

Medical Image Analysis with Deep Learning

- Uncertainties - Explainability

M. Sdika ¹

¹Univ.Lyon, INSA-Lyon, Université Claude Bernard Lyon 1, UJM-Saint Etienne, CNRS, Inserm, CREATIS

UMR 5220, U1206, F-69100, LYON, France

Deep learning

- ▶ Deep network:
 - high performance
 - model developped from the data
 - black box, lack of transparency
- ▶ decision on human people health → need safeguard
 - interpretation, explaination, localization
 - uncertainties estimation

Source of uncertainties:

Aleatoric uncertainties:
task, randomness

Epistemic uncertainties:
model, learning setup

Input Noise:

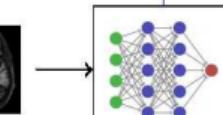
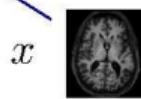
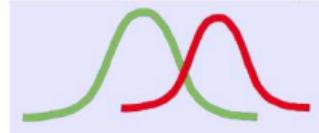
- Acquisition
- Pre/Post Processing

Architecture

Parameters:

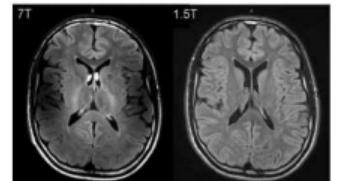
- improper training strategy
- lack of training instances

Feature Overlap



Annotation Noise

Domain Shift:
- train vs test



Source of uncertainties:

Aleatoric uncertainties:
task, randomness

Epistemic uncertainties:
model, learning setup

Input Noise:

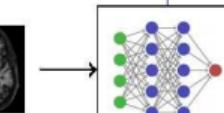
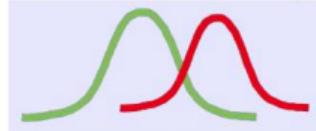
- Acquisition
- Pre/Post Processing

Architecture

Parameters:

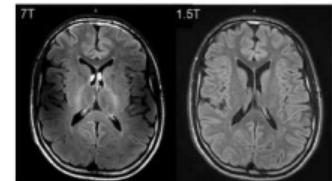
- improper training strategy
- lack of training instances

Feature Overlap

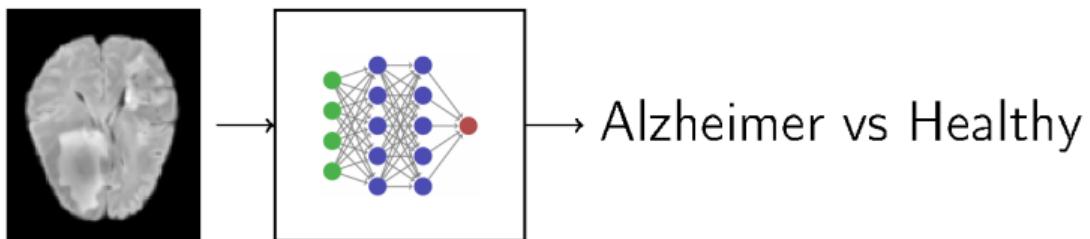


Annotation Noise

Domain Shift:
- train vs test



Domain Shift



Covariate shift:

