

Explainable AI summary

I. Papastratis, Prof. Ioannis Pitas
Aristotle University of Thessaloniki
pitas@csd.auth.gr
www.aiia.csd.auth.gr
Version 1.0.1

Explainable AI



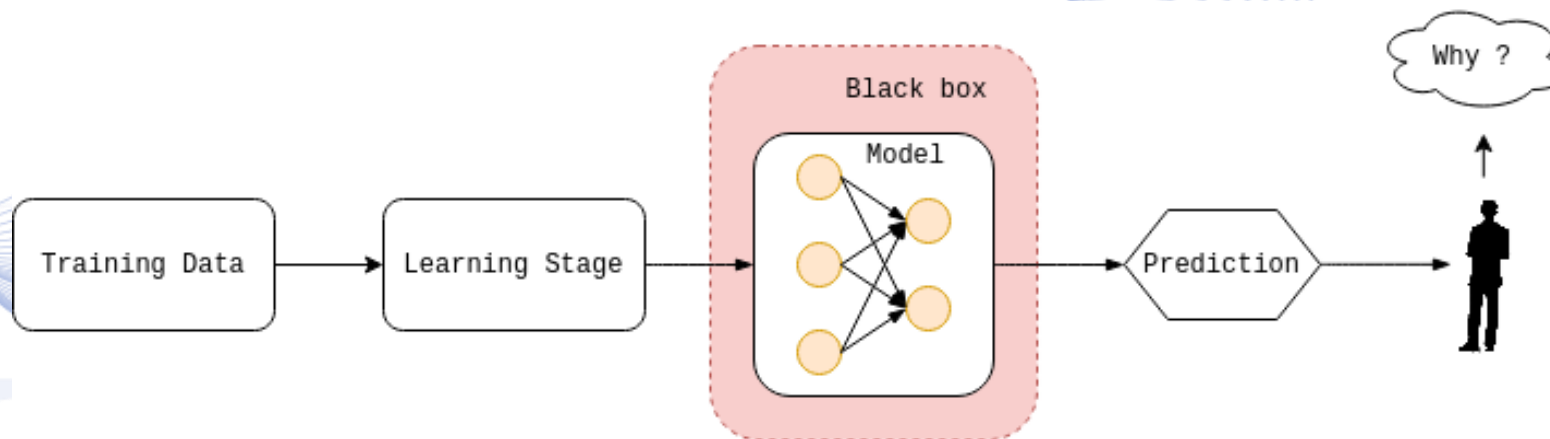
- Introduction to Explainable Artificial Intelligence
- Interpretability
- Types of Interpretability
 - Visual explanations
 - Image-based
 - Plot visualizations
 - Textual explanations
 - Numerical-Mathematical explanations
- Applications
- Frameworks

Introduction to Explainable AI

Machine and Deep Learning have surpassed humans in many tasks (image and speech recognition, recommendation systems, medical diagnosis)

Drawbacks:

- Machine learning (ML) architectures are usually considered as blackbox models



Introduction to Explainable AI



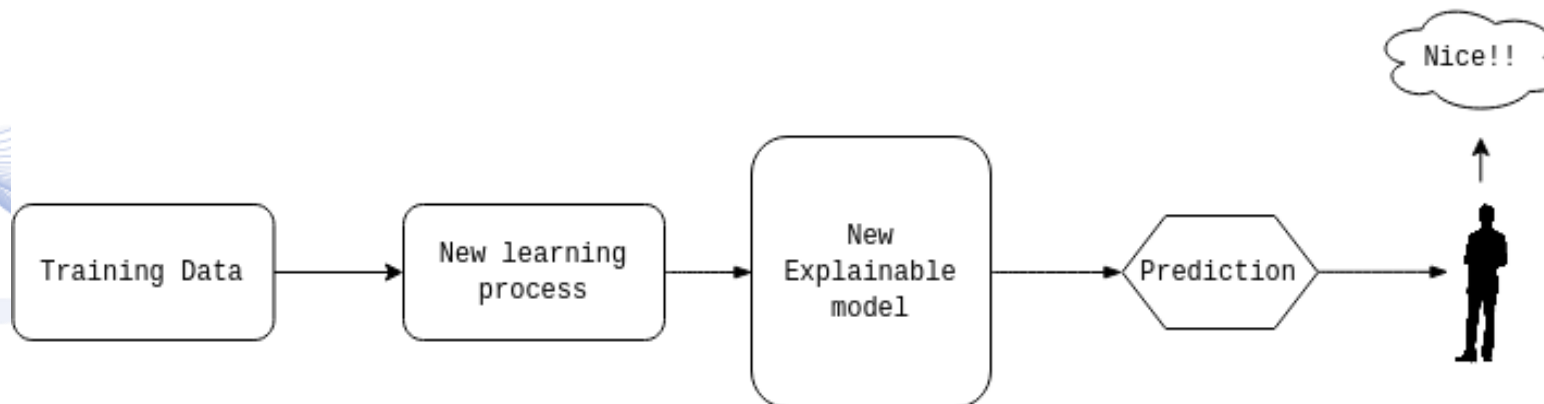
- Deep learning (DL) have achieved outstanding performance but they don't justify their reliability
- Failure
 - An error in any moment of a self-driving car can lead to a fatal crash.
 - in the medical area, human lives may be dependent on these decisions

Introduction to Explainable AI

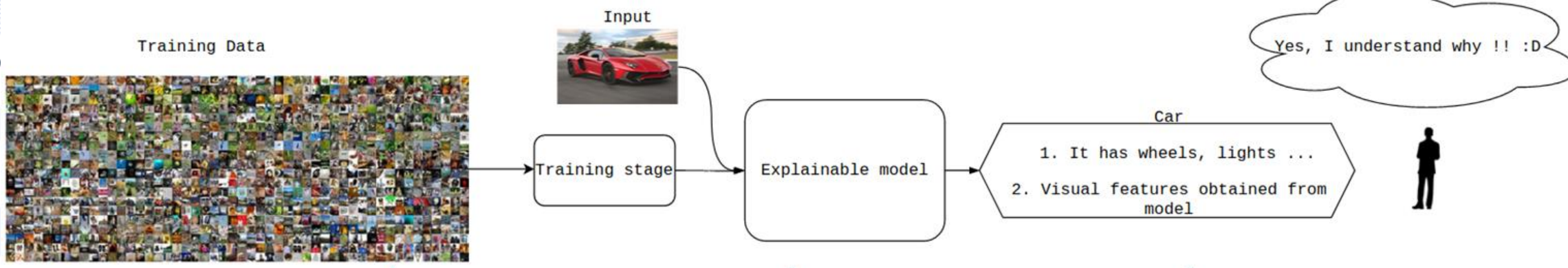
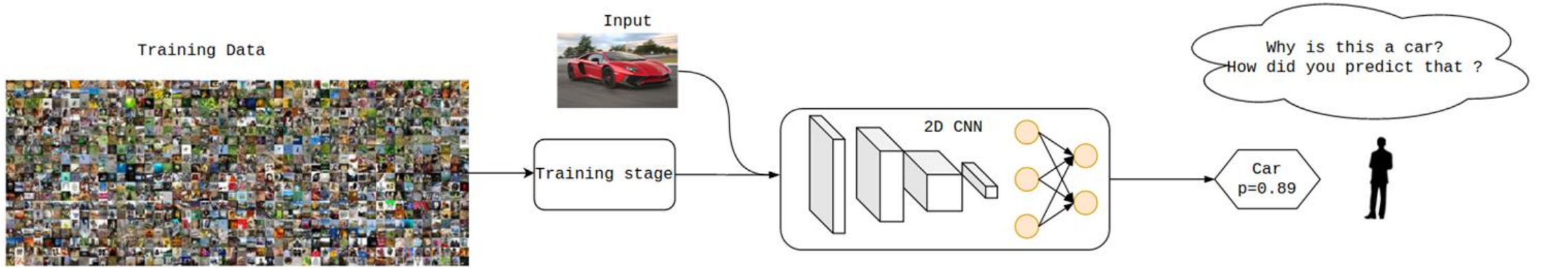
Issues

- responsibility for bad AI decisions
- explain errors of AI decisions
- improvement of AI models

Solution: Explainable AI models easily understandable by humans



Introduction to Explainable AI



Interpretability

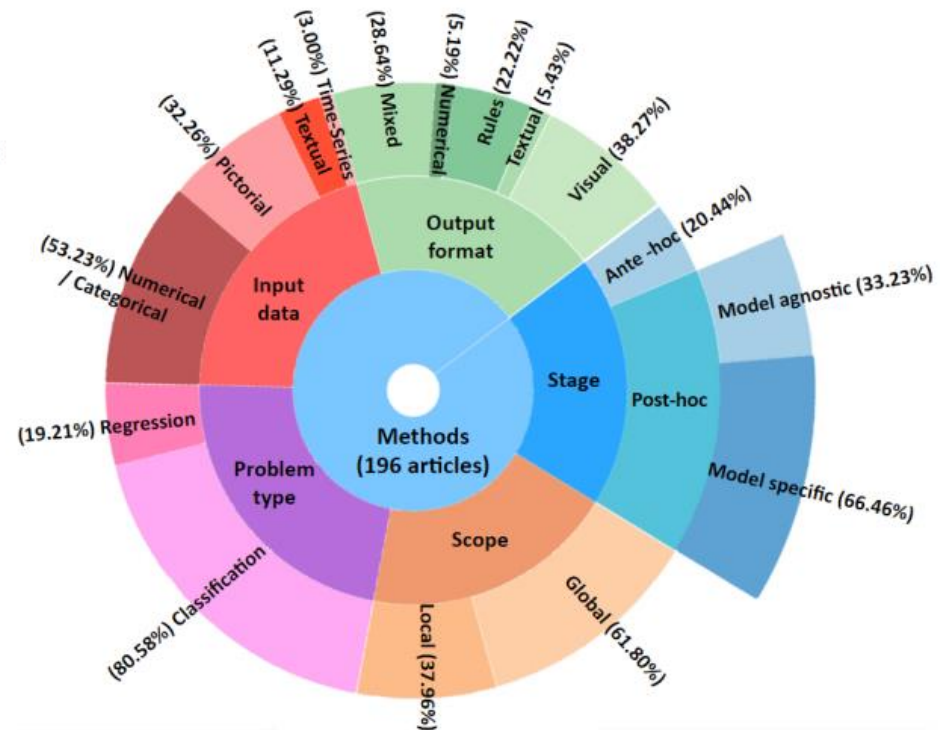
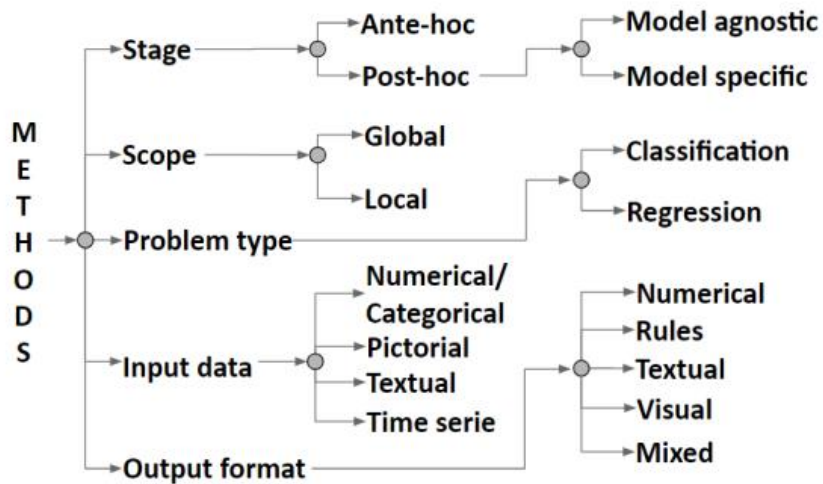


How easily we can understand the cause of an algorithm's decision or action
The categorization of Interpretability methods is based on how interpretable information is provided.

Categories:

- Visual Interpretability (“obviously” interpretable information, easily perceived from human eye)
- Textual Explanations (Given in form of text)
- Mathematical-Numerical Explanations

Types of Interpretability Methods



[Vilone2020] Categorization of explainability methods

Visual Explanations

- Visual explainable methods produce pictures or maps in order to provide information about the model's decision
 - Most common: **Saliency** methods explain results of model by producing outputs to show which components are responsible.
 - These values take the form of output probabilities or images like heatmaps.
- Plot visualization methods produce scatter plots to explain decisions or visualize the data

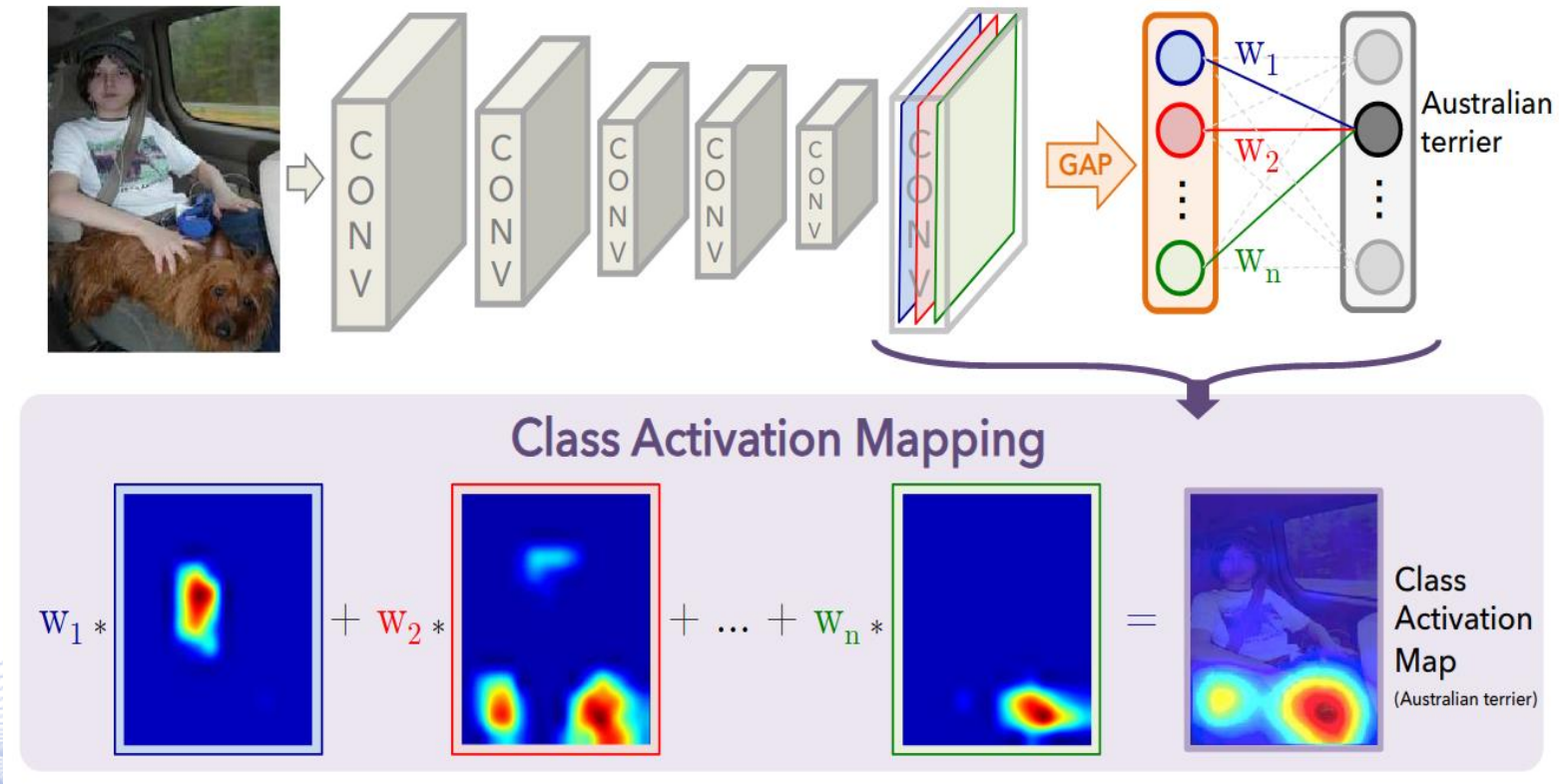
Visual Explanations

- Methods
 - CAM (Class Activation Maps)
 - Grad-CAM: Gradient based CAM
 - LRP (Layer-wise Relevance Propagation)
 - Peak Response Map (PRM)
 - CLass-Enhanced Attentive Response (CLEAR)
 - DeConvNet
 - DeepResolve
 - SCOUTER

Class Activation Maps (CAM)

- Generate activation maps for a decision
- Global average pooling layer is after convolutional layers
- The output features of the convolutional neural network are passed through a fully-connected layer that makes the prediction
- CAM indicates the region on the image that correspond to the prediction result

Class Activation Maps(CAM)



[Zhou2016] Overview of CAM

Class Activation Maps(CAM)



