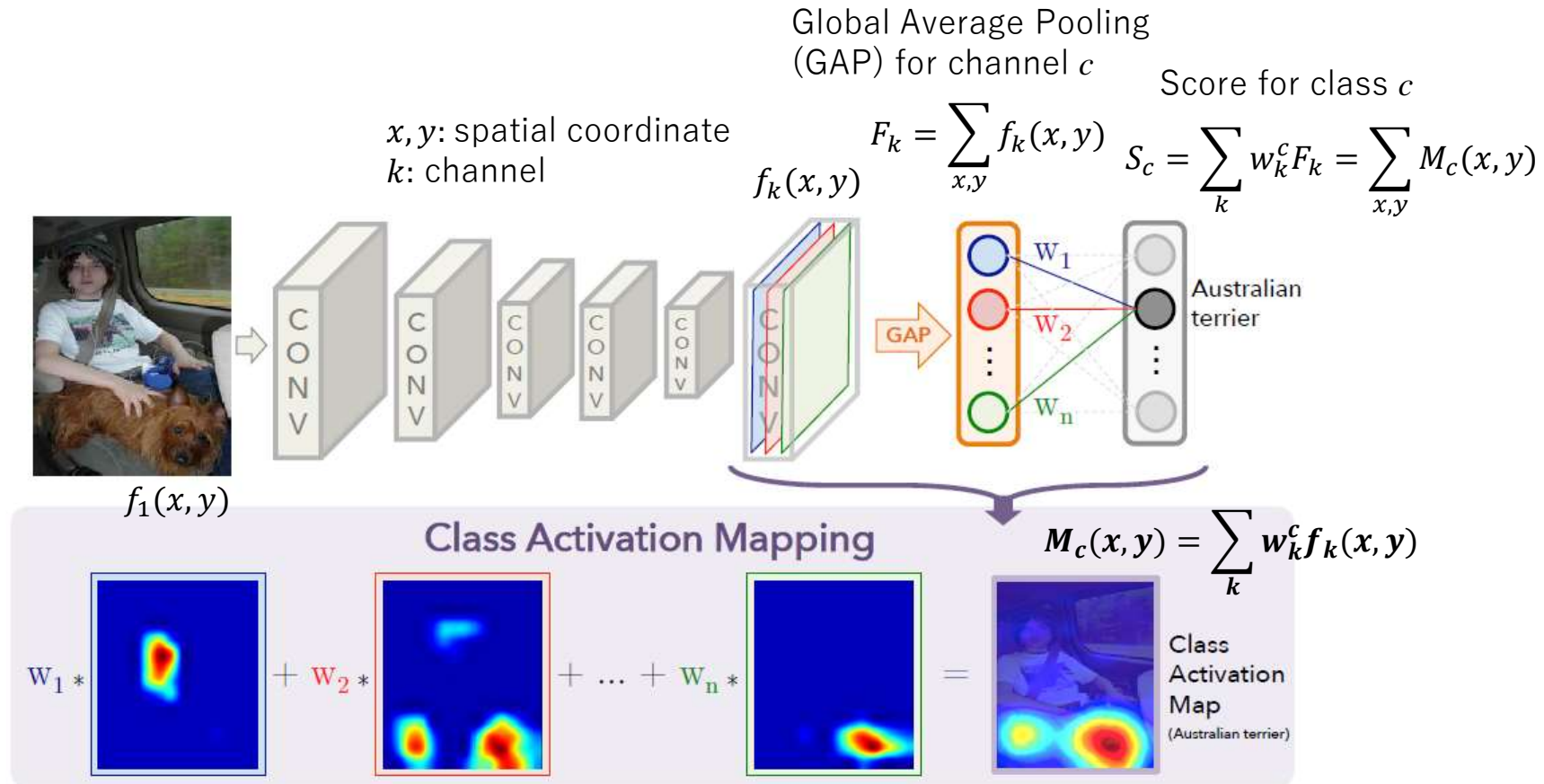  Gives **training data** which the prediction is based on