- network input
- change in acquisition
- demographics

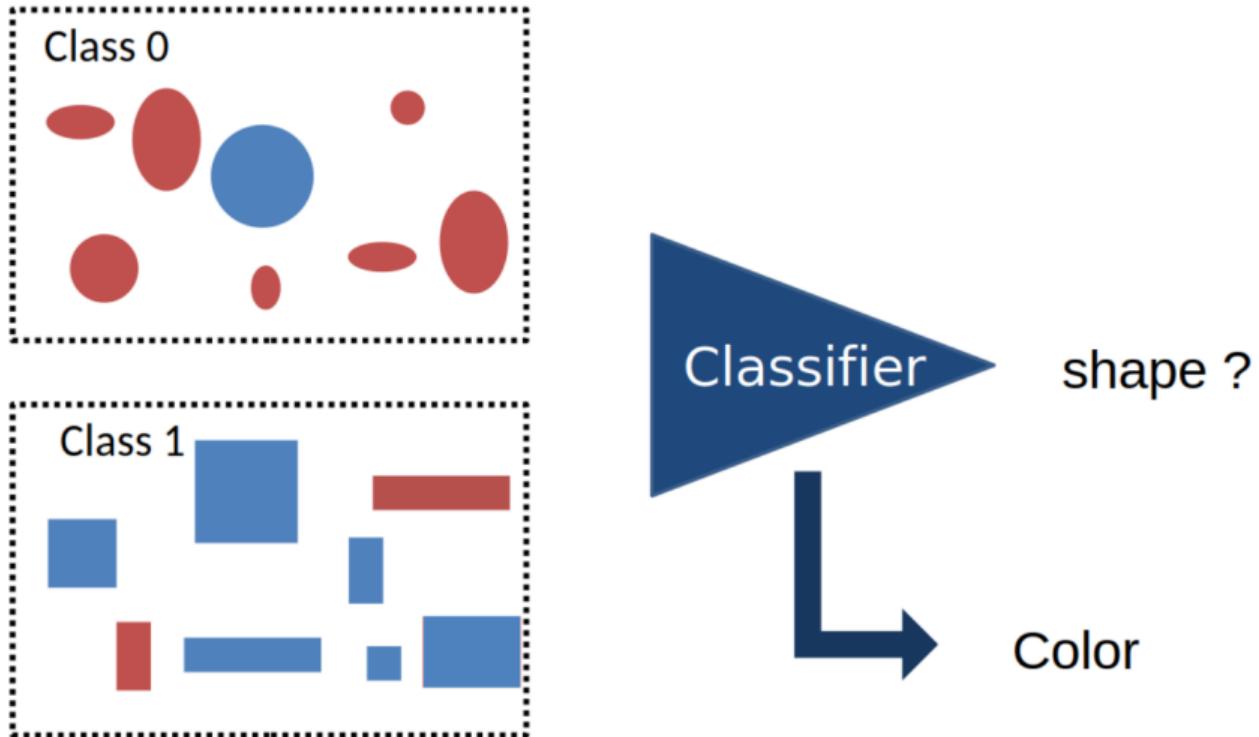
Label shift

- ### Concept Drift:
- labels repartition

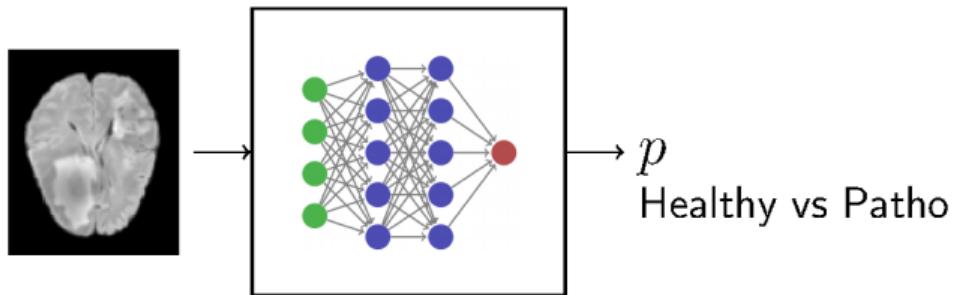
Semantic shift:

- network output
- input HIV+ patient
- OOD

Dataset biases, Confounder



Dataset biases, Confounder



Do the "Healthy" and "Patho" data match in:

- imaging parameters: scanner brand ? acquisition parameters ? ...
- population: age ? morphology ? ...
- something else ...

Content

Uncertainties Estimation

Calibration

Uncertainty Estimation for deep networks

Conformal Prediction

Out of Distribution Detection (OOD)

Right for the Right Reason ?

Attribution Maps

Counterfactual Explanation

Confounder

Content

Uncertainties Estimation

Calibration

Uncertainty Estimation for deep networks

Conformal Prediction

Out of Distribution Detection (OOD)

Right for the Right Reason ?

Attribution Maps

Counterfactual Explanation

Confounder

Content

Uncertainties Estimation

Calibration

Uncertainty Estimation for deep networks

Conformal Prediction

Out of Distribution Detection (OOD)

Right for the Right Reason ?

Attribution Maps

Counterfactual Explanation

Confounder

Scoring Rule

Scoring rule: $S(p, y)$

measures how well an estimated probability vector p explains the observed labels y

Strictly proper scoring rule:

- proper: minimal for the true distribution of y
- strictly proper: the true distribution of y is the only minimum

Example:

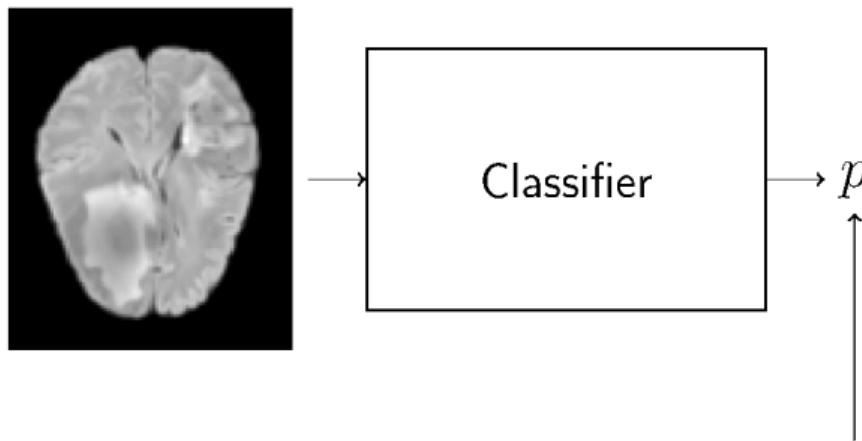
- cross entropy (CE),
negative log likelihood (NLL)

$$-\log(p_y)$$

- Brier Score:

$$\frac{1}{C} \sum_c (p_c - y_c)^2$$

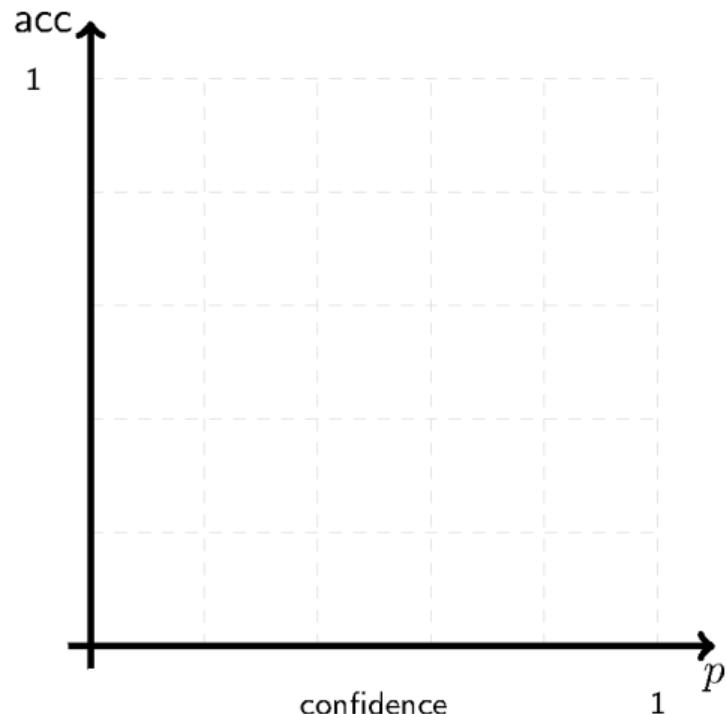
The network output as an uncertainty measure ??