[Zhou2016] Examples of CAM

Grad-CAM: Gradient-based CAM



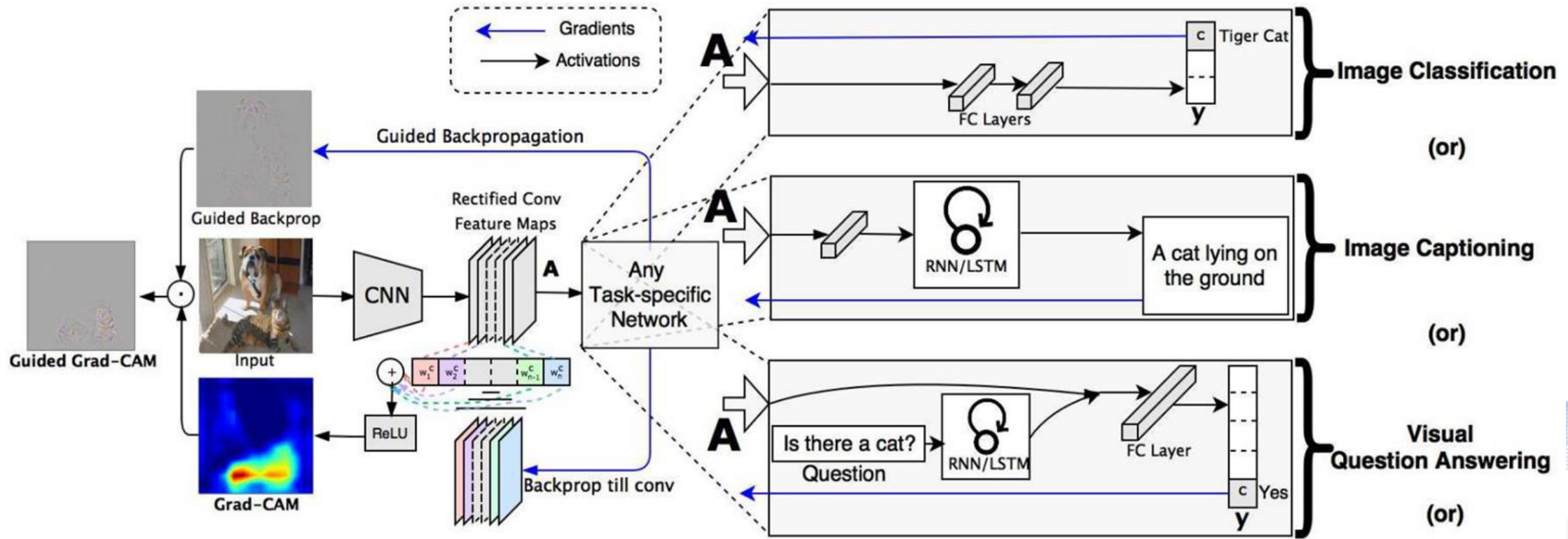
Gradient-weighted Class Activation Maps(Grad-CAM)

- an extended version of CAM by computing the gradients with respect to the target that flow to the final convolutional layer
- produces a map, which highlights the most useful pixels for classification

Processing steps

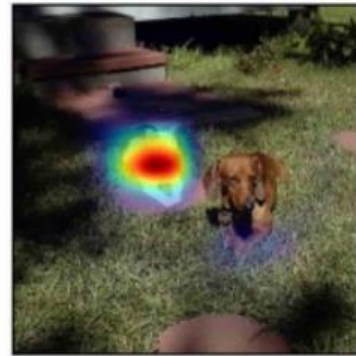
- Forward pass of the input image to produce the prediction
- Gradients of the target class
- The gradients of the target are back-propagated to last convolutional layer
- Find the important locations of the image

Grad-CAM: Gradient-based CAM

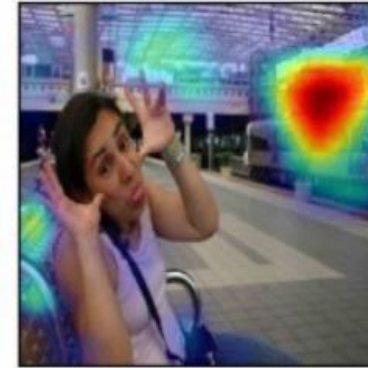


[Selvaraju2017] Overview of Grad-CAM

Grad-CAM: Gradient-based CAM



cat



train



dog



boat

[Selvaraju2017] Examples of Grad-CAM

Layer-wise Relevance Propagation (LRP)



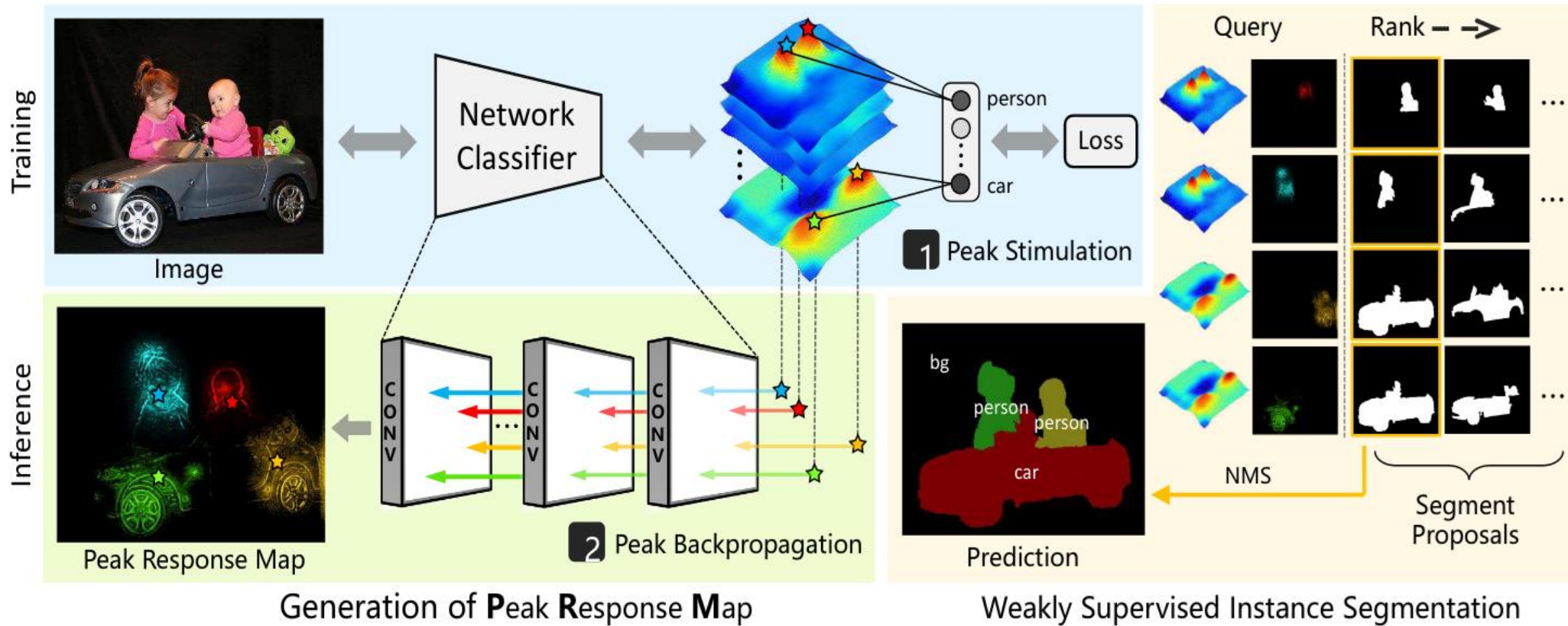
- Decomposition of the classification decision of a DNN model in using the input
- The classification layer is decomposed into several layers
- Backward pass to produce the pixel-wise contributions to the final output from last to first layers
- We denote as $w(i, j)$ the weight between of the connection from neuron i to neuron j

Layer-wise Relevance Propagation (LRP)



[Samek2016] LRP calculation for the input image

Peak Response Map (PRM)



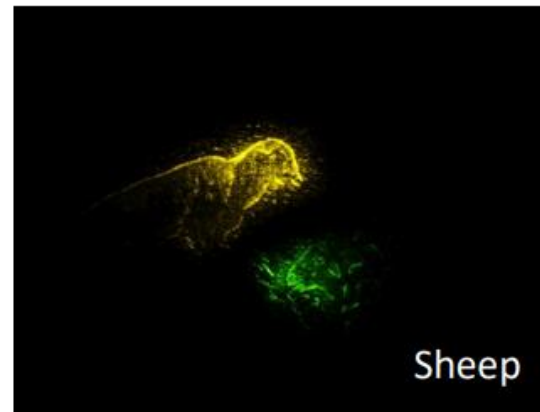
[Zhou2016] Overview of PRM

Peak Response Map (PRM)

Image



Peak Response Map



[Zhou2016] Examples of PRM

Class-Enhanced Attentive Response (CLEAR)



- Visualizes the decisions of image classification applications with attention maps produced by back-propagating the activations of the last layer
 - After forward pass using deconvolutions we obtain the deconvolved output

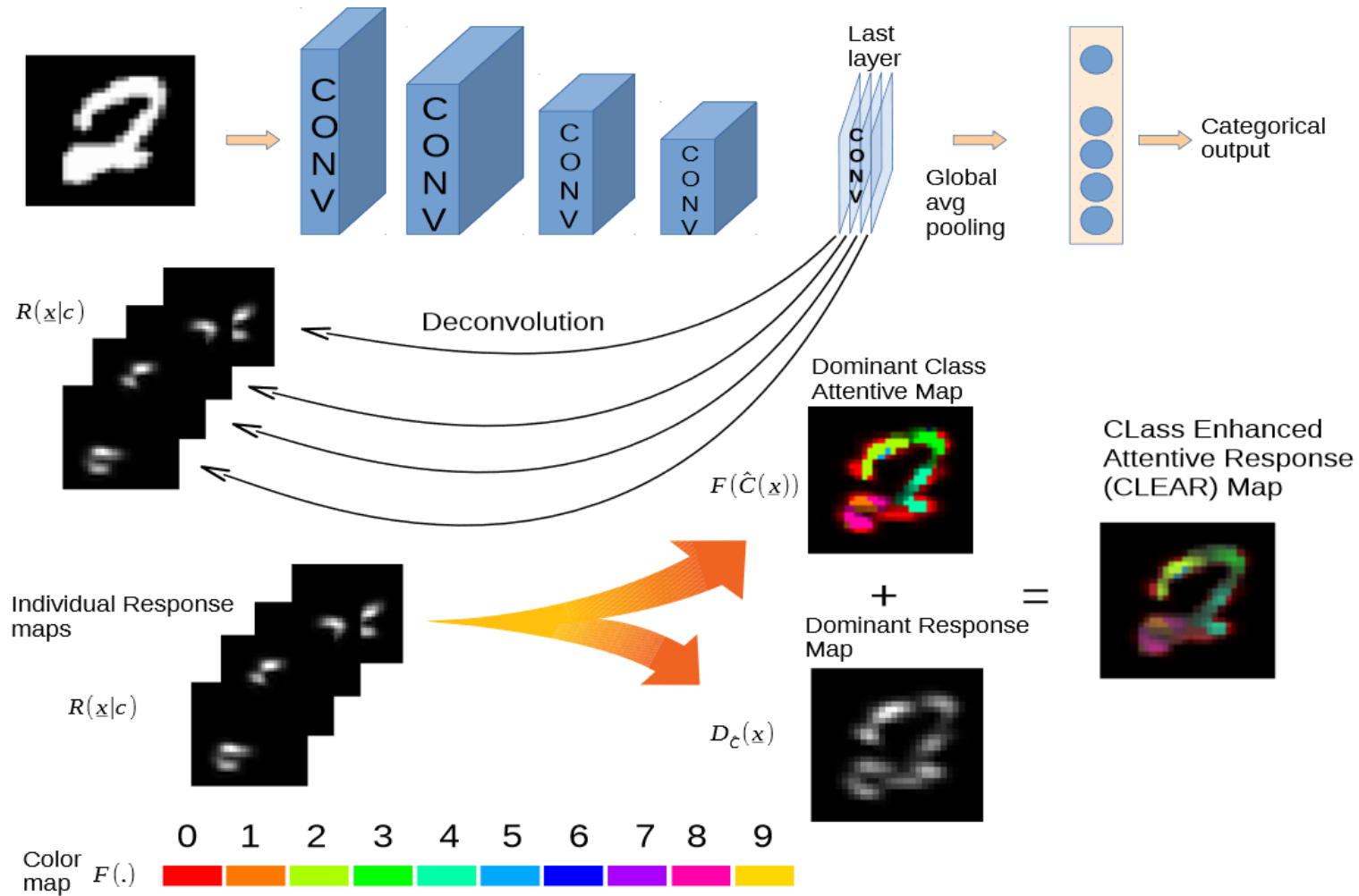
$$\mathbf{h}(l) = \sum_{k=1}^K z(k, l) * w(k, l)$$

k : kernel index $\mathbf{z}(l)$: feature maps of layer l , $\mathbf{w}(l)$: kernel weights, K kernels

- Final response of layer l is the product: $\mathbf{R}(l) = \mathbf{h}(1)\mathbf{h}(2) \dots \mathbf{h}(l)$
- Compute individual attention maps $\mathbf{R}(\mathbf{x}', c)$ of class c and back-projected input \mathbf{x}' from all L layers of the deconvolutional network as :

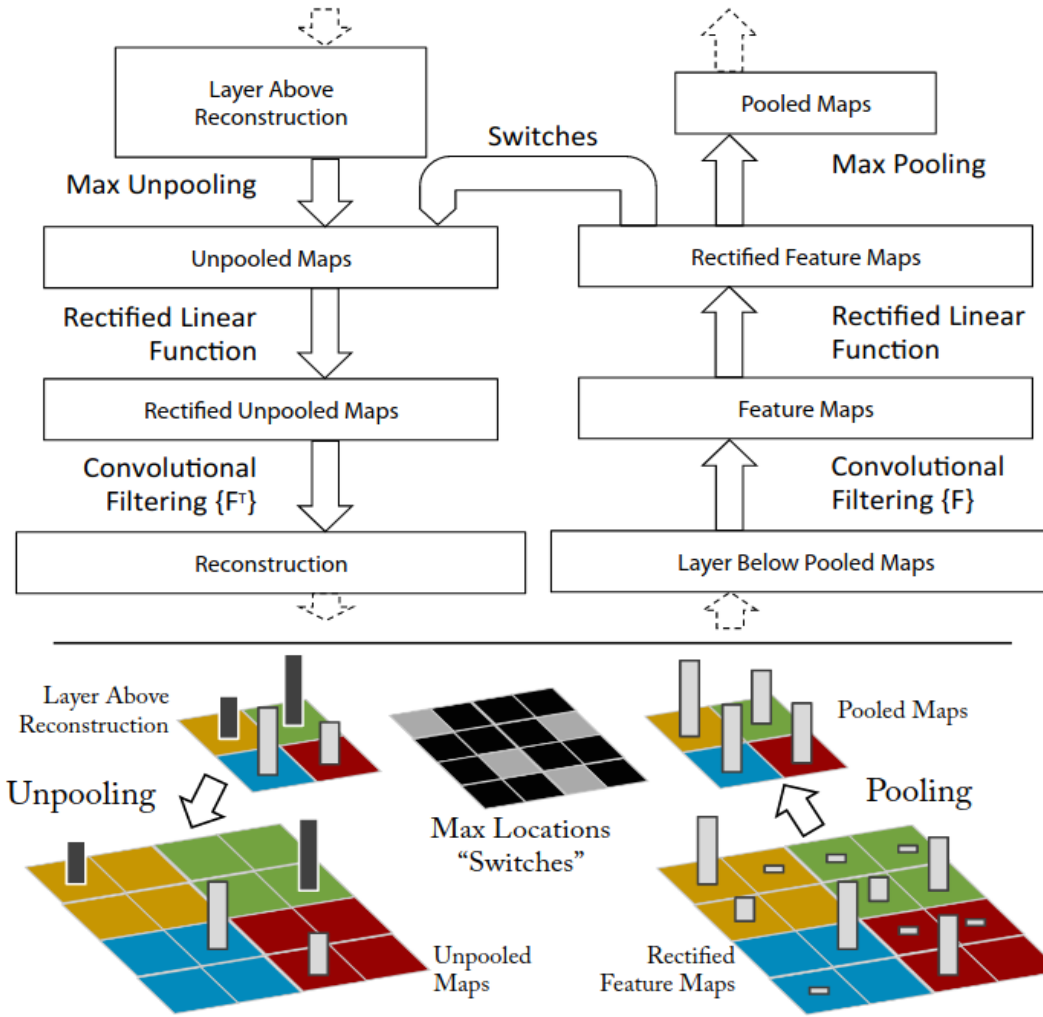
$$\mathbf{R}(\mathbf{x}', c) = \mathbf{h}(1)\mathbf{h}(2) \dots \mathbf{h}(L)$$

Class-Enhanced Attentive Response (CLEAR)



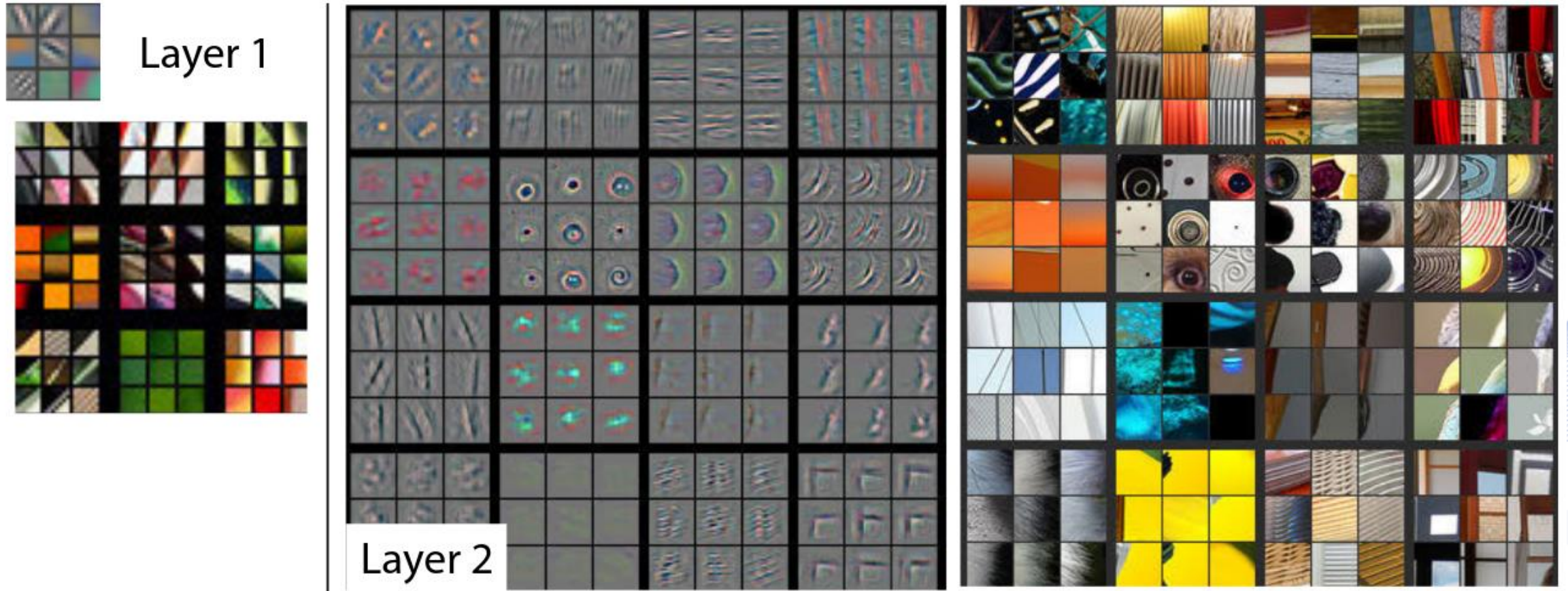
- A deconvnet is actually a convnet model with the same convolutional layers but in reverse since it maps pixels to features
- Deconvnet is used here to generate the input images from features
- Following operations are adopted:
- Trnsnpose Convolution (deconvolution):
 - Use transposed learned filters of cnn to reconstruct the input of the layer from the output
- Unpooling
 - record maximum values locations during maxpooling and use the locations to reconstruct input during deconvolution