Test image    Training images which affect the prediction of the test image

# Contents

**Class Activation Mapping (CAM)**

B. Zhou, et. al., arXiv:1512.04150

Global Average Pooling
(GAP) for channel $c$

Score for class $c$

$x, y$: spatial coordinate
$k$: channel

$f_k(x, y)$

$$F_k = \sum_{x,y} f_k(x, y)$$

$$S_c = \sum_k w_k^c F_k = \sum_{x,y} M_c(x, y)$$



$f_1(x, y)$

GAP

Australian terrier

W$_1$

W$_2$

W$_n$

**Class Activation Mapping**

$$M_c(x, y) = \sum_k w_k^c f_k(x, y)$$

W$_1$ * $+$ W$_2$ * $+$ ... $+$ W$_n$ * $=$
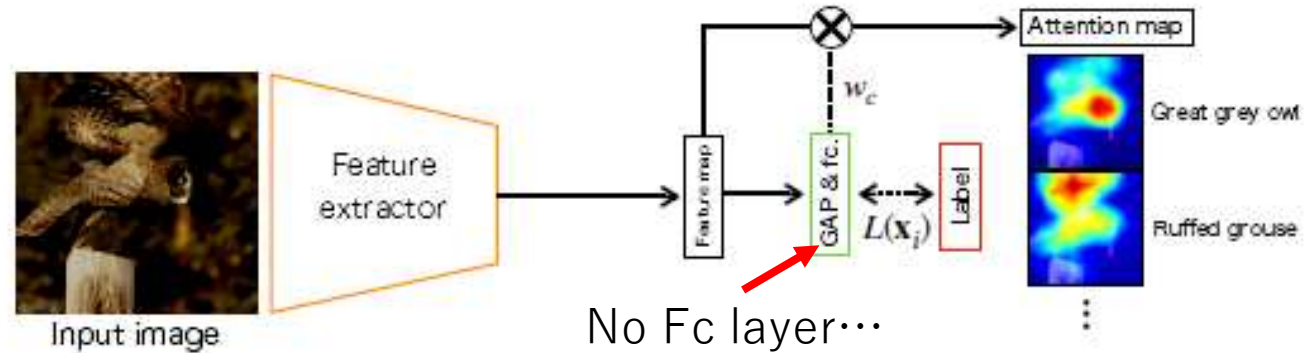
Class Activation Map
(Australian terrier)

☹

- Decrease classification accuracy because fully-connected (Fc) layer is replaced with GAP.
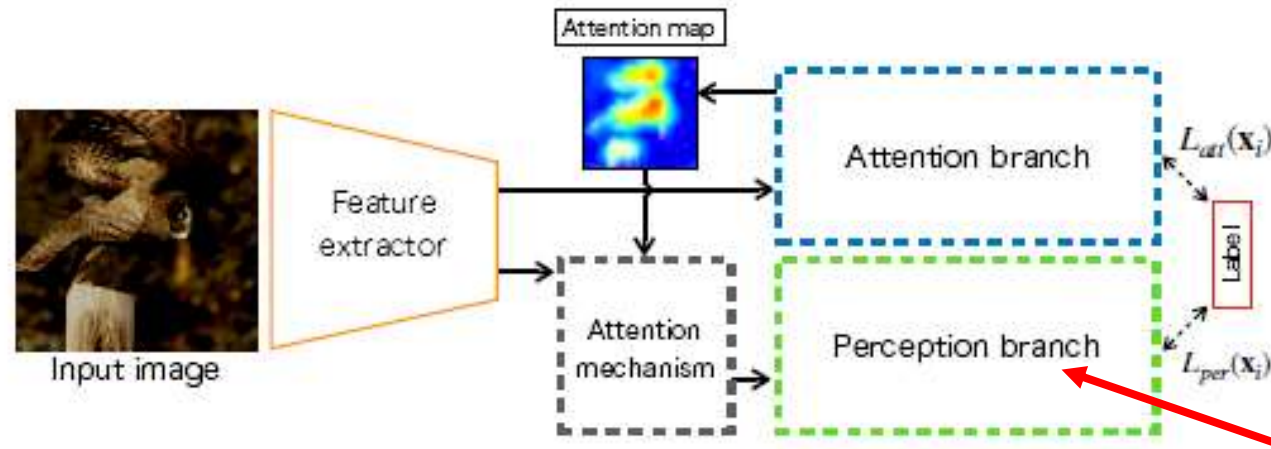
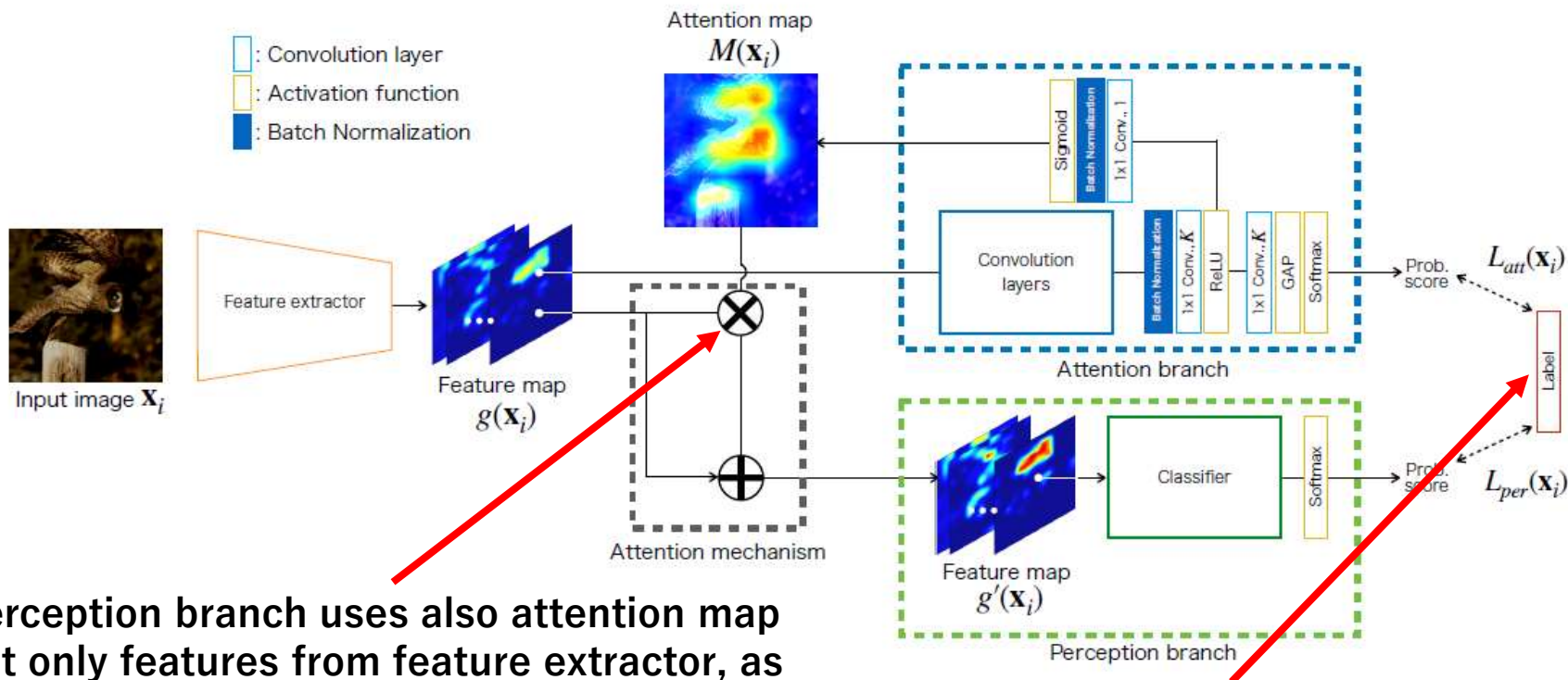**Attention Branch Network（ABN）**  H. Fukui, et. al., arXiv:1812.10025