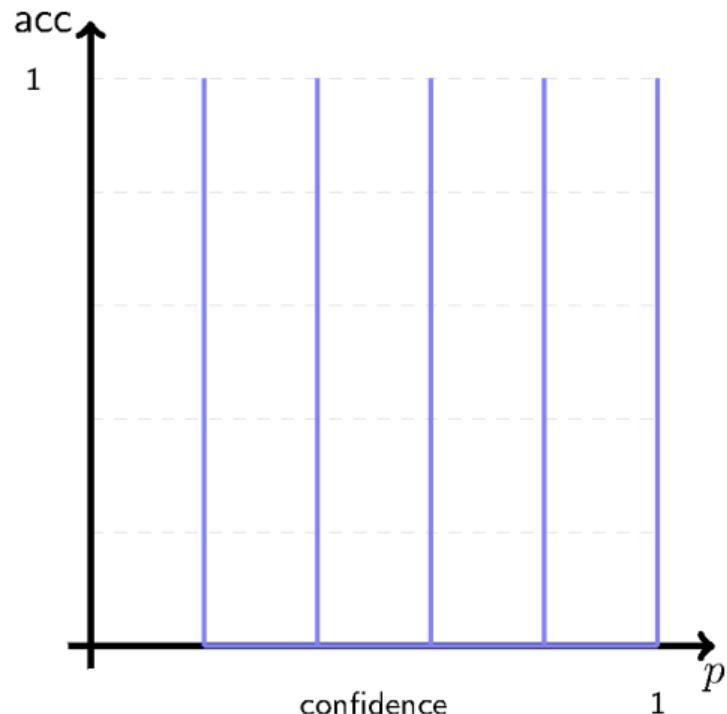
Is it the probability to give a correct output ?