DeConvNet



[Zeiler2014] Reconstruction of the convnet features up to the pixel space

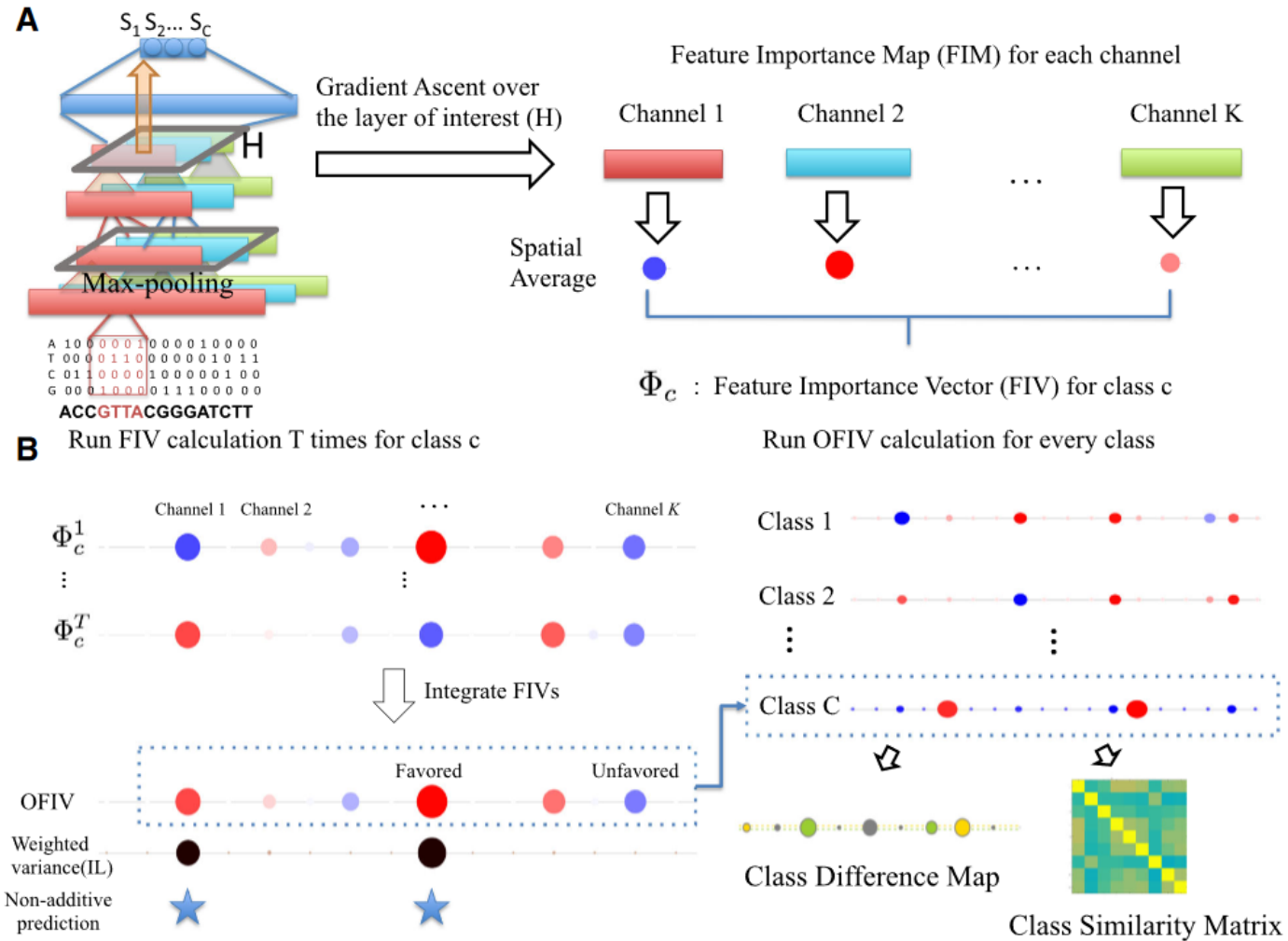
DeConvNet



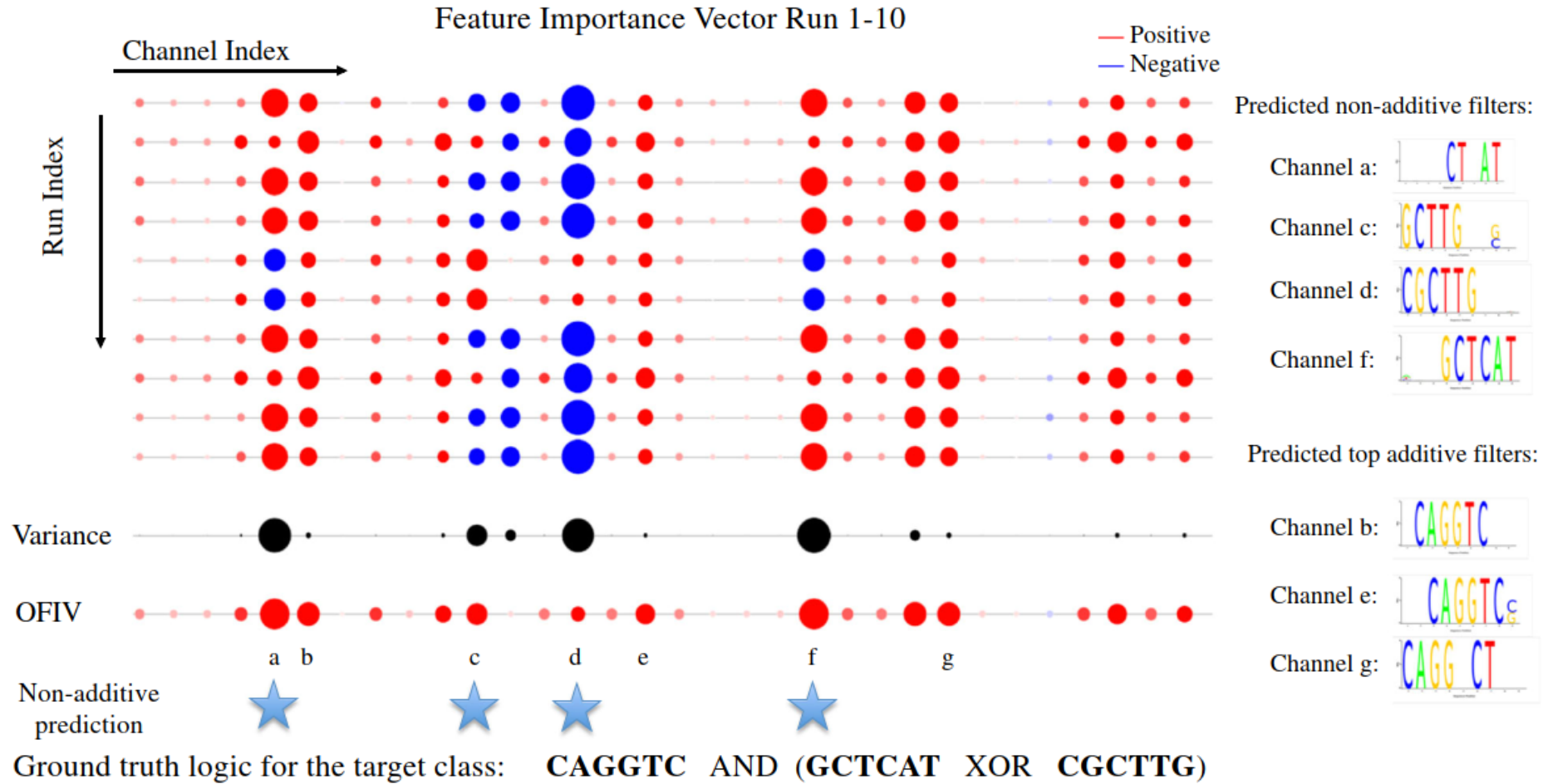
[Zeiler2014] Visualization of intermediate layers

- Generate intermediate layer heatmaps to show how the network combines features for classification
- Compute optimized feature map $\mathbf{H}_c = \operatorname{argmax}_{\mathbf{H}} S_c(\mathbf{H}) - \lambda \|\mathbf{H}\|_2^2$, S_c is the score of class c obtained from the last layer, λ tunable hyperparameter and $\mathbf{H} \in \mathbf{R}^{K \times W}$
- Global average of Feature Importance Maps (FIM) \mathbf{H}_c to obtain feature importance vector (FIV) $\Phi_c = (\varphi_c^1, \dots, \varphi_c^k)$ where $\varphi_c^k = \frac{1}{W} \sum_{i=1}^W (H^k(i))_c$
- This procedure is ran T times with different initial parameters to get several estimations of \mathbf{H}_c^t and Φ_c^t

DeepResolve



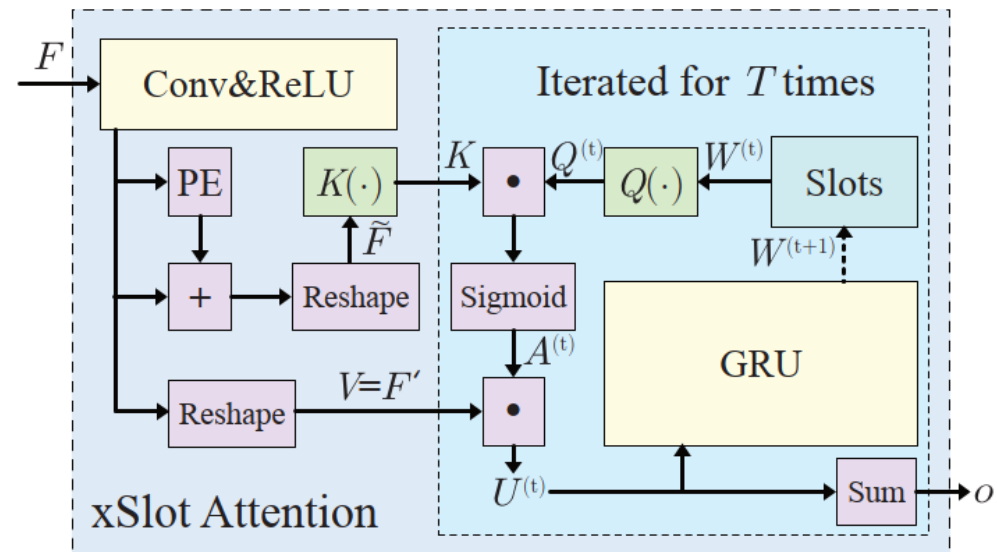
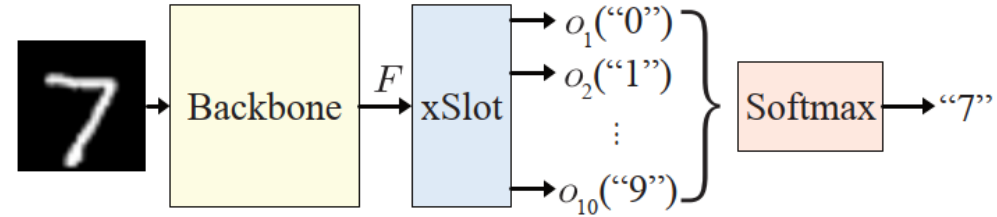
DeepResolve



[Liu2019] Generation of importance vectors

- Explain the reasons an image is classified or not to a specific class
- The cnn's classifier (fully connected layer) is replaced with a slot attention model
- Each slot produces a confidence score for each class
- The cnn's features \mathbf{F} are passed through a convolutional layer and a position embedding to model the spatial information
- Then, a self-attention mechanism is utilized to compute the weighted sum of features as: $\mathbf{A}^{(t)} = \sigma(Q(\mathbf{W}^{(t)})K(\mathbf{F}))$
- σ : sigmoid function Q, K : fully-connected networks $\mathbf{W}^{(t)}$: slot weights

SCOUTER



[Li2020] SCOUTER overview

SCOUTER



[Li2020] Examples of positive and negative decisions

Visual explanation of deep neural networks by interpretation



- Identification of relevant features to the predictions of the network F
- Forward-propagation of N training images to obtain m -dimensional features

$$\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2 \dots \mathbf{x}_N) \in \mathbb{R}^{N \times m}$$

- Target labels $\mathbf{L} = (\mathbf{l}_1, \mathbf{l}_2, \dots \mathbf{l}_N) \in \mathbb{R}^{C \times N}$

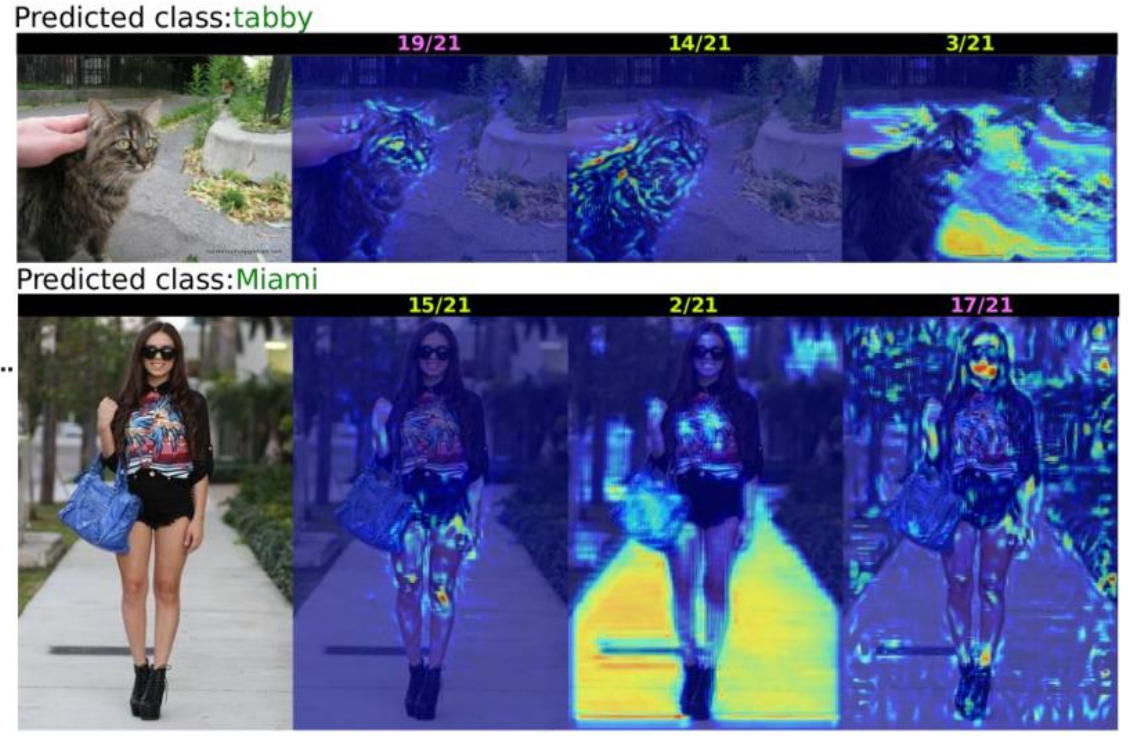
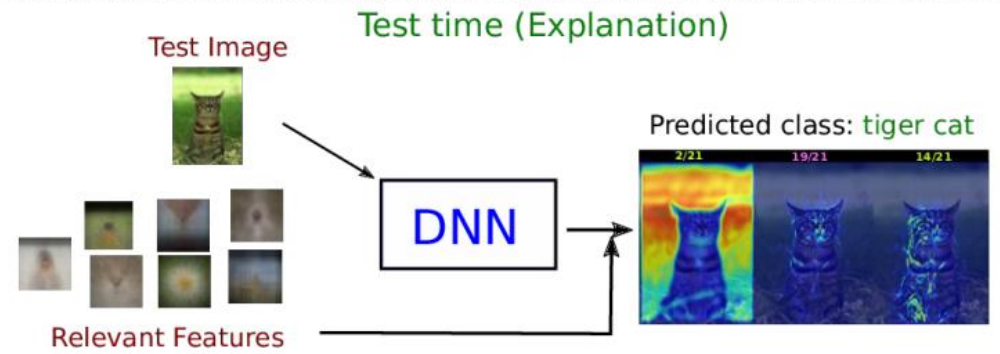
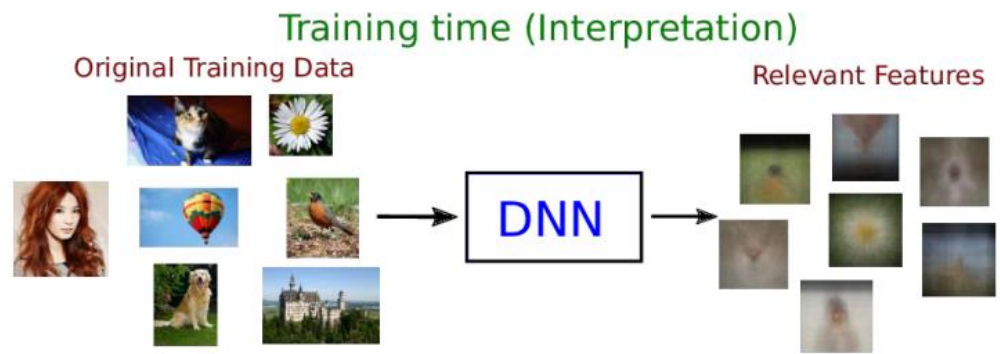
- Predict a linear combination of activations \mathbf{X} for each class using

$$\mathbf{W} = (\mathbf{w}_1, \mathbf{w}_2, \dots \mathbf{w}_N) \in \mathbb{R}^{C \times m}$$

by solving the optimization problem :

$$\mathbf{W}^* = \underset{\mathbf{W}}{\operatorname{argmin}} \left\| \mathbf{X}^T \mathbf{W} - \mathbf{L}^T \right\|_F^2$$