**Basic idea: Divide attention branch from classification branch so that Fc layers can be used in the latter branch.**



No Fc layer…

(a) Class Activation Mapping

(b) Attention Branch Network

Can add Fc layers!!

**Perception branch uses also attention map not only features from feature extractor, as its input.**

$$g'_c(\mathbf{x}_i) = M(\mathbf{x}_i) \cdot g_c(\mathbf{x}_i)$$

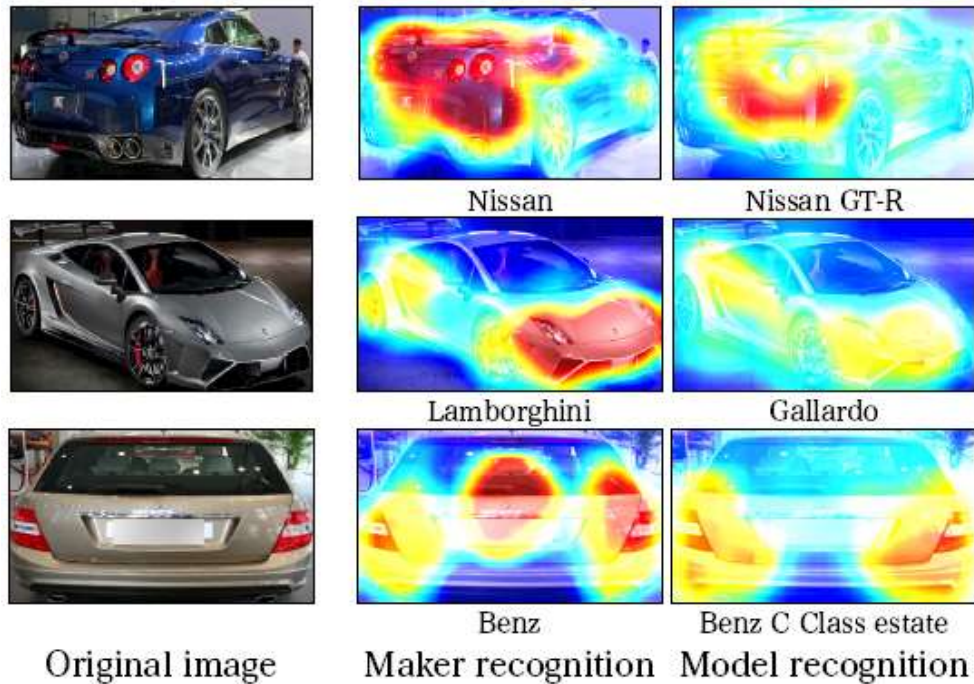**Loss function has attention/perception terms.**
**Both consists of cross entropy for classification task.**