→ **network confidence**

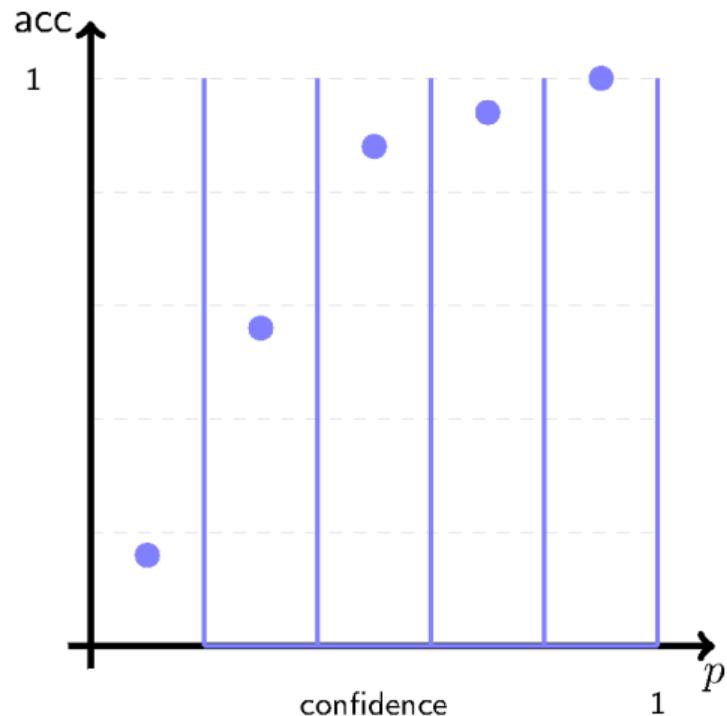
Reliability Diagram



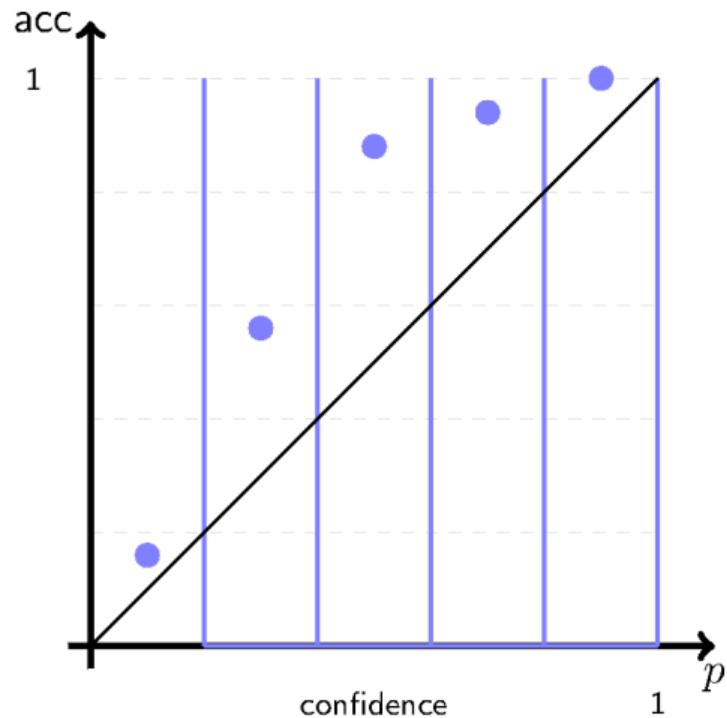
Reliability Diagram



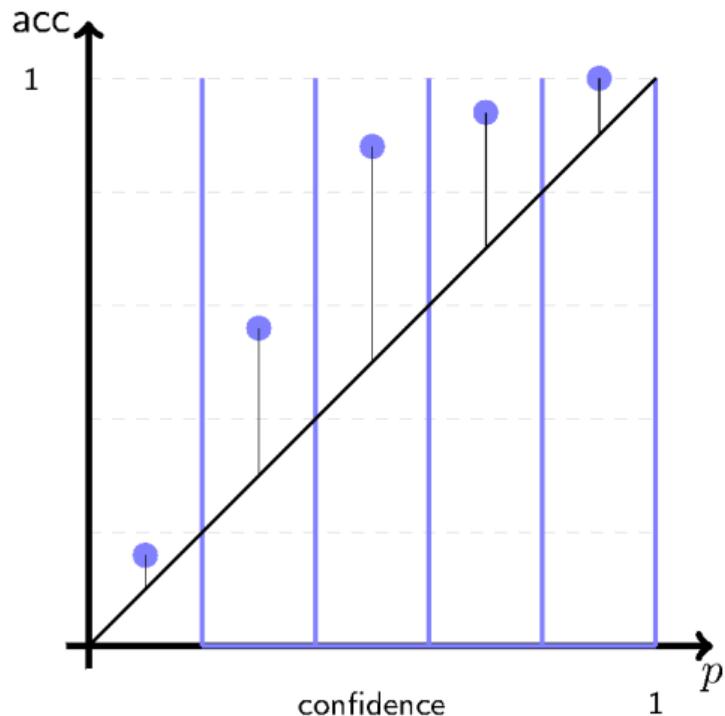
Reliability Diagram



Reliability Diagram



Reliability Diagram



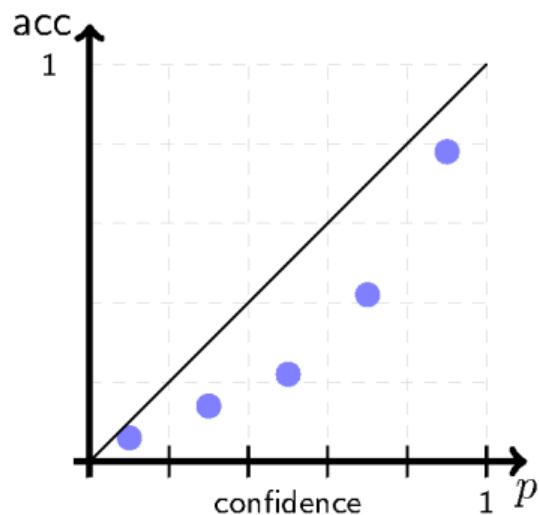
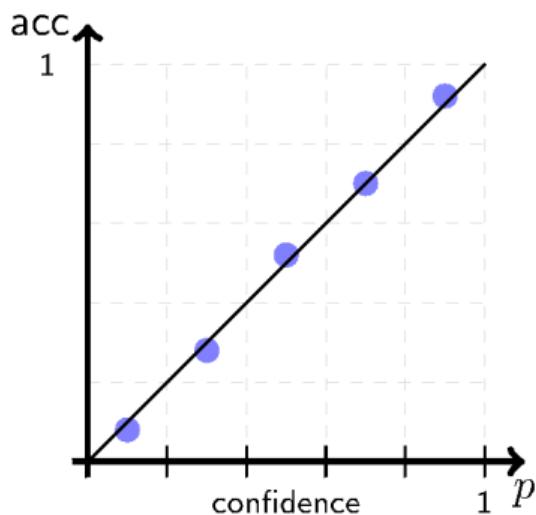
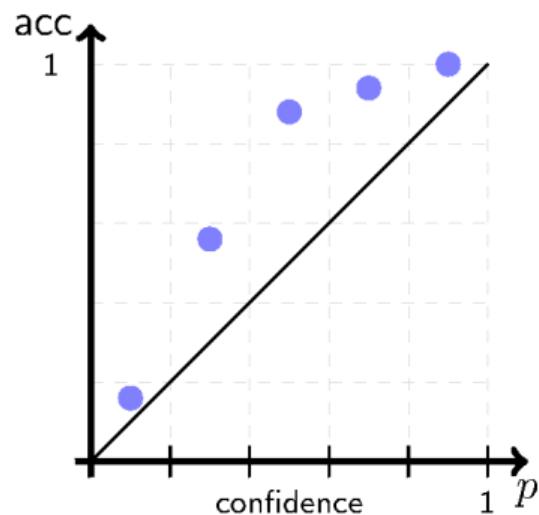
Expected Calibration Error (ECE)

$$\text{ECE} = \frac{1}{B} \sum_{b=1}^B |\text{conf}_b - \text{acc}_b|$$

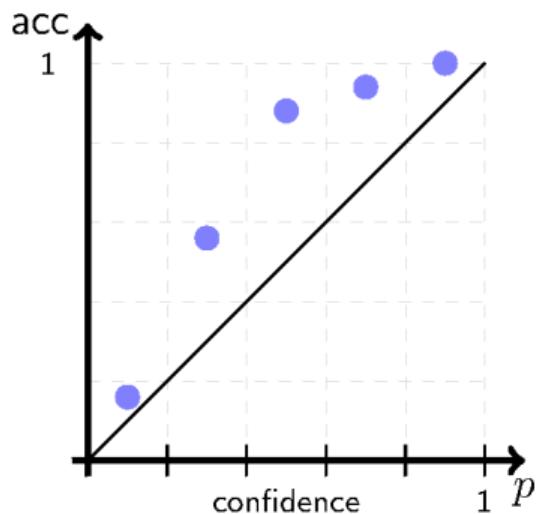
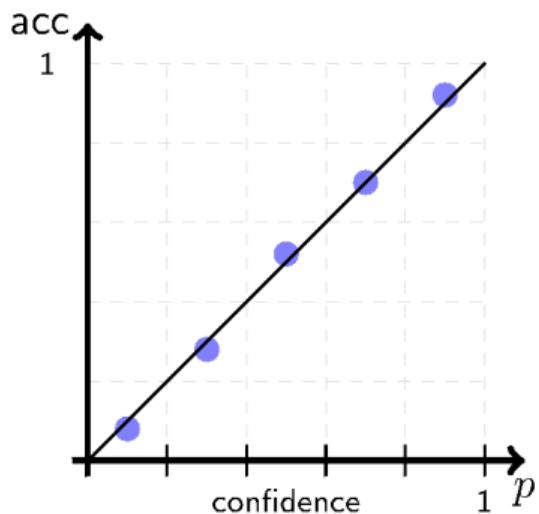
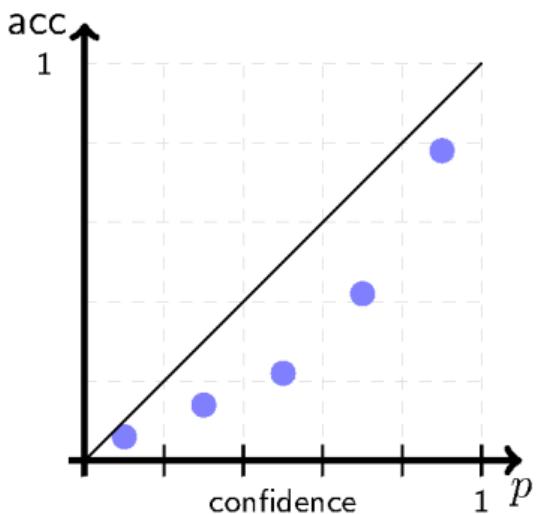
Multi-class, either:

- avg accuracy VS avg confidence
- classwize reliability diagram

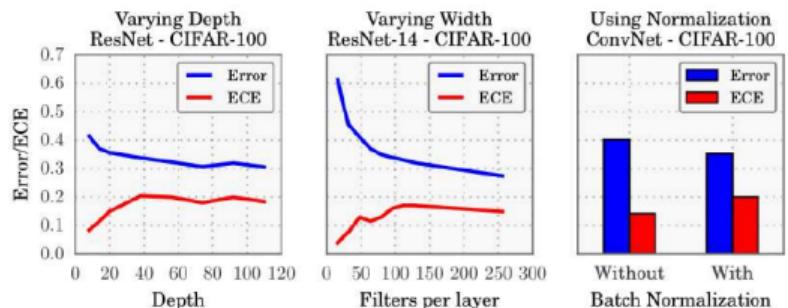
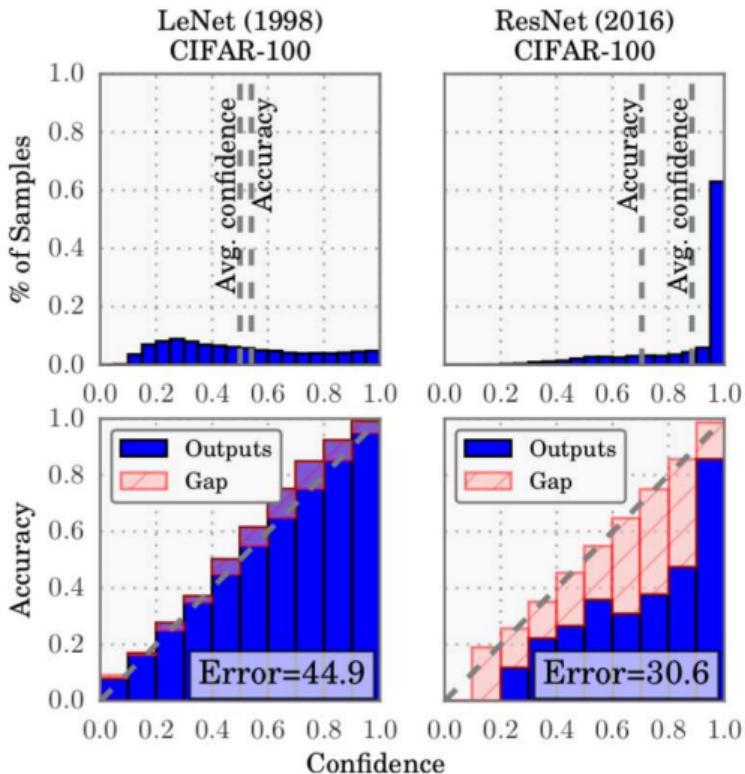
Reliability Diagram



Reliability Diagram

**under-confident****calibrated****over-confident**

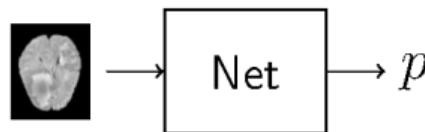
Modern deep networks are not calibrated



cause:

- depth/width/batch norm
- Cross Entropy loss:
 - $-\log(p_y)$

Label Smoothing



y GT label
 correct with probability $1 - \epsilon$

$$CE = - \sum_i y_i \log(p_i)$$

- standard : $y_i = 0/1$
- label smoothing: $y_i = \frac{\epsilon}{K-1} / 1 - \epsilon$

Szegedy et al, CVPR 2016, Rethinking the inception architecture for computer vision

Deep Learning, Goodfellow et al

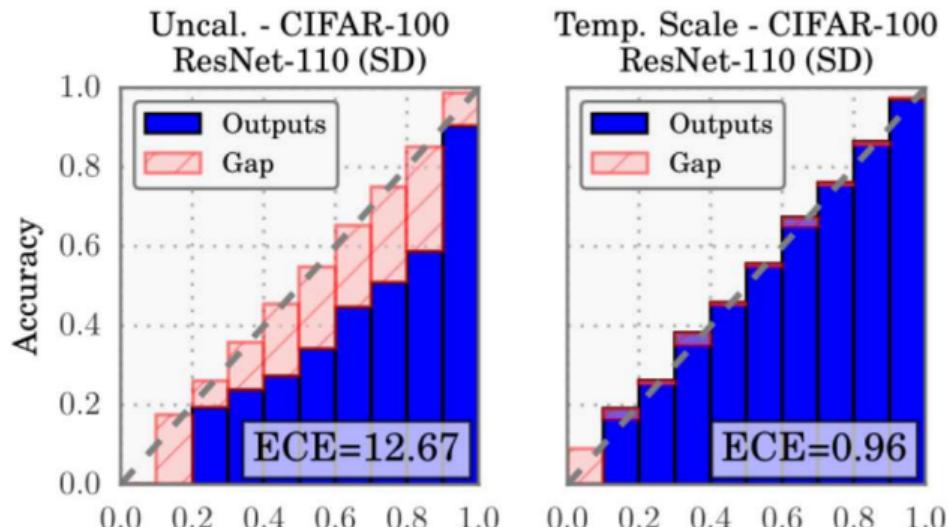
Temperature scaling

Softmax :