Visual explanation of deep neural networks by interpretation

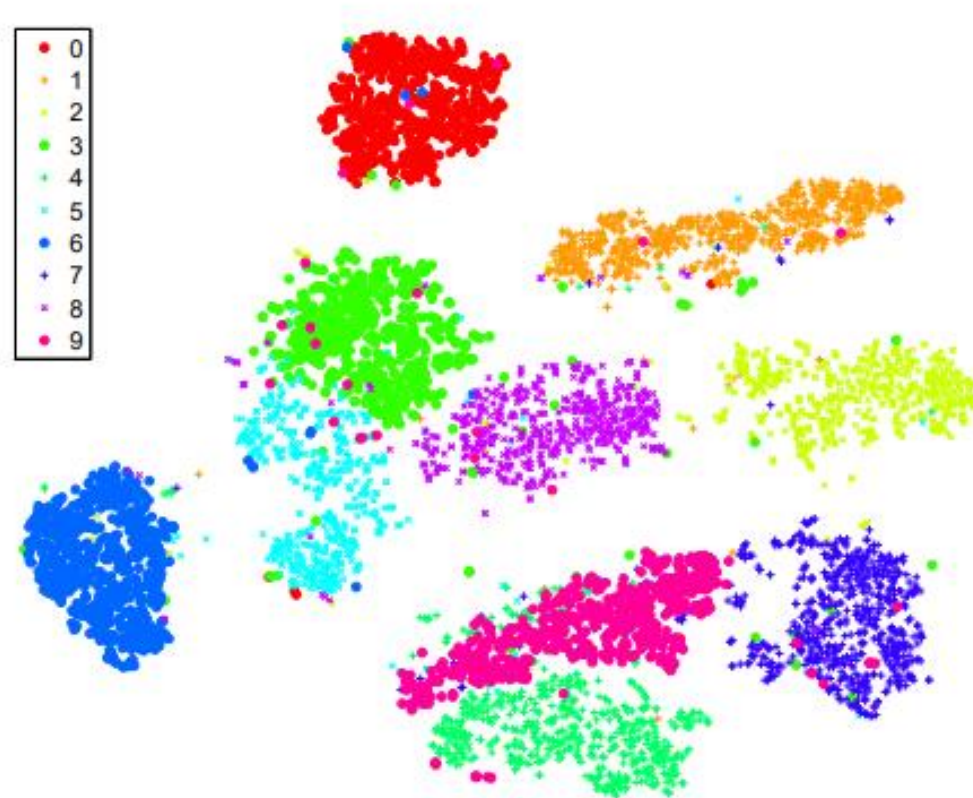


[Oramas2019] Overview of training and testing pipeline

Plot visualizations

- Produce scatter plots to explain decisions or visualize the data
- Methods
 - t-SNE
 - Understanding deep features using PCA
 - Visualization of hidden layers
 - TreeView

t-distributed stochastic neighbor embeddings (t-SNE)

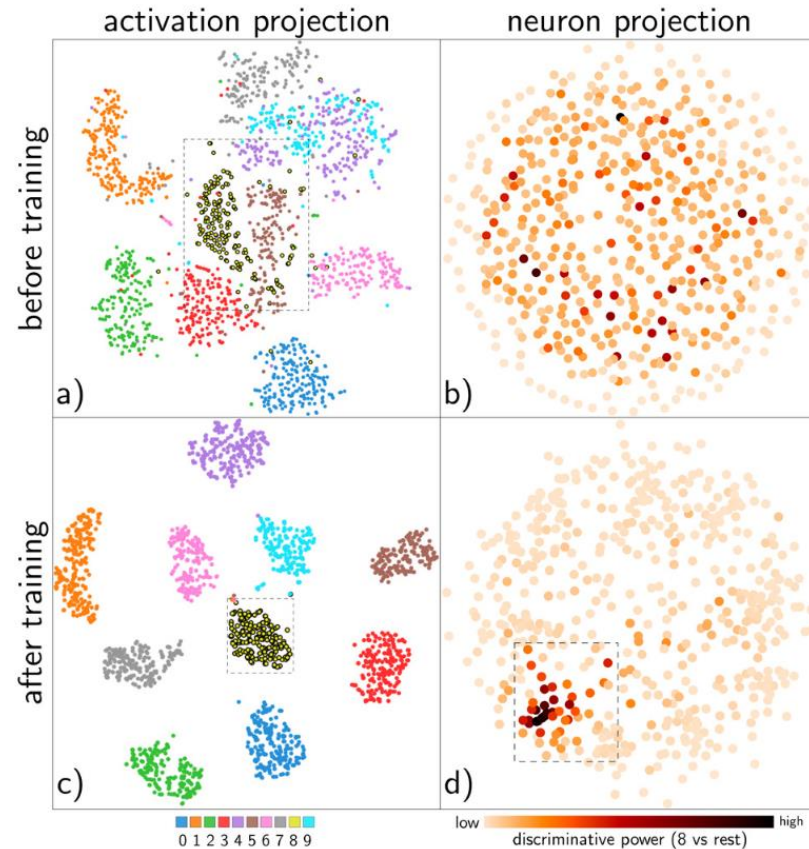


[Maaten2008] Visualization of MNIST

Visualization of hidden layers

- 2D scatter-plot the projections of the hidden neurons' activation coloured according to the class
- Dimensionality reduction and visualization of:
 - observations
 - neuron's relationships
- Clustering of activations in groups to explain predictions

Visualization of hidden layers



[Rauber2016] T-SNE projection of neurons and classes on MNIST test set

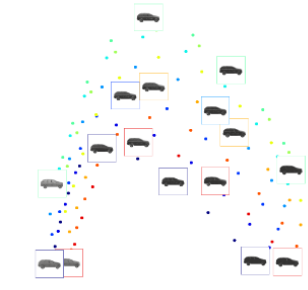
Understanding deep features using PCA



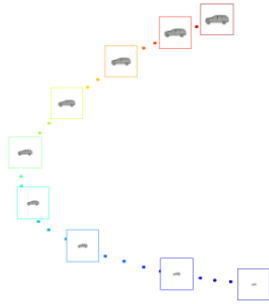
- analyzes CNN feature responses of layers with decomposition into a linear combination of principal components using knowledge from the input scene
- Given
 - an image r_θ from a set Ω with Θ images indexed from $\theta \in [1, \Theta]$
 - features $\hat{\mathbf{F}}^L(r_\theta)$ of the L_{th} layer of CNN
- Calculate centered features $\mathbf{F}^L(r_\theta) = \hat{\mathbf{F}}^L(r_\theta) - \frac{1}{\Theta} \sum_{t=1}^{\Theta} \hat{\mathbf{F}}^L(r_t)$
- r_t : t image from set Ω
- Compute eigenvectors of the covariance matrix:

$$\frac{1}{\Theta} \sum_{\theta=1}^{\Theta} \mathbf{F}^L(r_\theta) \mathbf{F}^L(r_\theta)^T$$

Understanding deep features using PCA



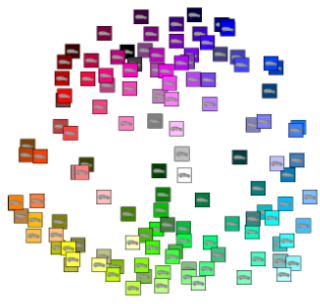
(a) Lighting



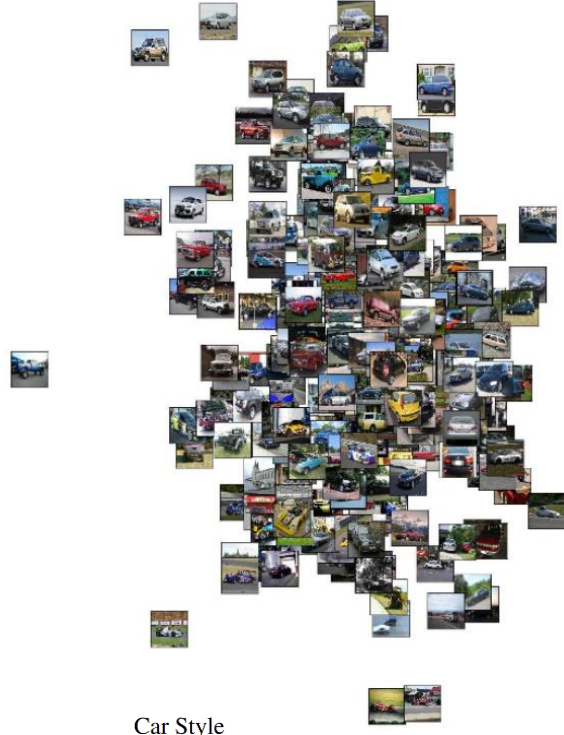
(b) Scale



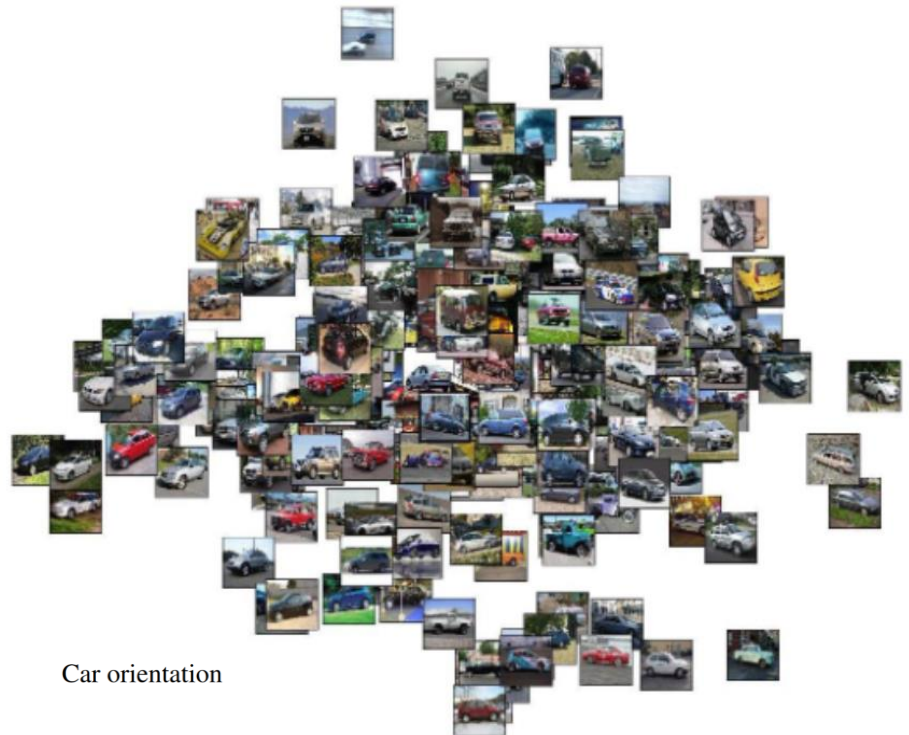
(c) Object color



(d) Background color



Car Style



Car orientation

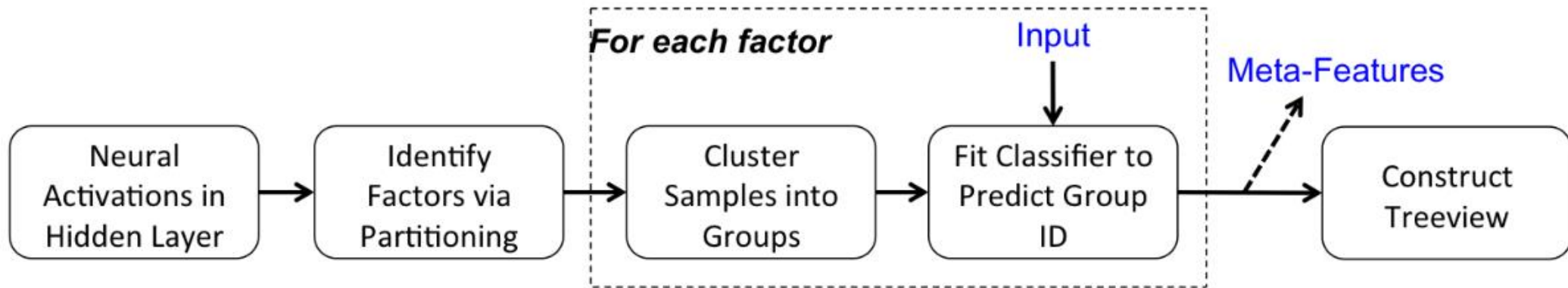
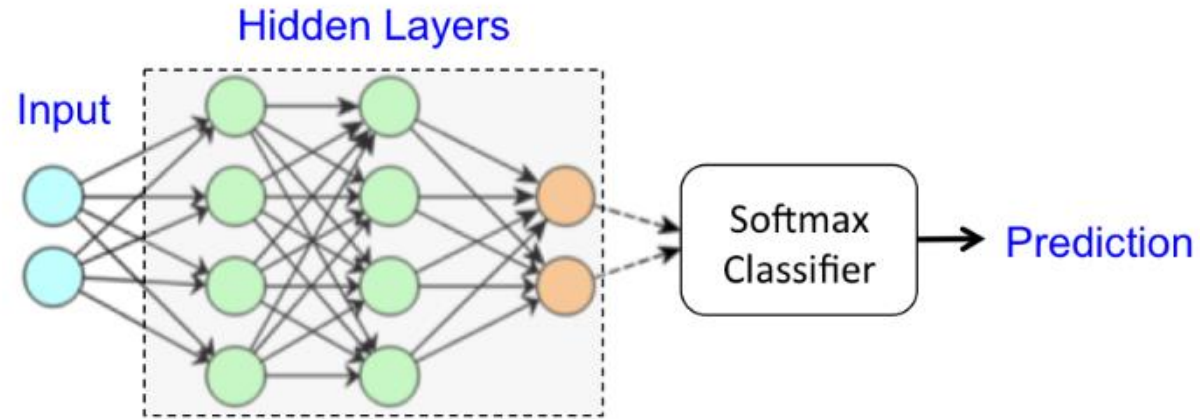
[Aubry2015] PCA embeddings based on different factors

TreeView

- Hierarchical decomposition of the feature subspaces
- Transformation $T_1: \mathbf{X} \rightarrow \mathbf{Y}$ input from space: \mathbf{X} into new space of features: \mathbf{Y} ,
- Transformation $T_2: \mathbf{Y} \rightarrow \mathbf{Z}$ classify features \mathbf{Y} to label space \mathbf{Z}

- Partition space of features: \mathbf{Y} into K subspaces with similar activations of the hidden layers
- Each cluster i describes a specific factor S_i
- A new K -dimensional vector from cluster labels is constructed and used for visualization

TreeView

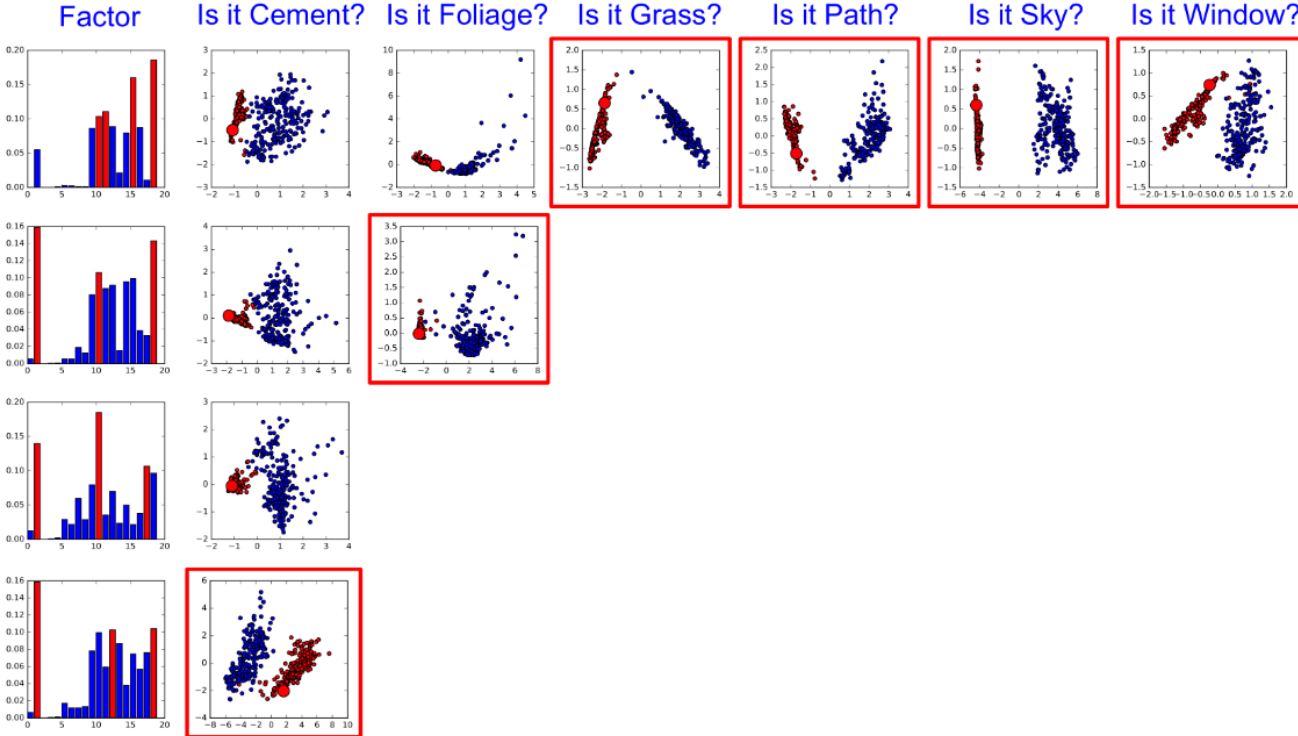


TreeView



Test Sample – True Label: Brickface

••• True Label
••• Hypothesis



[Thiagarajan2016] Visualization of a correctly classified example for each factor