$$L(\mathbf{x}_i) = L_{att}(\mathbf{x}_i) + L_{per}(\mathbf{x}_i)$$

- Improved classification accuracy because it can use Fc layers.
- Actually, using attention map in the input of perception/loss function improves accuracy.
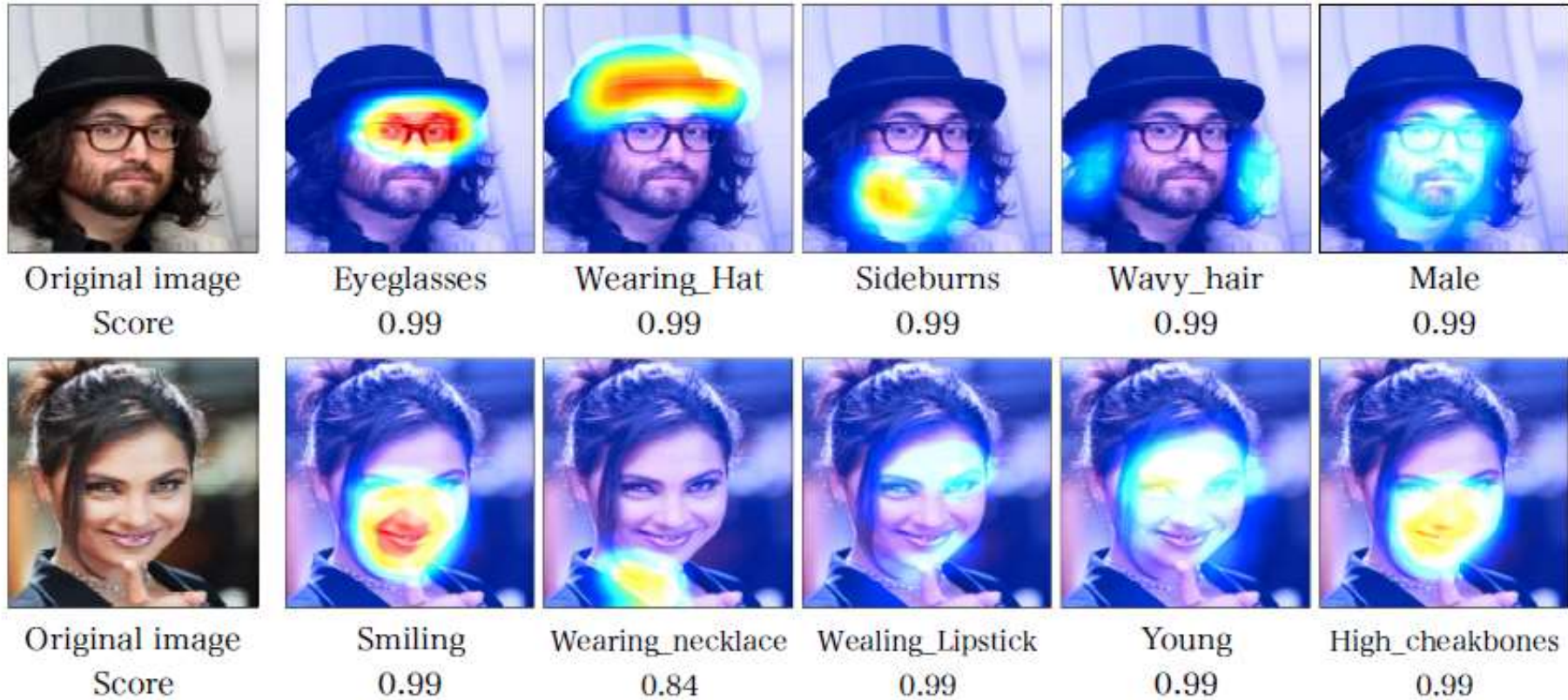
**Results of ABN**



| task | model [%] | maker [%] |
|------|-----------|-----------|
| VGG16 | 85.9 | 90.4 |
| ResNet101 | 90.2 | 90.1 |
| VGG16+ABN | 90.7 | 92.9 |
| ResNet101+ABN | **97.1** | **98.1** |

**Attention map improves accuracy!**

As I do not have domain knowledge,
I can not judge whether the model is focusing on correct features…
Are the highlighted parts characteristic for each maker/model?