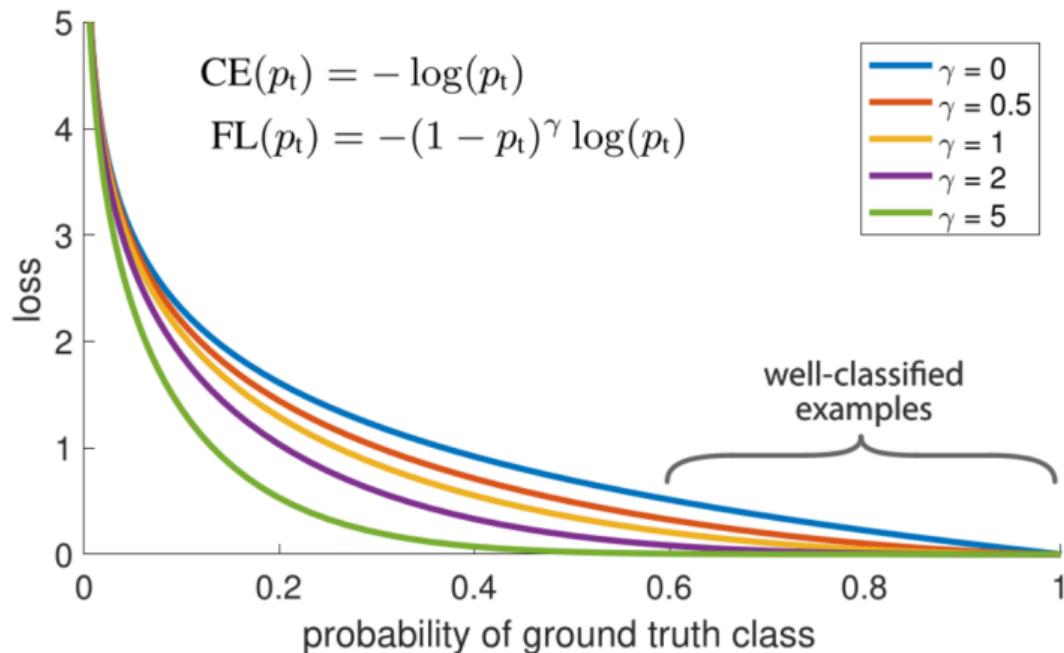
$$\frac{e^{z_k}}{\sum e^{z_i}} \rightarrow \frac{e^{z_k/T}}{\sum e^{z_i/T}}$$

On the validation dataset:

$$\min_T \text{NLL}(T)$$

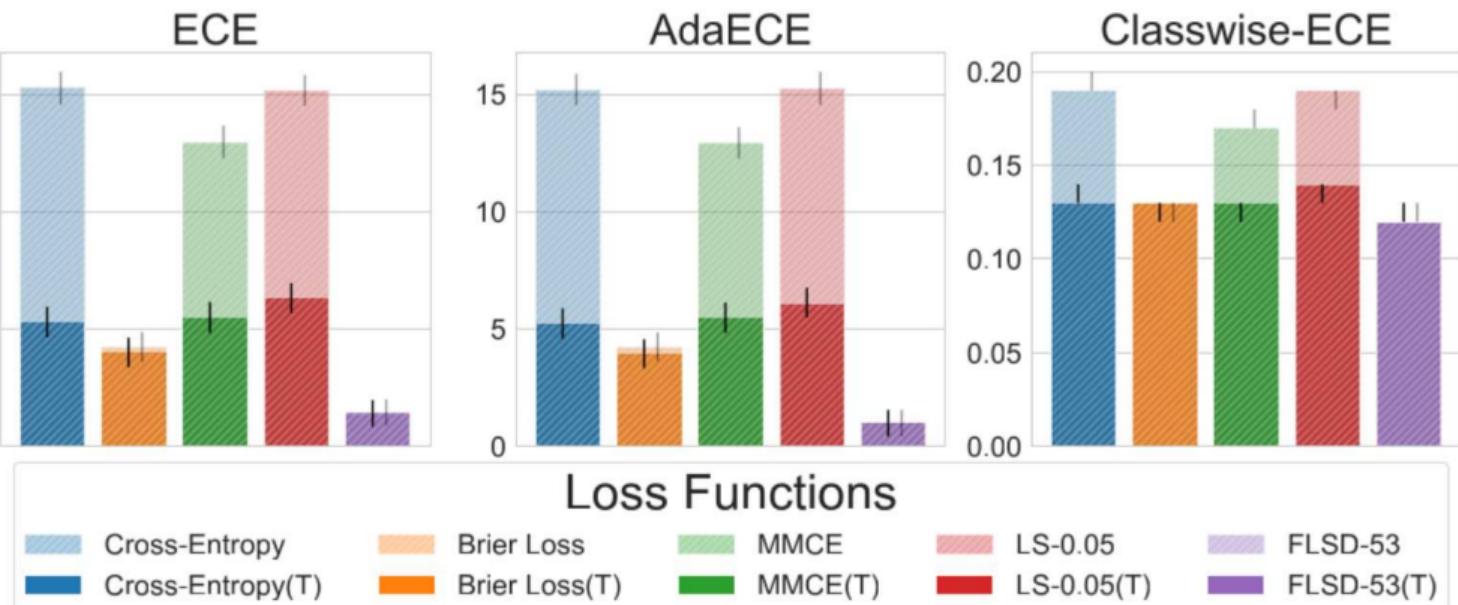


Focal Loss



Lin et al, ICCV 2017. Focal loss for dense object detection
Mukhoti et al Neurips 2020. Calibrating deep neural networks using focal loss

Focal Loss



Content

Uncertainties Estimation

Calibration

Uncertainty Estimation for deep networks

Conformal Prediction

Out of Distribution Detection (OOD)