Textual explanations

- Produce natural language-text to explain the decisions of the algorithm
- Find semantic words that provide qualitative explanations
- Methods
 - Cell Activation Value
 - InterpNET
 - Hierarchical Question and Image Co-Attention for Visual Question Answering
 - Visual Dialog
 - Explain Deep Neural Networks with Semantic Information

Cell Activation Value



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.
```

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."
```

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;
}
```

A large portion of cells are not easily interpretable. Here is a typical example:

```
/* Unpack a filter field's string representation from user-space
 * buffer. */
char *audit_unpack_string(void **bufp, size_t *remain, size_t len)
{
    char *str;
    if (!*bufp || (len == 0) || (len > *remain))
        return ERR_PTR(-EINVAL);
    /* Of the currently implemented string fields, PATH_MAX
     * defines the longest valid length.
     */
}
```

Cell that turns on inside comments and quotes:

```
/* Duplicate LSM field information. The lsm_rule is opaque, so
 * re-initialized. */
static inline int audit_dupe_lsm_field(struct audit_field *df,
                                     struct audit_field *sf)
{
    int ret = 0;
    char *lsm_str;
    /* Our own copy of lsm_str */
    lsm_str = kstrdup(sf->lsm_str, GFP_KERNEL);
    if (unlikely(!lsm_str))
        return -ENOMEM;
    df->lsm_str = lsm_str;
    /* Our own (refreshed) copy of lsm_rule */
    ret = security_audit_rule_init(df->type, df->op, df->lsm_str,
                                  (void **)&df->lsm_rule);
    /* Keep currently invalid fields around in case they
     * become valid after a policy reload. */
    if (ret == -EINVAL) {
        pr_warn("audit rule for LSM '%s' is invalid\n",
              df->lsm_str);
        ret = 0;
    }
    return ret;
}
```

Cell that is sensitive to the depth of an expression:

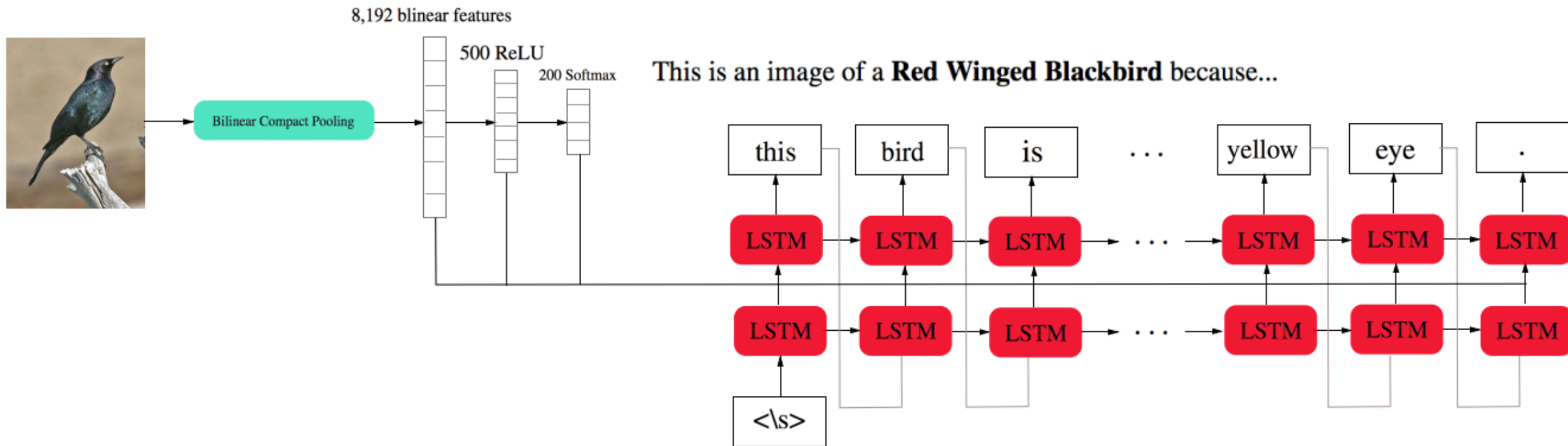
```
#ifdef CONFIG_AUDITSYSCALL
static inline int audit_match_class_bits(int class, u32 *mask)
{
    int i;
    if (classes[class]) {
        for (i = 0; i < AUDIT_BITMASK_SIZE; i++)
            if (mask[i] & classes[class][i])
                return 0;
    }
    return 1;
}
```

Cell that might be helpful in predicting a new line. Note that it only turns on for some "":