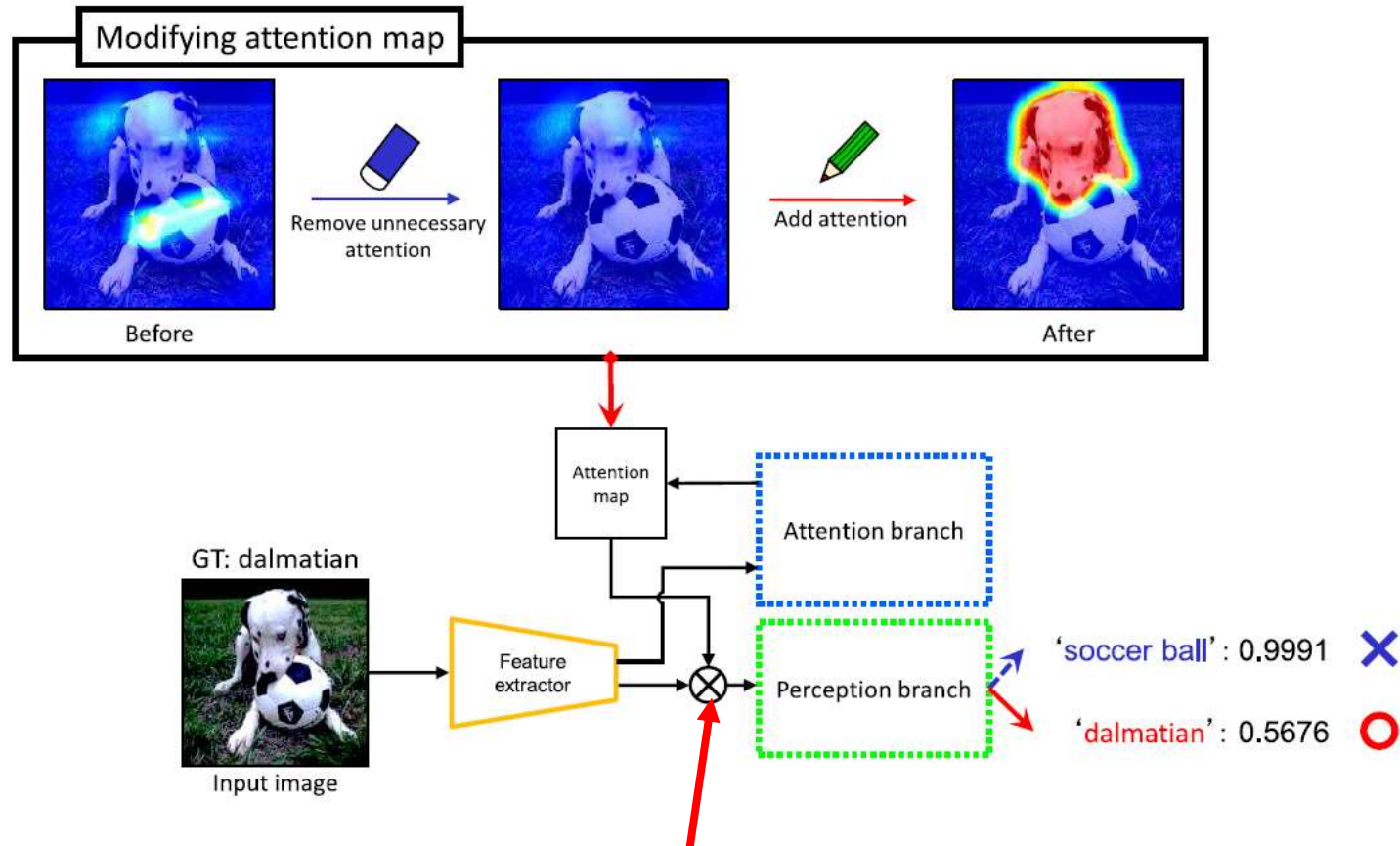**It is interesting that, the attention is paid to different parts depending on the task.**

The model seems to be focusing on correct features.