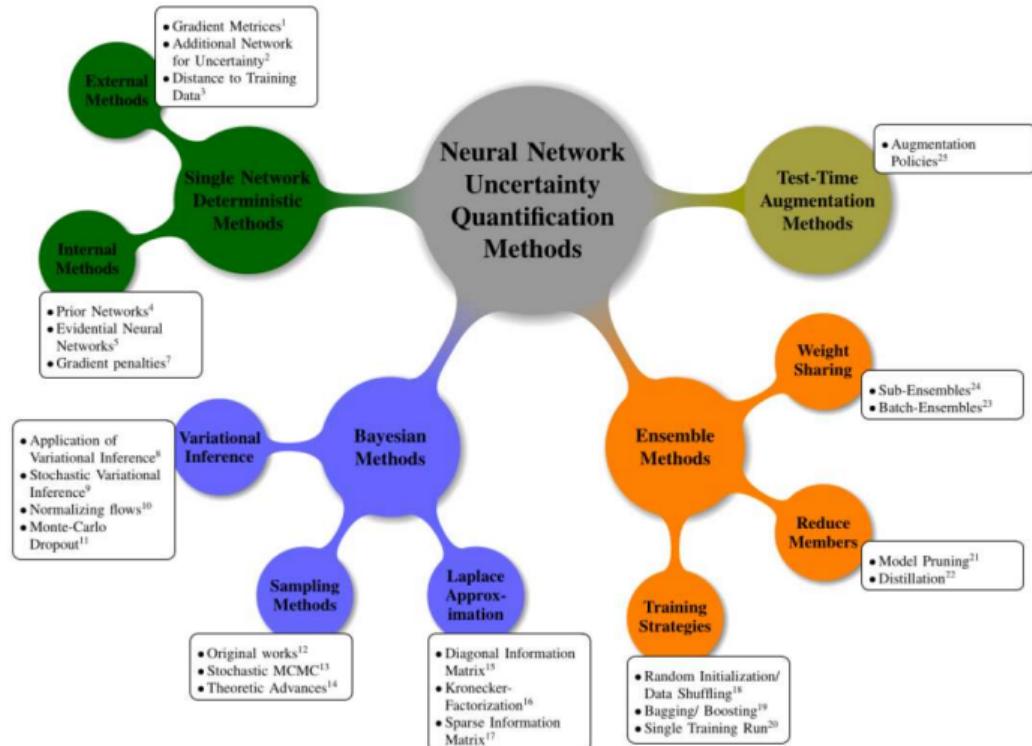
Right for the Right Reason ?

Attribution Maps

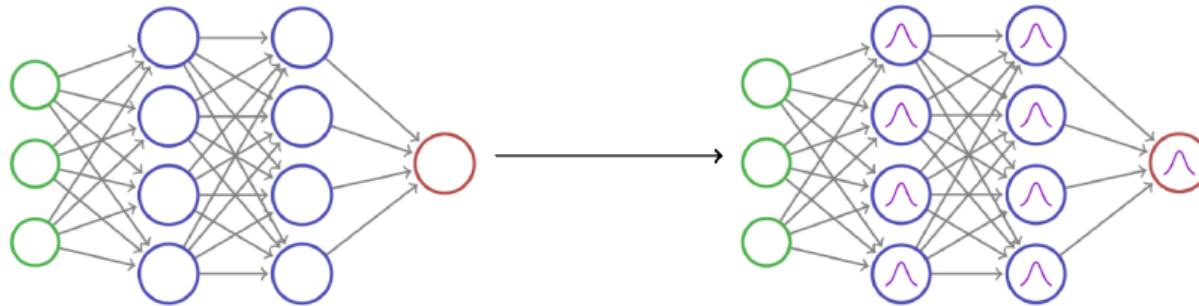
Counterfactual Explanation

Confounder

Uncertainties approaches



Bayesian Networks

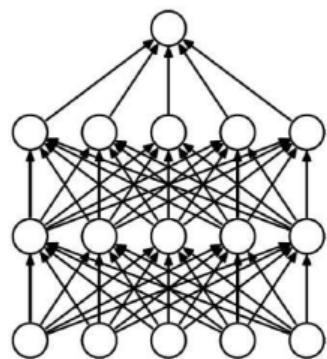
**Inference:** MC sampling

- D : training dataset
- draw random $w_k \sim q(w)$
- $p(y|x, D) \approx \frac{1}{K} \sum_{k=1}^K p(y|x, w_k)$

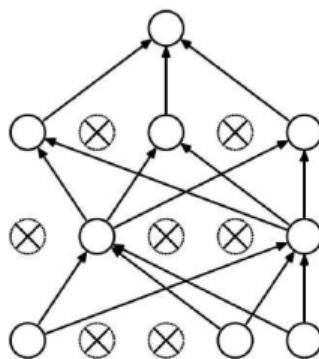
Uncertainty Estimation:

- $p^c = p(y = c|x, w_k)$
- $H[y|x, D] \approx - \sum_c p^c \log(p^c)$

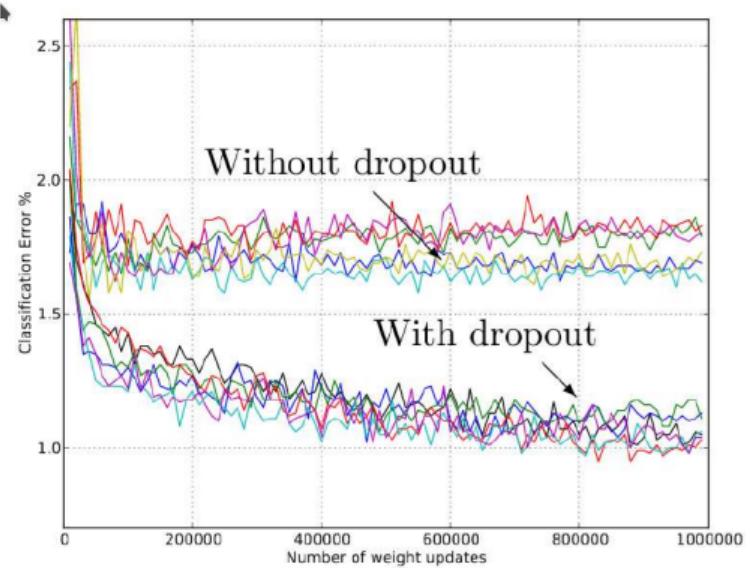
Dropout Regularization



(a) Standard Neural Net



(b) After applying dropout.