```
char *audit_unpack_string(void **bufp, size_t *remain, size_t len)
{
    char *str;
    if (!*bufp || (len == 0) || (len > *remain))
        return ERR_PTR(-EINVAL);
    /* Of the currently implemented string fields, PATH_MAX
     * defines the longest valid length.
     */
    if (len > PATH_MAX)
        return ERR_PTR(-ENAMETOOLONG);
    str = kmalloc(len + 1, GFP_KERNEL);
    if (unlikely(!str))
        return ERR_PTR(-ENOMEM);
    memcpy(str, *bufp, len);
    str[len] = 0;
    *bufp += len;
    *remain -= len;
    return str;
}
```

[Karpathy2015] activations of LSTM. Text color corresponds to activation value

InterpNET



[Barratt2017] Overview of Interpnet



This is an image of a **Tree Sparrow** because...

InterpNET(0) this bird has a **brown crown, brown primaries, and a brown belly.**
InterpNET(1) this bird has **wings that are brown** and has a **white belly.**
InterpNET(2) this bird has **wings that are brown** and has a **white belly.**
InterpNET(3) this bird has a **brown crown, brown primaries, and a brown belly.**
Captioning this bird has a **brown crown, brown primaries, and a brown belly.**



This is an image of a **Philadelphia Vireo** because...

InterpNET(0) this bird has **wings that are grey** and has a **yellow belly.**
InterpNET(1) this bird has **wings that are grey** and has a **yellow belly.**
InterpNET(2) this bird has **wings that are brown** and has a **yellow belly.**
InterpNET(3) this bird has a **yellow belly** and breast with a **short pointy bill.**
Captioning this bird has a **yellow belly** and breast with a **gray crown** and **white wingbars.**



This is an image of an **Ivory Gull** because...

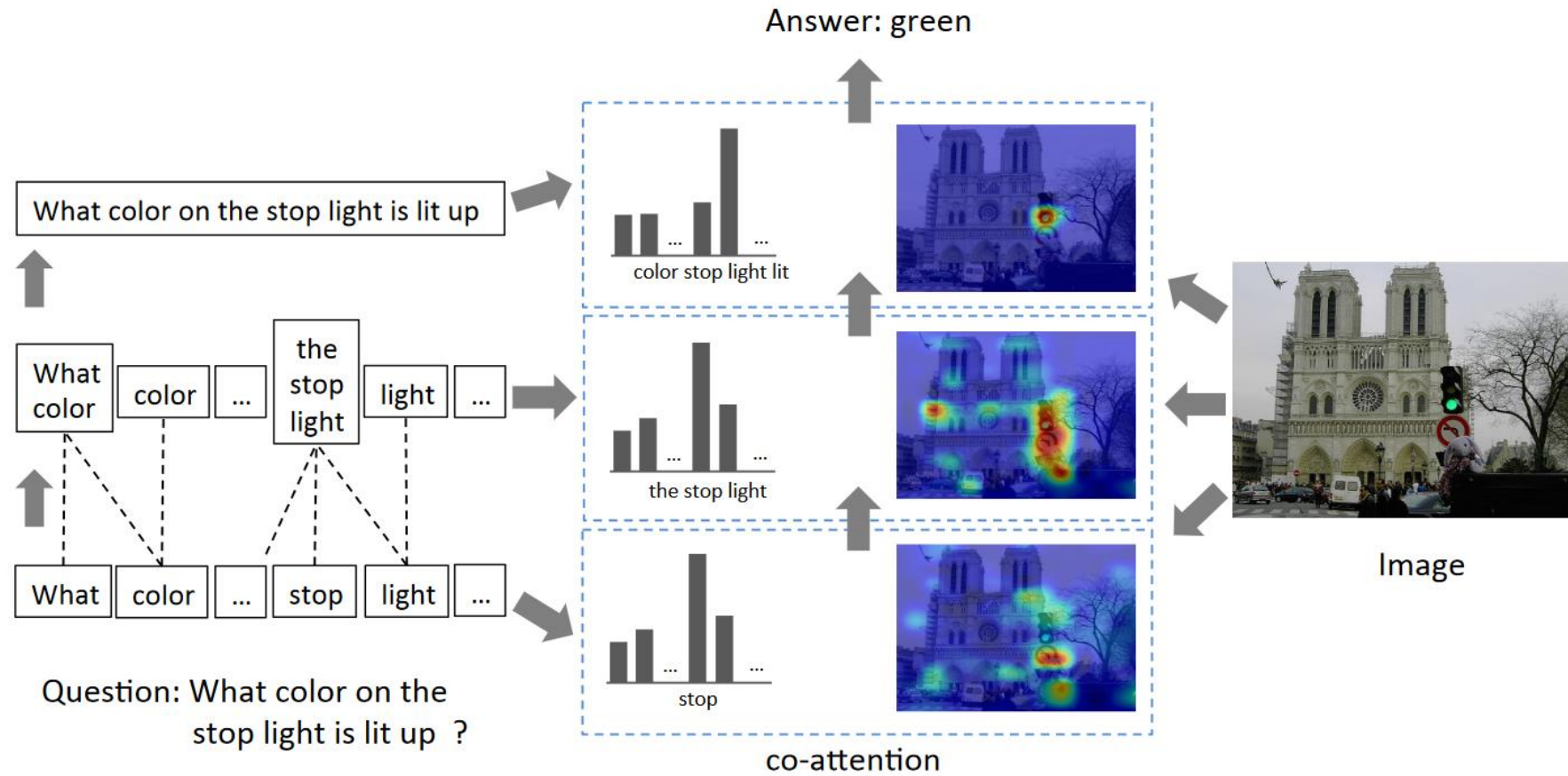
InterpNET(0) this bird has **wings that are white** and has a **yellow bill.**
InterpNET(1) this bird has **wings that are black** and has a **white belly.**
InterpNET(2) this bird has **wings that are grey** and has a **white belly.**
InterpNET(3) this bird has **wings that are grey** and has a **white belly.**
Captioning this bird has a **white belly and breast** with a **black wing** and **long hooked bill.**



This is an image of a **Scott Oriole** because...

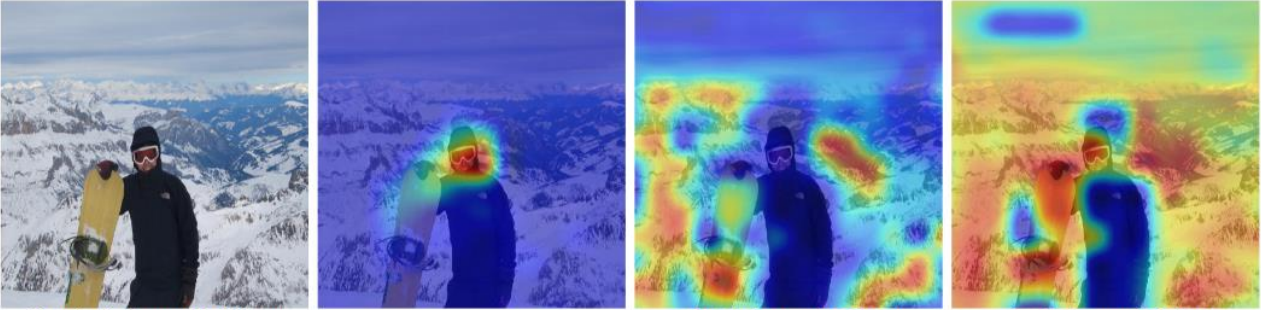
InterpNET(0) this bird has a **yellow belly and breast** with a **black superciliary** and **white wingbars.**
InterpNET(1) this bird has a **black crown, black primaries, and a yellow belly.**
InterpNET(2) this bird has **wings that are black** and has a **yellow belly.**
InterpNET(3) this bird has a **black crown, a black bill, and a black breast.**
Captioning this bird has a **black crown, black primaries, and a white belly.**

Hierarchical Question-Image Attention for Visual Question Answering

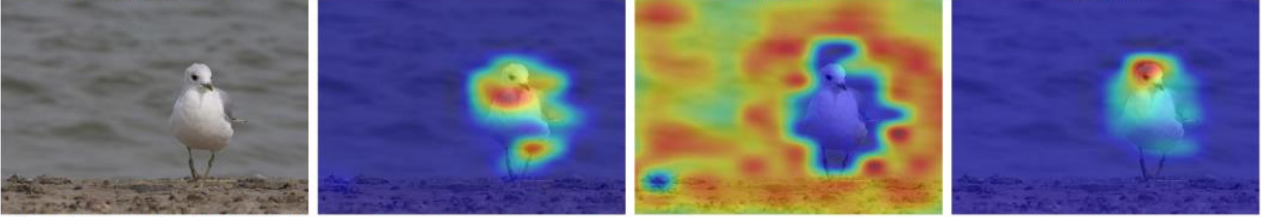


[Lu2016] Types of attention

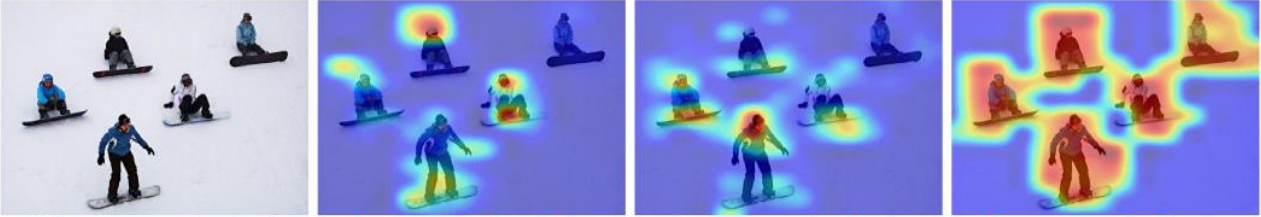
Hierarchical Question-Image Attention for Visual Question Answering



Q: what is the man holding a snowboard on top of a snow covered? A: mountain
what is the man holding a snowboard on top of a snow covered
what is the man holding a snowboard on top of a snow covered ?
what is the man holding a snowboard on top of a snow covered ?



Q: what is the color of the bird? A: white
what is the color of the bird ?
what is the color of the bird ?
what is the color of the bird ?



Q: how many snowboarders in formation in the snow, four is sitting? A: 5
how many snowboarders in formation in the snow , four is sitting ?
how many snowboarders in formation in the snow , four is sitting ?
how many snowboarders in formation in the snow , four is sitting ?

[Lu2016] Example of attention maps

Visual dialog



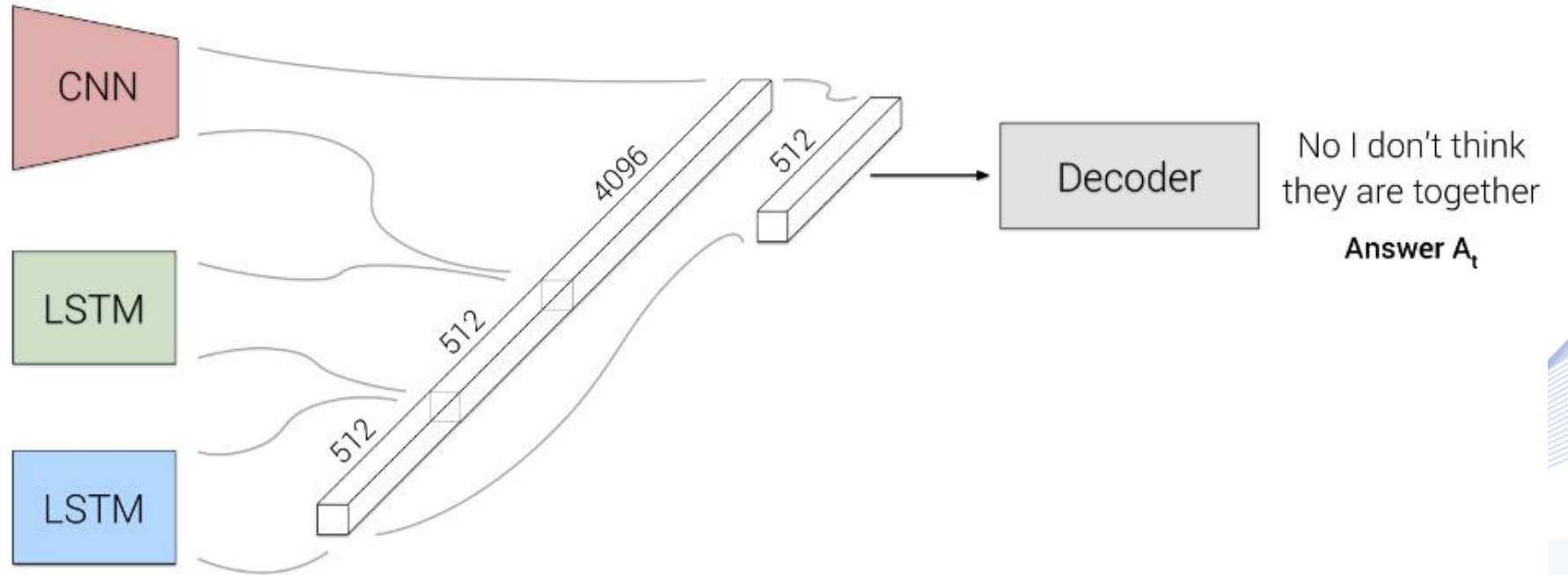
Image I

Do you think the woman is with him?

Question Q_t

The man is riding his bicycle on the sidewalk. Is the man wearing a helmet? No he does not have a helmet on. ... Are there any people nearby? Yes there's a woman walking behind him.


t rounds of history
(concatenated)



[Das2018] Visual dialog system to predict the answer

Visual dialog

Visual Dialog



A cat drinking water out of a coffee mug.

What color is the mug?

White and red

Are there any pictures on it?

No, something is there can't tell what it is

Is the mug and cat on a table?

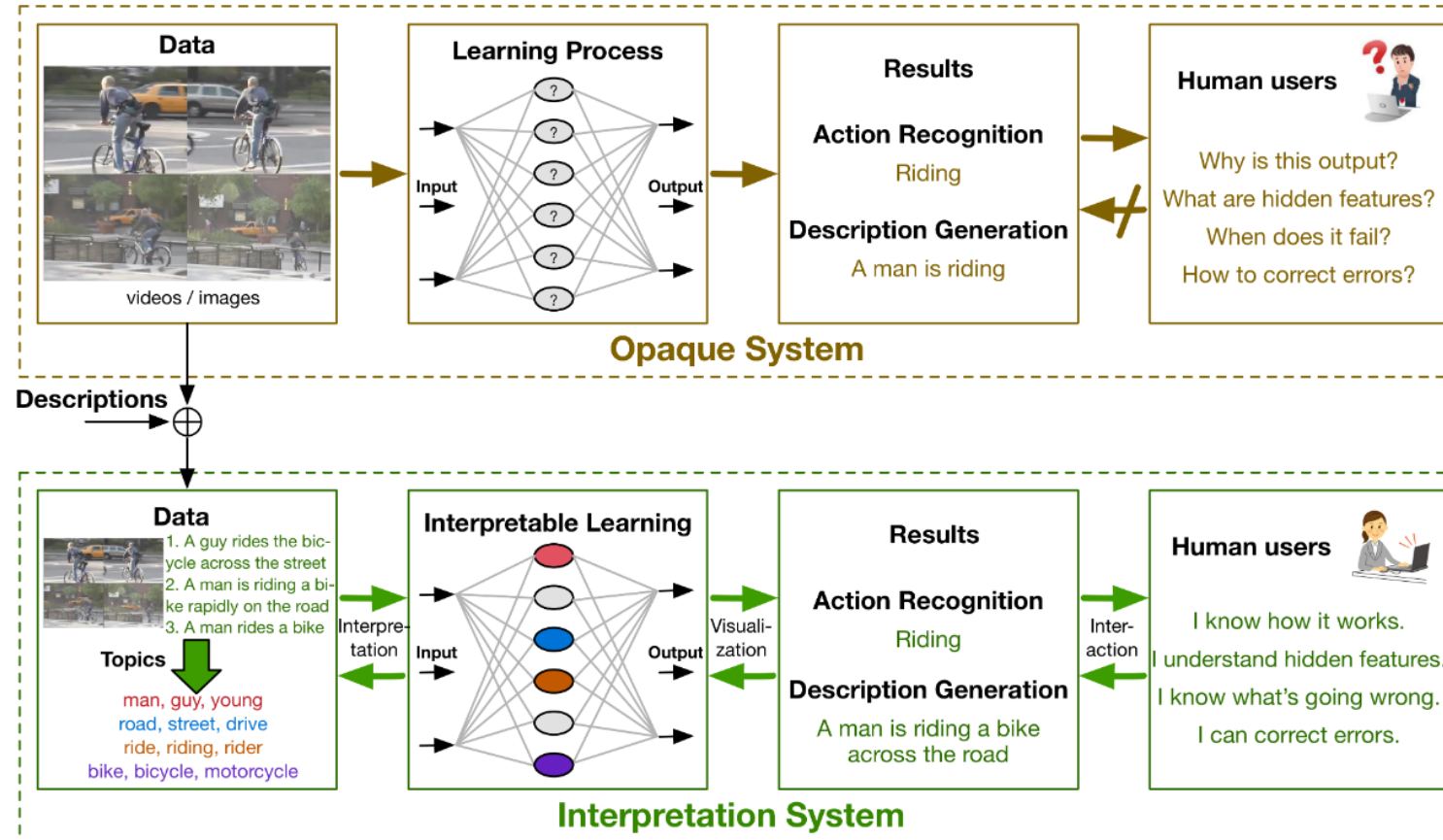
Yes, they are

Are there other items on the table?

Yes, magazines, books, toaster and basket, and a plate

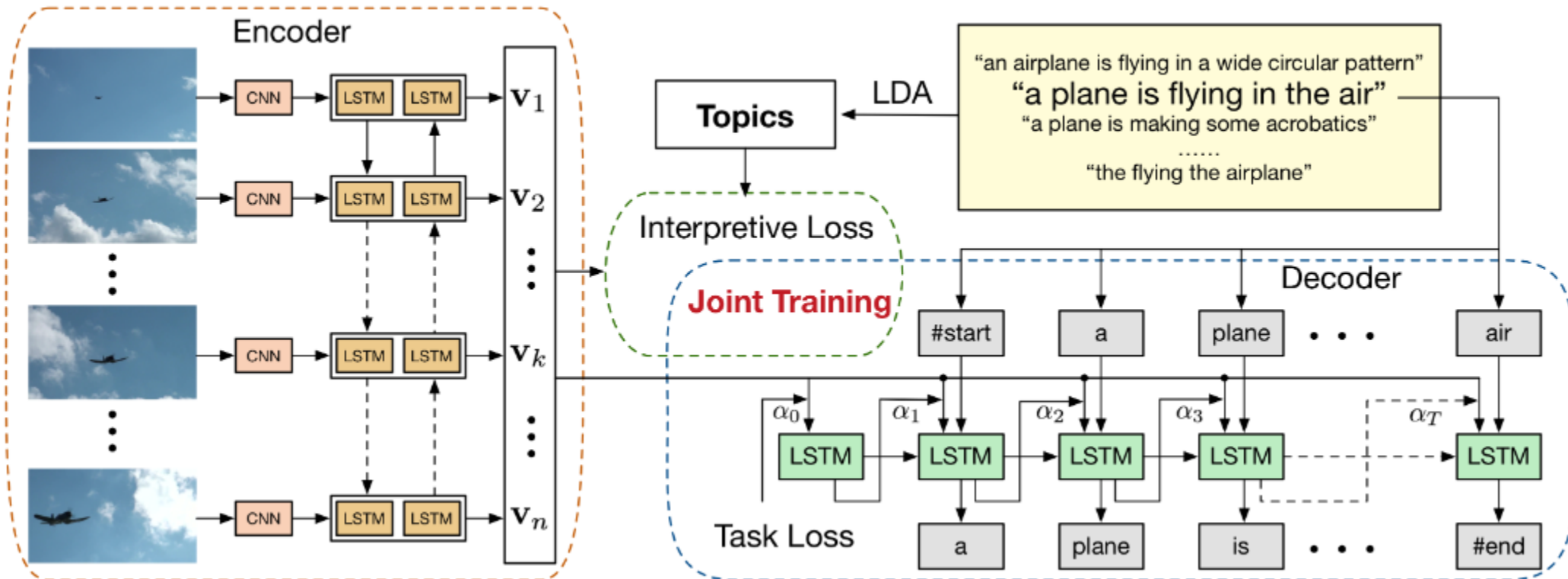
Start typing question here ...

Semantic explanation of Deep Neural Networks



[Dong2017] Overview of the method

Semantic explanation of Deep Neural Networks

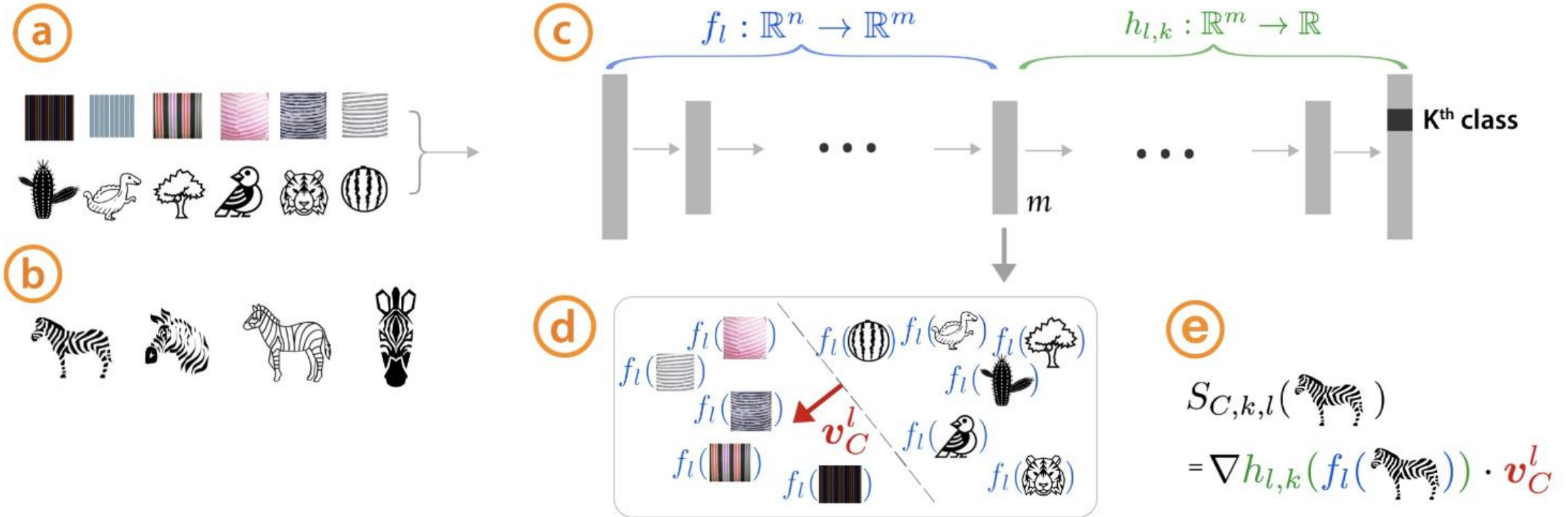


[Dong2017] Video captioning encoder's training process along with human descriptions decoder

Numerical Explanations

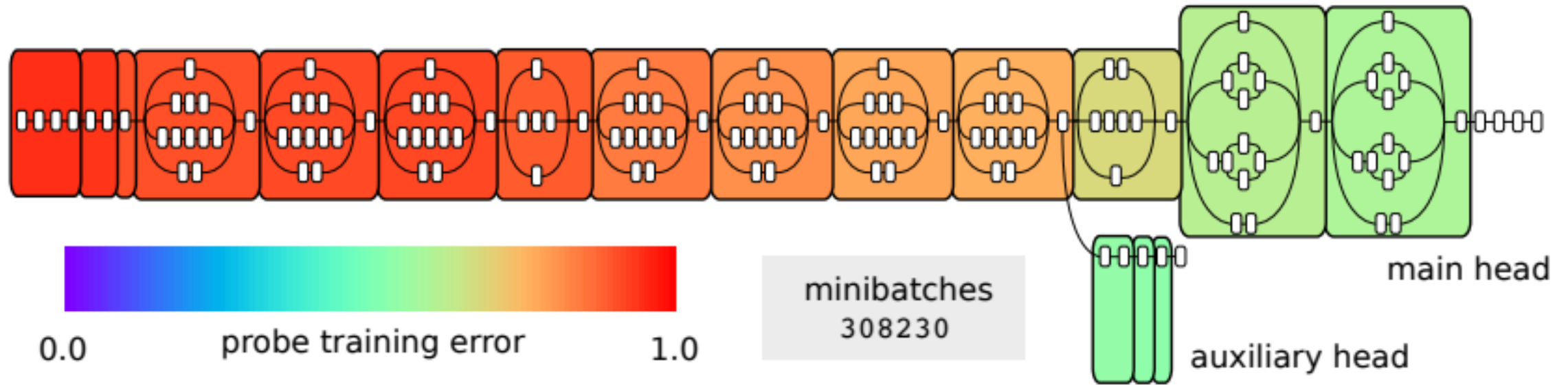
- Provide numerical outputs to interpret models
- Train classifiers that explain the model
 - Concept Activation Vectors
 - Linear classifier probes
 - LIME

Concept Activation Vectors (CAV)



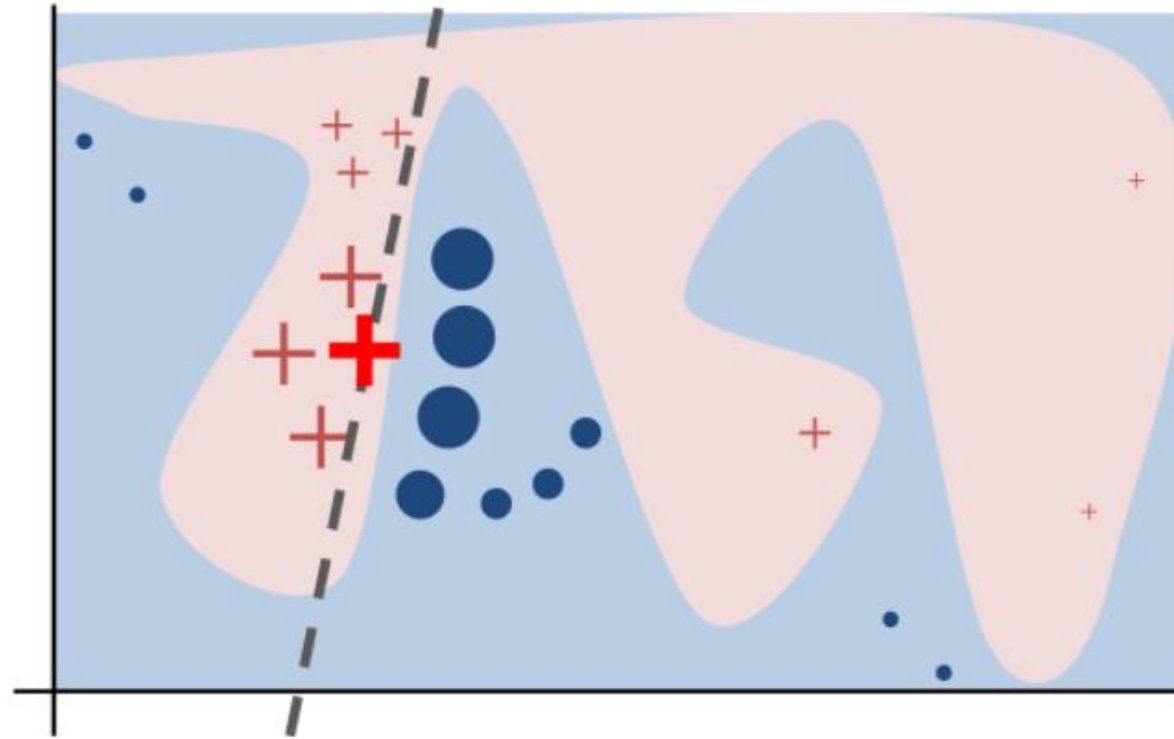
[Kim2018] Overview of CAVs learning process

Linear classifier probes



[Alain2016] Prediction error of each layer using probe

LIME (Local Interpretable Model-agnostic Explanations)



[Ribeiro2016] The data represented with the red cross is explained locally using the dashed line.

Applications



of explainable methods in important tasks

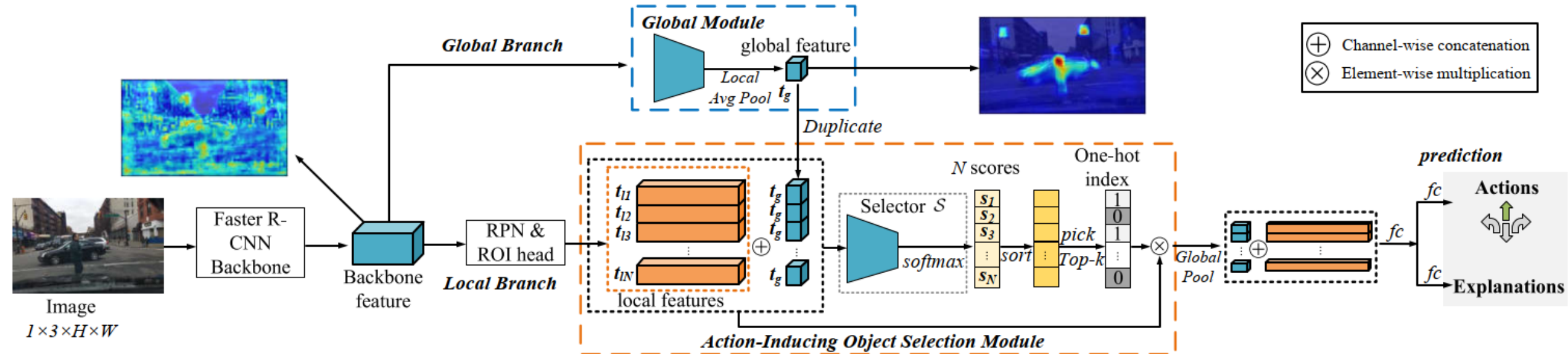
- Autonomous Driving
 - Advisable Learning for Self-driving Vehicles by Internalizing Observation-to-Action Rules
 - Explaining Autonomous Driving by Learning End-to-End Visual Attention
- Medical Applications
 - COVID detection

Explainable Object-induced Action Decision for Autonomous Vehicles



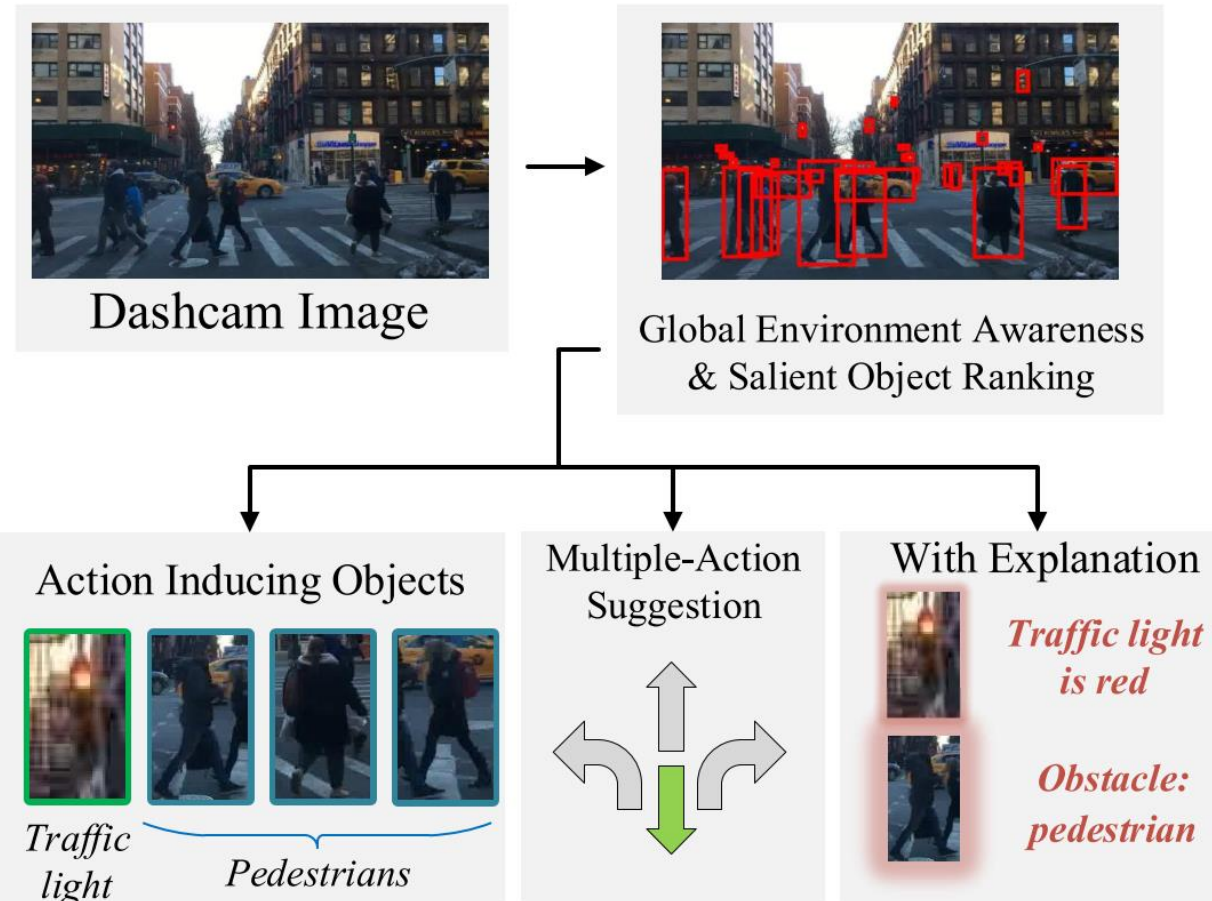
- Explainable autonomous driving architecture
- Global module to generate scene context
- Local module to predict actions and explanations
- Concatenation of the two modules to improve decision accuracy
- Predict next action such as (stop, turn left, go forward) and the explanations e.g. Stop, the traffic light is red , an obstacle is in front stop the car)

Explainable Object-induced Action Decision for Autonomous Vehicles



[Yu2020] Overview of the proposed explainable self-driving method

Explainable Object-induced Action Decision for Autonomous Vehicles



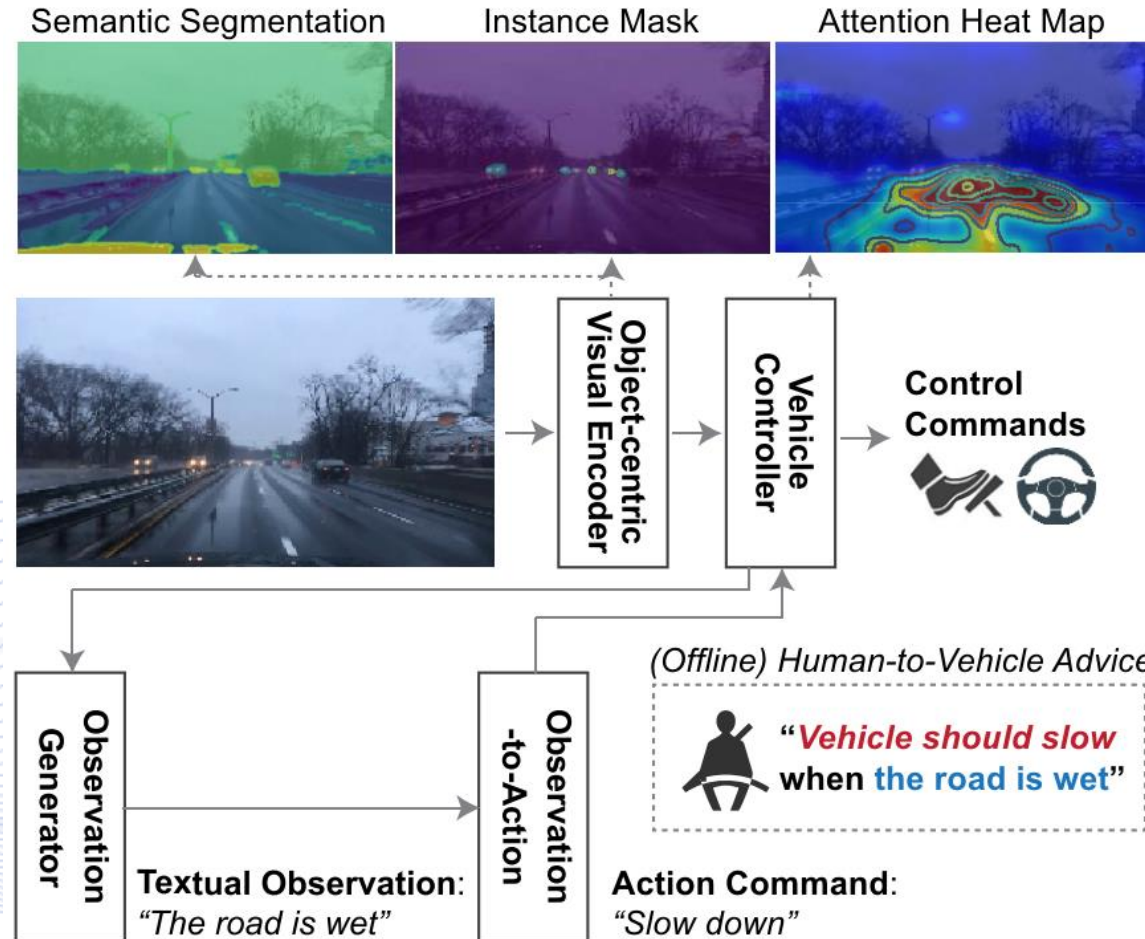
[Yu2020] Examples of action and explanations of the system

Advisable Learning for Self-driving Vehicles by Internalizing Observation-to-Action Rules



- Explainable self-driving system
- Human advice integrated on the architecture
- Visual (attention maps) and textual explanations (sentences)
- Main modules
 1. Object segmentation encoder
 2. Vehicle controller
 3. Observation generator
 4. Action generator

Advisable Learning for Self-driving Vehicles by Internalizing Observation-to-Action Rules

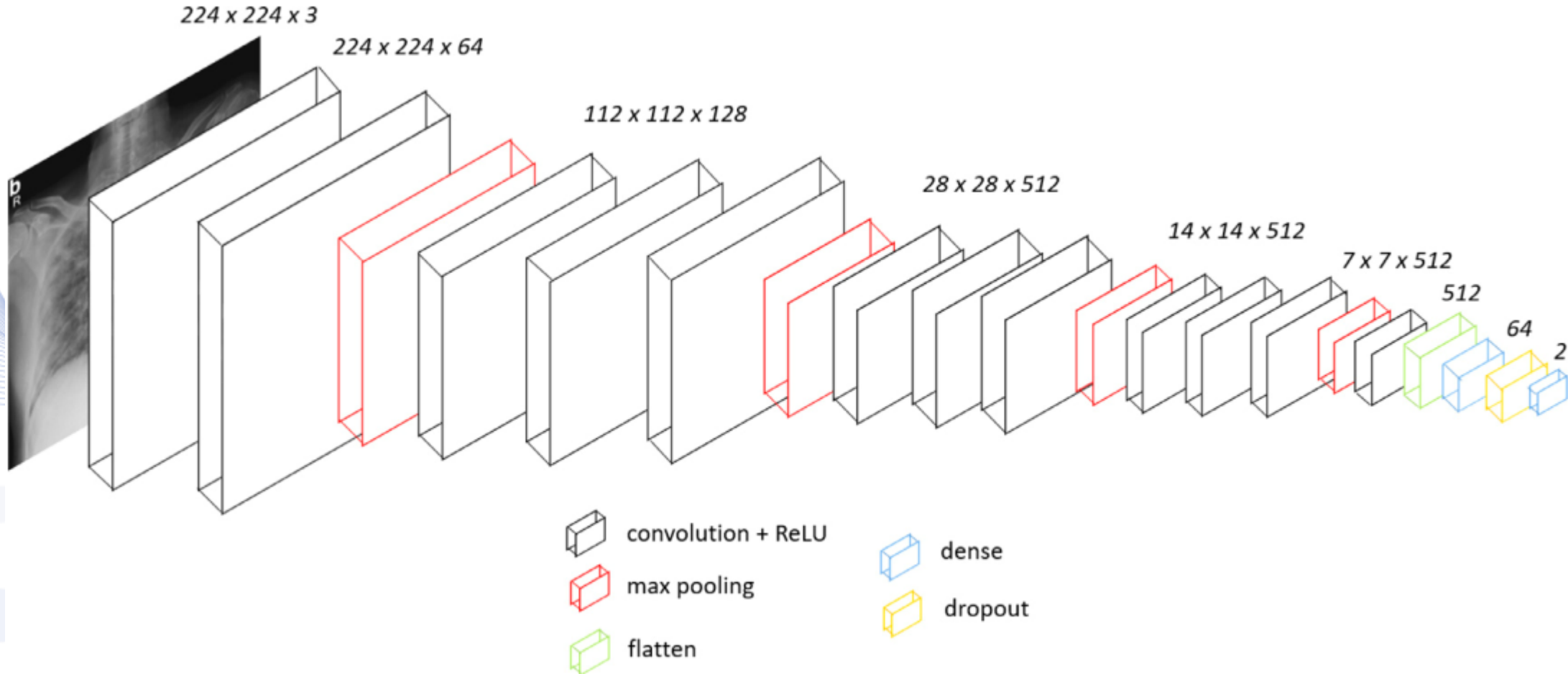


Explainable Pulmonary Disease and COVID-19 Detection from X-rays



- Deep convolutional network VGG-16 is used to distinguish between healthy lungs, pneumonia or covid-19
- Grad-CAM provides interprets decisions by visualizing feature maps
- Localize areas responsible for the detection of pneumonia or covid-19

Explainable Pulmonary Disease and COVID-19 Detection from X-rays



Explainable Pulmonary Disease and COVID-19 Detection from X-rays



[Brunese2020] Response maps of input chest x-ray

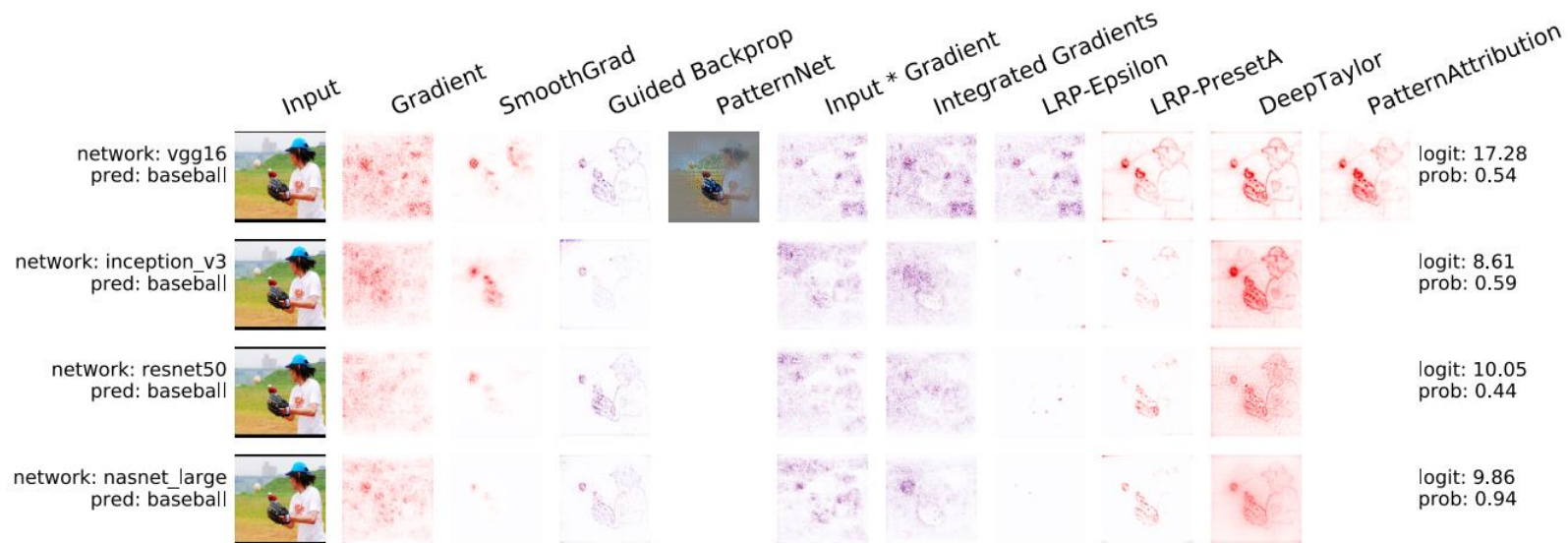
Explainable AI frameworks

with implemented explainable methods can be used for interpretation

- iNNvestigate Neural Networks
- Explainer
- InterpretML

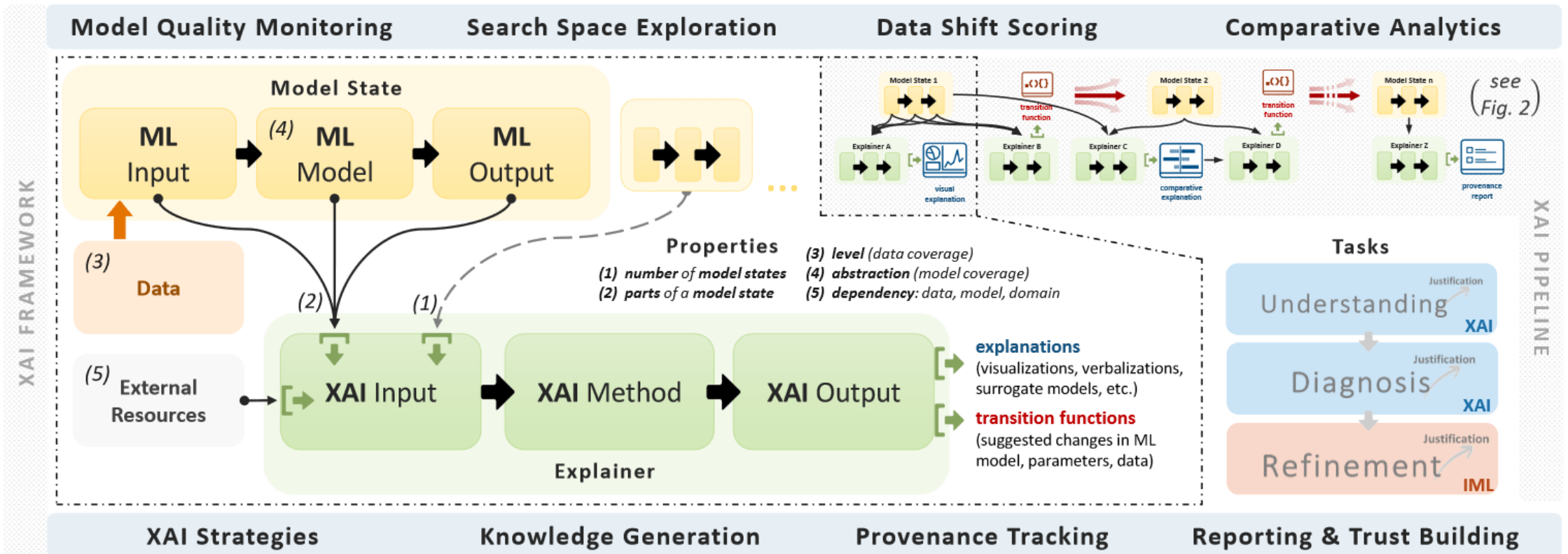
iNNvestigate Neural Networks