**ABN and human-in-the-loop**
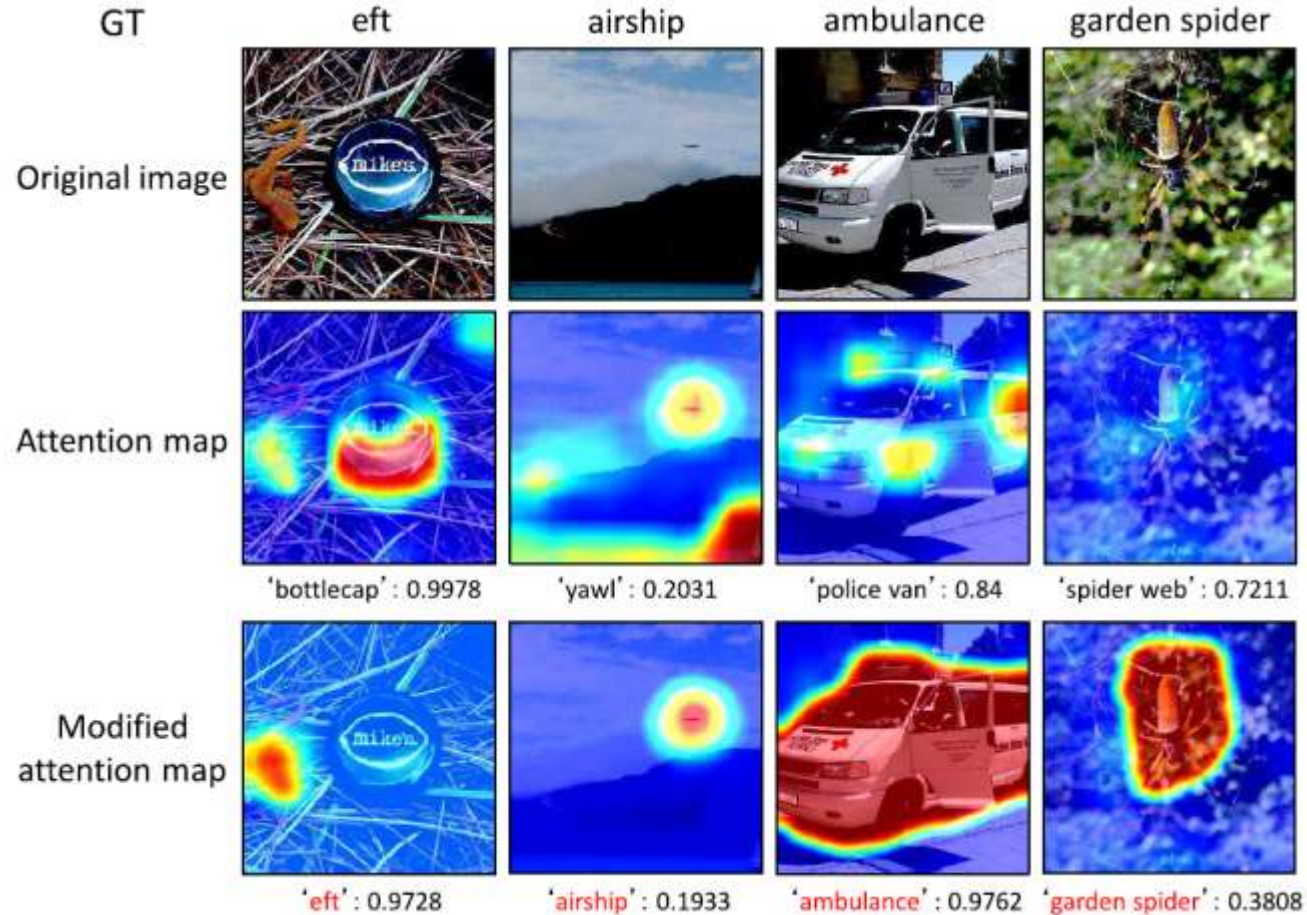
M. Mitsuhara et. al., arXiv: 1905.03540

**Basic idea: By modifying attention map in ABN using human knowledge, try to improve the accuracy of image classifier.**

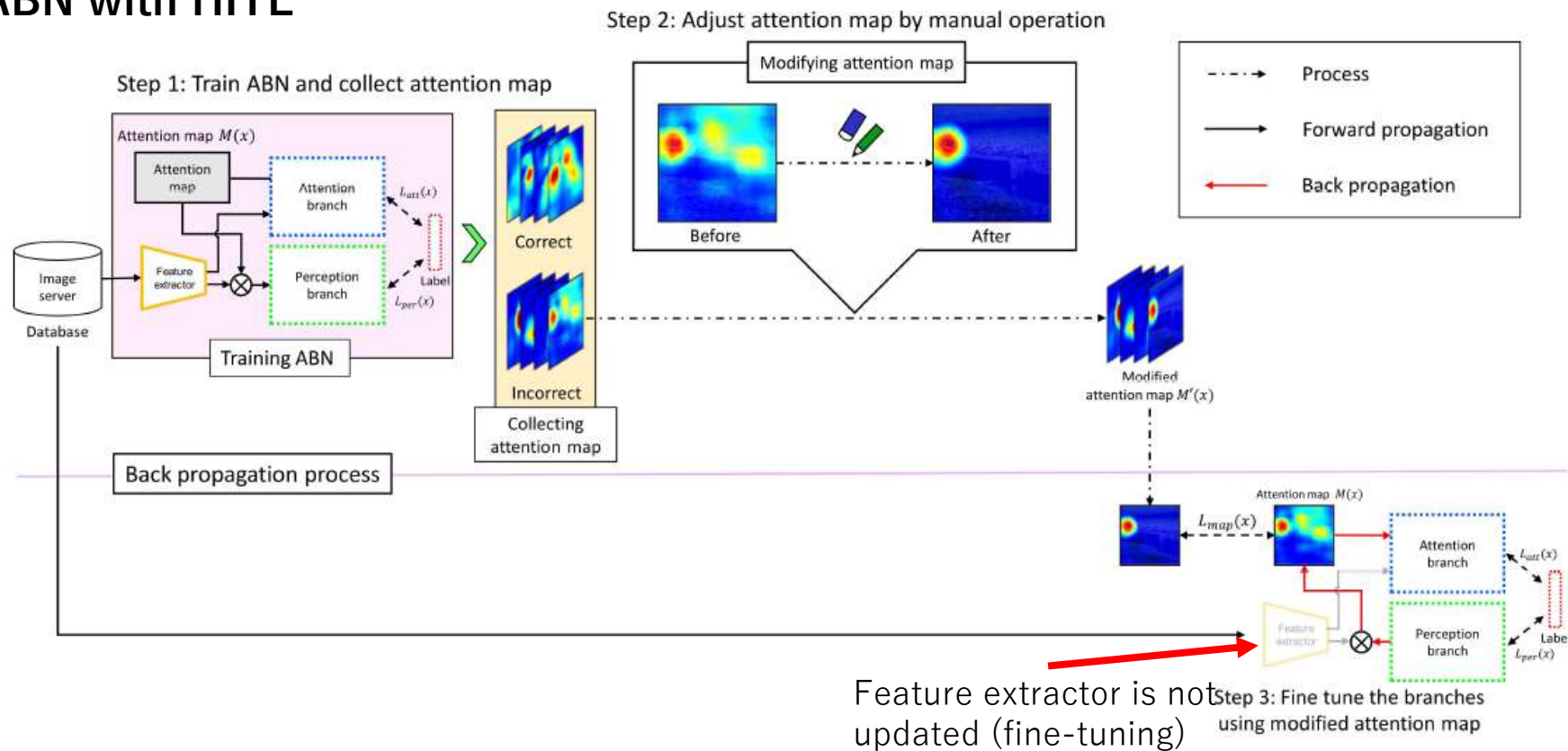**Modifying attention map manually helps classification**