```
import innvestigate
model = create_a_keras_model()
analyzer = innvestigate.create_analyzer('analyzer_name', model)
analyzer.fit(X_train) # if needed
analysis = analyzer.analyze(X_test)
```



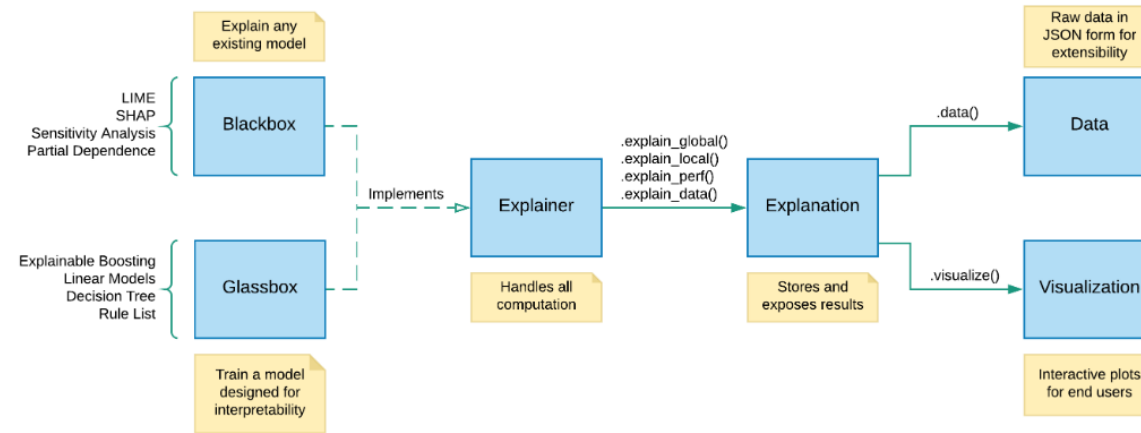
[Alber2019] Usage Example of iNNvestigate Neural Networks

ExplAIner



[Spinner2019] General approach of explAIner

InterpretML



Glassbox

Blackbox

```
1 from interpret import show
2 from interpret.glassbox import LogisticRegression
3
4
5 clf = LogisticRegression()
6 clf.fit(X, y)
7
8 global_exp = clf.explain_global()
9 local_exp = clf.explain_local(X, y)
10
11 show([global_exp, local_exp])
```

```
1 from interpret import show
2 from interpret.blackbox import PartialDependence
3 from sklearn.neural_network import MLPClassifier
4
5 blackbox = MLPClassifier()
6 blackbox.fit(X, y)
7
8 pdp = PartialDependence(blackbox.predict_proba, X)
9 global_exp = pdp.explain_global()
10
11 show(global_exp)
```

[Nori2019] Example of InterpretML usage

References



- [Vilone2020]** Vilone, Giulia, and Luca Longo. "Explainable Artificial Intelligence: a Systematic Review." *arXiv preprint arXiv:2006.00093* (2020).
- [Tjoa2020]** Tjoa, Erico, and Cuntai Guan. "A Survey on Explainable Artificial Intelligence (XAI): Toward Medical XAI." *IEEE Transactions on Neural Networks and Learning Systems* (2020).
- [Adadi2018]** Adadi, Amina, and Mohammed Berrada. "Peeking inside the black-box: A survey on Explainable Artificial Intelligence (XAI)." *IEEE Access* 6 (2018): 52138-52160.
- [Zhou2016]** Zhou, Bolei, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. "Learning deep features for discriminative localization." In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2921-2929. 2016.
- [Selvaraju2017]** Selvaraju, Ramprasaath R., Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. "Grad-cam: Visual explanations from deep networks via gradient-based localization." In *Proceedings of the IEEE international conference on computer vision*, pp. 618-626. 2017.
- [Samek2016]** Samek, Wojciech, Grégoire Montavon, Alexander Binder, Sebastian Lapuschkin, and Klaus-Robert Müller. "Interpreting the predictions of complex ml models by layer-wise relevance propagation." *arXiv preprint arXiv:1611.08191* (2016).
- [Ribeiro2016]** Ribeiro, Marco Tulio, Sameer Singh, and Carlos Guestrin. "" Why should I trust you?" Explaining the predictions of any classifier." In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pp. 1135-1144. 2016.

References



- [Zhou2016]** Zhou, Yanzhao, Yi Zhu, Qixiang Ye, Qiang Qiu, and Jianbin Jiao. "Weakly supervised instance segmentation using class peak response." In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3791-3800. 2018.
- [Zeiler2014]** Zeiler, Matthew D., and Rob Fergus. "Visualizing and understanding convolutional networks." In *European conference on computer vision*, pp. 818-833. Springer, Cham, 2014.
- [Dosovitskiy2015]** Dosovitskiy, Alexey, and Thomas Brox. "Inverting convolutional networks with convolutional networks." *arXiv preprint arXiv:1506.02753* 4 (2015).
- [Kumar2017]** Kumar, Devinder, Alexander Wong, and Graham W. Taylor. "Explaining the unexplained: A class-enhanced attentive response (clear) approach to understanding deep neural networks." In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 36-44. 2017.
- [Liu2019]** Liu G, Zeng H, Gifford DK. Visualizing complex feature interactions and feature sharing in genomic deep neural networks. *BMC bioinformatics*. 2019 Dec;20(1):1-4.
- [Aubry2015]** Aubry, Mathieu, and Bryan C. Russell. "Understanding deep features with computer-generated imagery." In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 2875-2883. 2015.
- [Rauber2016]** Rauber, Paulo E., Samuel G. Fadel, Alexandre X. Falcao, and Alexandru C. Telea. "Visualizing the hidden activity of artificial neural networks." *IEEE transactions on visualization and computer graphics* 23, no. 1 (2016): 101-110.

References



- [Thiagarajan2016]** Thiagarajan, Jayaraman J., Bhavya Kailkhura, Prasanna Sattigeri, and Karthikeyan Natesan Ramamurthy. "Treeview: Peeking into deep neural networks via feature-space partitioning." *arXiv preprint arXiv:1611.07429* (2016).
- [Maaten2008]** Maaten, Laurens van der, and Geoffrey Hinton. "Visualizing data using t-SNE." *Journal of machine learning research* 9, no. Nov (2008): 2579-2605.
- [Karpathy2015]** Karpathy, Andrej, Justin Johnson, and Li Fei-Fei. "Visualizing and understanding recurrent networks." *arXiv preprint arXiv:1506.02078* (2015).
- [Barratt2017]** Barratt, Shane. "Interpnet: Neural introspection for interpretable deep learning." *arXiv preprint arXiv:1710.09511* (2017).
- [Lu2016]** Lu, Jiasen, Jianwei Yang, Dhruv Batra, and Devi Parikh. "Hierarchical question-image co-attention for visual question answering." In *Advances in neural information processing systems*, pp. 289-297. 2016.
- [Dong2017]** Dong, Yinpeng, Hang Su, Jun Zhu, and Bo Zhang. "Improving interpretability of deep neural networks with semantic information." In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4306-4314. 2017
- [Das2018]** Das, A., S. Kottur, K. Gupta, A. Singh, D. Yadav, S. Lee, J. Moura, D. Parikh, and D. Batra. "Visual Dialog." *IEEE transactions on pattern analysis and machine intelligence* (2018).

References



- [Kim2018]** Kim, Been, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, and Fernanda Viegas. "Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav)." In *International conference on machine learning*, pp. 2668-2677. PMLR, 2018.
- [Alain2016]** Alain, Guillaume, and Yoshua Bengio. "Understanding intermediate layers using linear classifier probes." *arXiv preprint arXiv:1610.01644* (2016).
- [Li2020]** Li, Liangzhi, Bowen Wang, Manisha Verma, Yuta Nakashima, Ryo Kawasaki, and Hajime Nagahara. "SCOUTER: Slot attention-based classifier for explainable image recognition." *arXiv preprint arXiv:2009.06138* (2020).
- [Oramas2019]** Oramas Mogrovejo, J. A., Wang, K., & Tuytelaars, T. (2019). Visual explanation by interpretation: Improving visual feedback capabilities of deep neural networks. In <https://iclr.cc/Conferences/2019/AcceptedPapersInitial>. openReview.
- [Kim2020]** Kim, Jinkyu, Suhong Moon, Anna Rohrbach, Trevor Darrell, and John Canny. "Advisable Learning for Self-Driving Vehicles by Internalizing Observation-to-Action Rules." In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9661-9670. 2020.
- [Yu2020]** Xu, Yiran, Xiaoyin Yang, Lihang Gong, Hsuan-Chu Lin, Tz-Ying Wu, Yunsheng Li, and Nuno Vasconcelos. "Explainable Object-induced Action Decision for Autonomous Vehicles." In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9523-9532. 2020.

References



[Brunese2020] Brunese, Luca, Francesco Mercaldo, Alfonso Reginelli, and Antonella Santone. "Explainable deep learning for pulmonary disease and coronavirus COVID-19 detection from X-rays." *Computer Methods and Programs in Biomedicine* 196 (2020): 105608.

[Alber2019] Alber, Maximilian, Sebastian Lapuschkin, Philipp Seegerer, Miriam Hägele, Kristof T. Schütt, Grégoire Montavon, Wojciech Samek, Klaus-Robert Müller, Sven Dähne, and Pieter-Jan Kindermans. "iNNvestigate neural networks!." *J. Mach. Learn. Res.* 20, no. 93 (2019): 1-8.

[Spinner2019] Spinner, Thilo, Udo Schlegel, Hanna Schäfer, and Mennatallah El-Assady. "explAIner: A visual analytics framework for interactive and explainable machine learning." *IEEE transactions on visualization and computer graphics* 26, no. 1 (2019): 1064-1074.

[Nori2019] Nori, Harsha, Samuel Jenkins, Paul Koch, and Rich Caruana. "InterpretML: A Unified Framework for Machine Learning Interpretability." *arXiv preprint arXiv:1909.09223* (2019).

Bibliography



- [Vilone2020]** Vilone, Giulia, and Luca Longo. "Explainable Artificial Intelligence: a Systematic Review." *arXiv preprint arXiv:2006.00093* (2020).
- [Tjoa2020]** Tjoa, Erico, and Cuntai Guan. "A Survey on Explainable Artificial Intelligence (XAI): Toward Medical XAI." *IEEE Transactions on Neural Networks and Learning Systems* (2020).
- [Adadi2018]** Adadi, Amina, and Mohammed Berrada. "Peeking inside the black-box: A survey on Explainable Artificial Intelligence (XAI)." *IEEE Access* 6 (2018): 52138-52160.
- [Zhou2016]** Zhou, Bolei, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. "Learning deep features for discriminative localization." In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2921-2929. 2016.
- [Selvaraju2017]** Selvaraju, Ramprasaath R., Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. "Grad-cam: Visual explanations from deep networks via gradient-based localization." In *Proceedings of the IEEE international conference on computer vision*, pp. 618-626. 2017.
- [Samek2016]** Samek, Wojciech, Grégoire Montavon, Alexander Binder, Sebastian Lapuschkin, and Klaus-Robert Müller. "Interpreting the predictions of complex ml models by layer-wise relevance propagation." *arXiv preprint arXiv:1611.08191* (2016).
- [Ribeiro2016]** Ribeiro, Marco Tulio, Sameer Singh, and Carlos Guestrin. "" Why should I trust you?" Explaining the predictions of any classifier." In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pp. 1135-1144. 2016.

Bibliography



- [Zhou2016]** Zhou, Yanzhao, Yi Zhu, Qixiang Ye, Qiang Qiu, and Jianbin Jiao. "Weakly supervised instance segmentation using class peak response." In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3791-3800. 2018.
- [Zeiler2014]** Zeiler, Matthew D., and Rob Fergus. "Visualizing and understanding convolutional networks." In *European conference on computer vision*, pp. 818-833. Springer, Cham, 2014.
- [Dosovitskiy2015]** Dosovitskiy, Alexey, and Thomas Brox. "Inverting convolutional networks with convolutional networks." *arXiv preprint arXiv:1506.02753* 4 (2015).
- [Kumar2017]** Kumar, Devinder, Alexander Wong, and Graham W. Taylor. "Explaining the unexplained: A class-enhanced attentive response (clear) approach to understanding deep neural networks." In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 36-44. 2017.
- [Liu2019]** Liu G, Zeng H, Gifford DK. Visualizing complex feature interactions and feature sharing in genomic deep neural networks. *BMC bioinformatics*. 2019 Dec;20(1):1-4.
- [Aubry2015]** Aubry, Mathieu, and Bryan C. Russell. "Understanding deep features with computer-generated imagery." In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 2875-2883. 2015.
- [Rauber2016]** Rauber, Paulo E., Samuel G. Fadel, Alexandre X. Falcao, and Alexandru C. Telea. "Visualizing the hidden activity of artificial neural networks." *IEEE transactions on visualization and computer graphics* 23, no. 1 (2016): 101-110.

Bibliography



- [Thiagarajan2016]** Thiagarajan, Jayaraman J., Bhavya Kailkhura, Prasanna Sattigeri, and Karthikeyan Natesan Ramamurthy. "Treeview: Peeking into deep neural networks via feature-space partitioning." *arXiv preprint arXiv:1611.07429* (2016).
- [Maaten2008]** Maaten, Laurens van der, and Geoffrey Hinton. "Visualizing data using t-SNE." *Journal of machine learning research* 9, no. Nov (2008): 2579-2605.
- [Karpathy2015]** Karpathy, Andrej, Justin Johnson, and Li Fei-Fei. "Visualizing and understanding recurrent networks." *arXiv preprint arXiv:1506.02078* (2015).
- [Barratt2017]** Barratt, Shane. "Interpnet: Neural introspection for interpretable deep learning." *arXiv preprint arXiv:1710.09511* (2017).
- [Lu2016]** Lu, Jiasen, Jianwei Yang, Dhruv Batra, and Devi Parikh. "Hierarchical question-image co-attention for visual question answering." In *Advances in neural information processing systems*, pp. 289-297. 2016.
- [Dong2017]** Dong, Yinpeng, Hang Su, Jun Zhu, and Bo Zhang. "Improving interpretability of deep neural networks with semantic information." In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4306-4314. 2017
- [Das2018]** Das, A., S. Kottur, K. Gupta, A. Singh, D. Yadav, S. Lee, J. Moura, D. Parikh, and D. Batra. "Visual Dialog." *IEEE transactions on pattern analysis and machine intelligence* (2018).

Bibliography



[Kim2018] Kim, Been, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, and Fernanda Viegas. "Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav)." In *International conference on machine learning*, pp. 2668-2677. PMLR, 2018.

[Alain2016] Alain, Guillaume, and Yoshua Bengio. "Understanding intermediate layers using linear classifier probes." *arXiv preprint arXiv:1610.01644* (2016).

[Li2020] Li, Liangzhi, Bowen Wang, Manisha Verma, Yuta Nakashima, Ryo Kawasaki, and Hajime Nagahara. "SCOUTER: Slot attention-based classifier for explainable image recognition." *arXiv preprint arXiv:2009.06138* (2020).

[Oramas2019] Oramas Mogrovejo, J. A., Wang, K., & Tuytelaars, T. (2019). Visual explanation by interpretation: Improving visual feedback capabilities of deep neural networks. In <https://iclr.cc/Conferences/2019/AcceptedPapersInitial>. openReview.

[Kim2020] Kim, Jinkyu, Suhong Moon, Anna Rohrbach, Trevor Darrell, and John Canny. "Advisable Learning for Self-Driving Vehicles by Internalizing Observation-to-Action Rules." In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9661-9670. 2020.

[Yu2020] Xu, Yiran, Xiaoyin Yang, Lihang Gong, Hsuan-Chu Lin, Tz-Ying Wu, Yunsheng Li, and Nuno Vasconcelos. "Explainable Object-induced Action Decision for Autonomous Vehicles." In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9523-9532. 2020.

Bibliography



[Brunese2020] Brunese, Luca, Francesco Mercaldo, Alfonso Reginelli, and Antonella Santone. "Explainable deep learning for pulmonary disease and coronavirus COVID-19 detection from X-rays." *Computer Methods and Programs in Biomedicine* 196 (2020): 105608.

[Alber2019] Alber, Maximilian, Sebastian Lapuschkin, Philipp Seegerer, Miriam Hägele, Kristof T. Schütt, Grégoire Montavon, Wojciech Samek, Klaus-Robert Müller, Sven Dähne, and Pieter-Jan Kindermans. "iNNvestigate neural networks!." *J. Mach. Learn. Res.* 20, no. 93 (2019): 1-8.

[Spinner2019] Spinner, Thilo, Udo Schlegel, Hanna Schäfer, and Mennatallah El-Assady. "explAIner: A visual analytics framework for interactive and explainable machine learning." *IEEE transactions on visualization and computer graphics* 26, no. 1 (2019): 1064-1074.

[Nori2019] Nori, Harsha, Samuel Jenkins, Paul Koch, and Rich Caruana. "InterpretML: A Unified Framework for Machine Learning Interpretability." *arXiv preprint arXiv:1909.09223* (2019).

Bibliography

- [PIT2021] I. Pitas, “Computer vision”, Createspace/Amazon, in press.
- [PIT2017] I. Pitas, “Digital video processing and analysis ” , China Machine Press, 2017 (in Chinese).
- [PIT2013] I. Pitas, “Digital Video and Television ” , Createspace/Amazon, 2013.
- [NIK2000] N. Nikolaidis and I. Pitas, 3D Image Processing Algorithms, J. Wiley, 2000.
- [PIT2000] I. Pitas, “Digital Image Processing Algorithms and Applications”, J. Wiley, 2000.

Q & A

Thank you very much for your attention!

**More material in
<http://icarus.csd.auth.gr/cvml-web-lecture-series/>**

**Contact: Prof. I. Pitas
pitass@csd.auth.gr**