Perception branch uses also attention map not only features from feature extractor, as its input.

**Results**



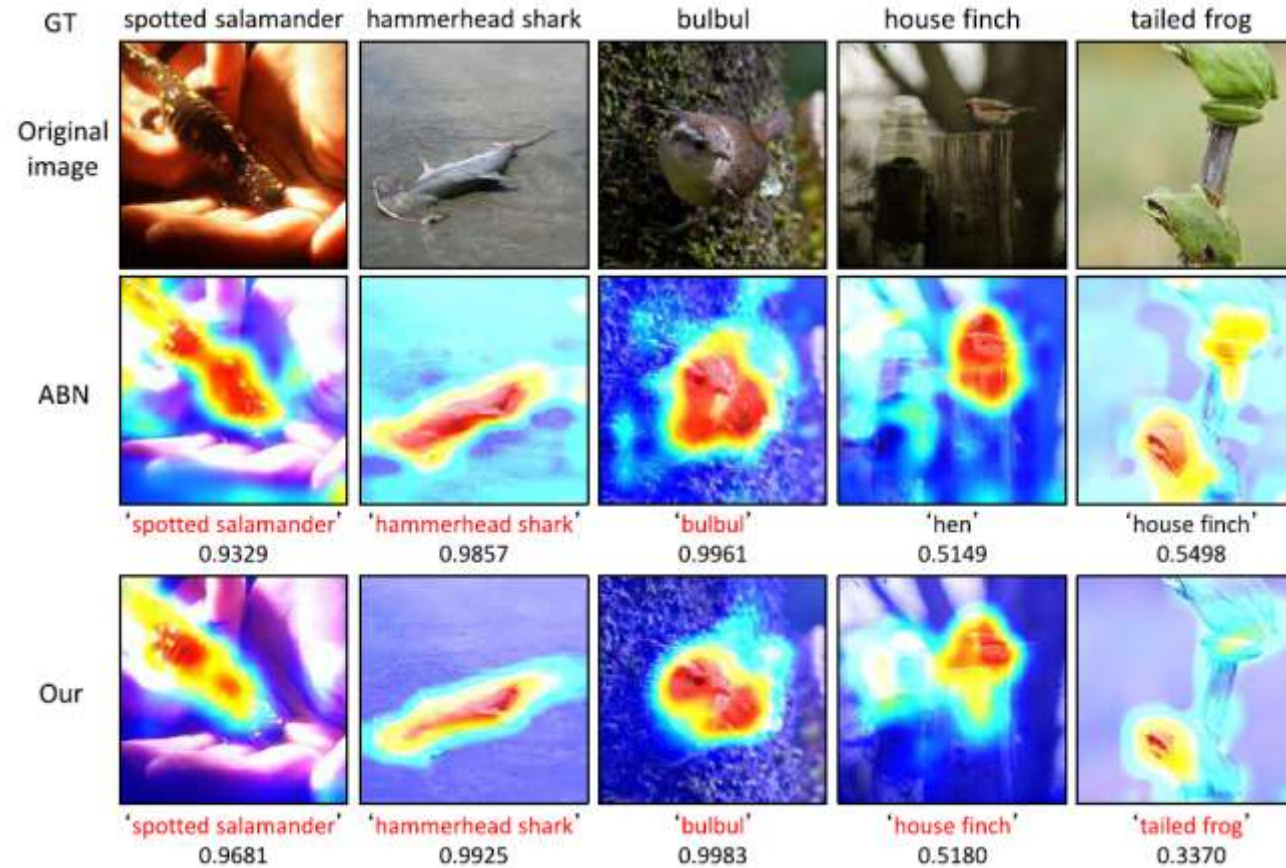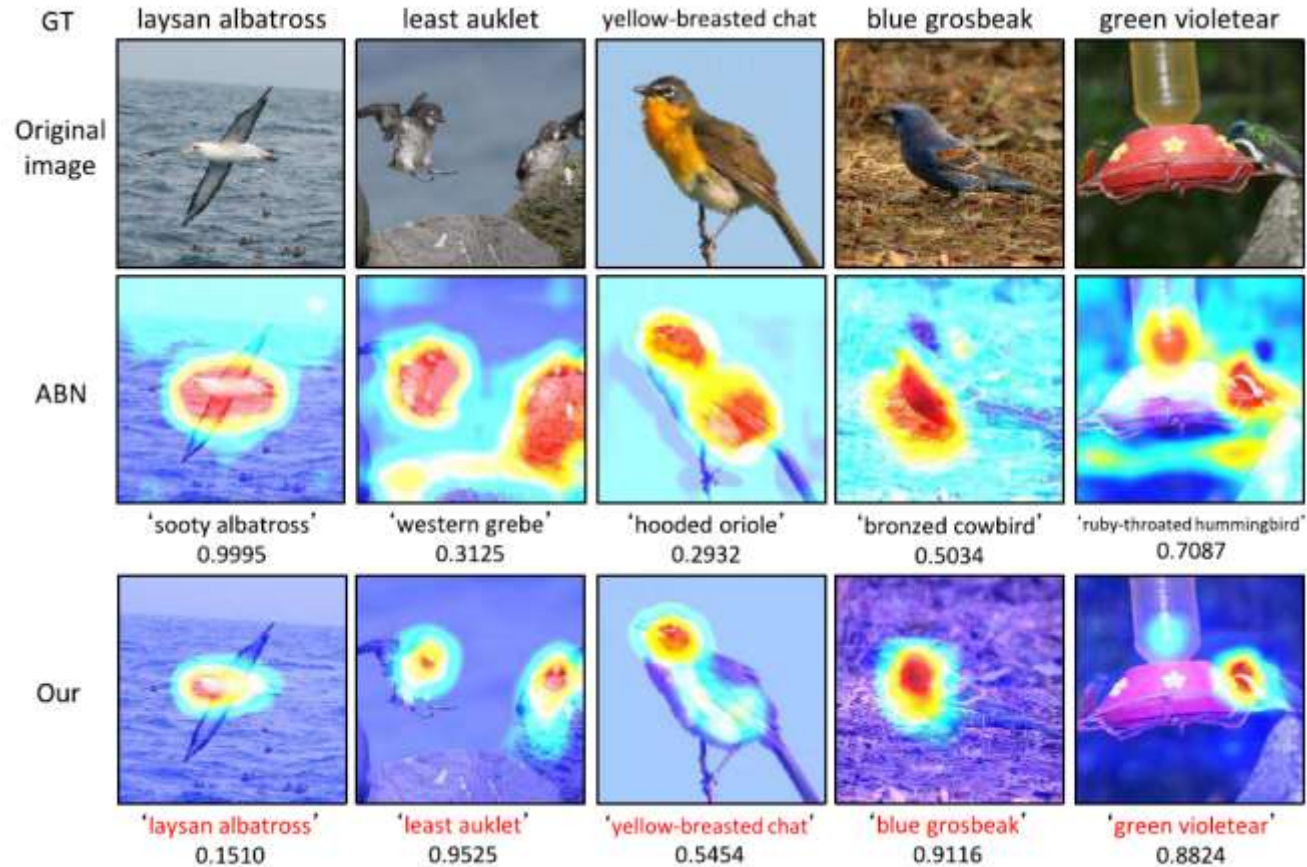**Modifying the attention map improved the classification accuracy**

**ABN with HITL**



$$L_{all}(\mathbf{x}) = L_{att}(\mathbf{x}) + L_{per}(\mathbf{x}) + L_{map}(\mathbf{x})$$

$$L_{map}(\mathbf{x}) = \gamma \|M'(\mathbf{x}) - M(\mathbf{x})\|_2$$

**Results**



**HITL improves the attention map（not very apparent）, and also classification accuracy.**

**HITL improves the attention map（focusing on relevant parts, not the whole body）, and also classification accuracy.**

HITL is possible because the explanation (attention map) is given as a part of their model in ABN, and it is reused as input of perception branch. Unlike LIME, ABN is not model-agnostic.

**So maybe, being model-agnostic is not always useful:
Sometimes it is better to be able to touch the model.**

I introduced a few recent methods on XAI.

1. **LIME**
   for
   - Classification of structured data
   - Regression of structured data
   - Classification of image

2. **ABN**
   for
   - Classification of image

I also introduced **application of human-in-the-loop to ABN**

人間に、愛を。
未来に、AIを。

**Arithmer 株式会社**

〒106-6040

東京都港区六本木一丁目 6 番 1 号 泉ガーデンタワー 38/40F( 受付 )

03-5579-6683

https://arithmer.co.jp/