Random set units to 0 during training

Test error with/without dropout

MC Dropout for uncertainty estimation

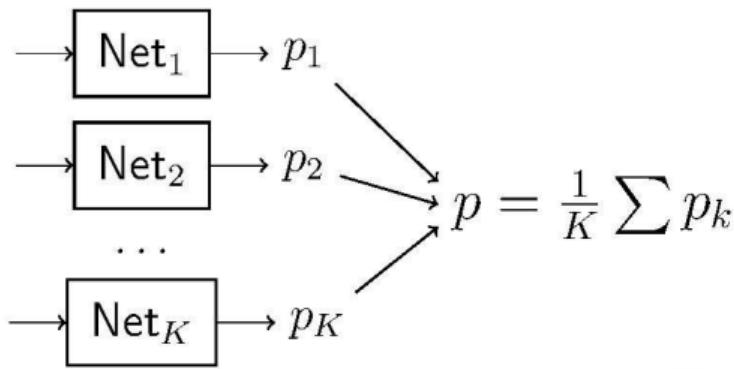
network: $x \rightarrow y = y(x)$

- ▶ dropout during inference: sample several outputs $(y_i(x))_{1,N}$
- ▶ Gal 2016: these samples are \approx drawn from $p(\hat{y}|x^*)$
- ▶ estimate first and second moments of $p(y|x)$

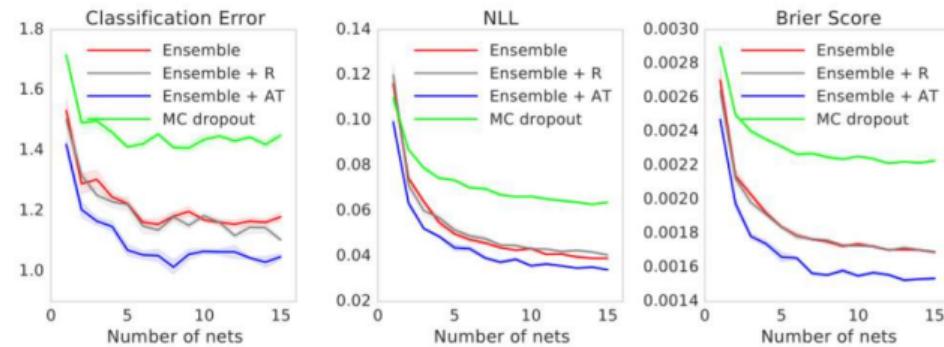
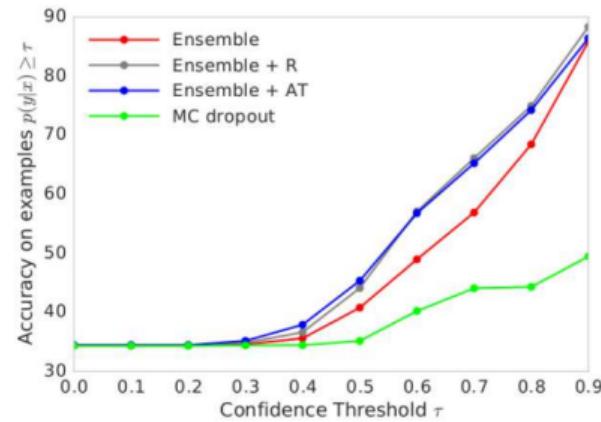
$$E(y|x) \approx \text{SampleMean}(y_i)$$

$$\text{Var}(y|x) \approx \tau I + \text{SampleVar}(y_i)$$

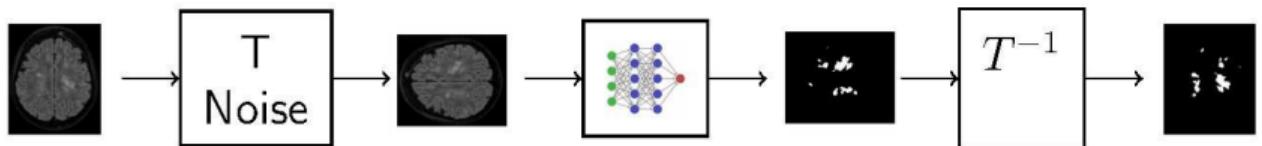
Deep ensemble



- train a set of network
- stochasticity from random init
- in practice, $K=5$



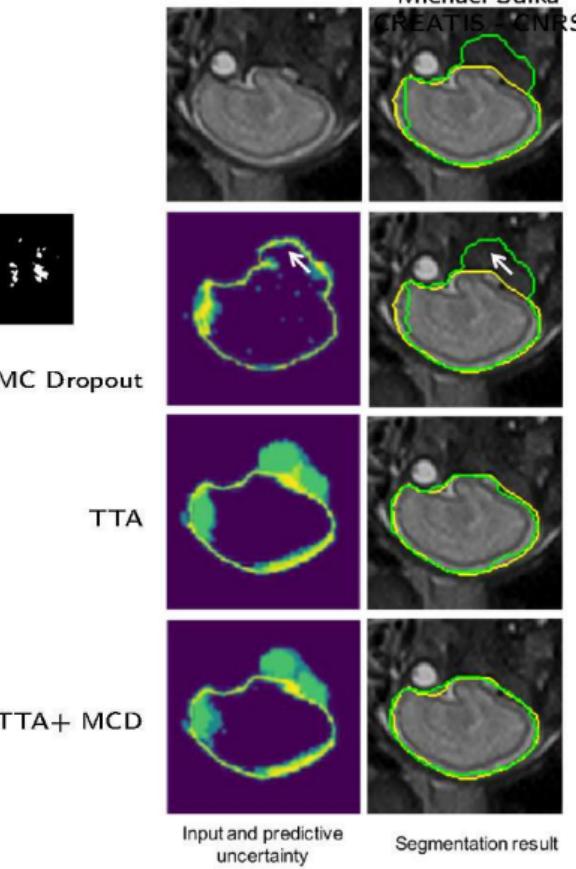
Test Time Augmentation (TTA)



During test:

- generate multiple output
- prediction: mean / mode
- uncertainty: std-dev / entropy

aleatoric uncertainty related to
image transformations and noise



Content

Uncertainties Estimation

Calibration

Uncertainty Estimation for deep networks

Conformal Prediction

Out of Distribution Detection (OOD)

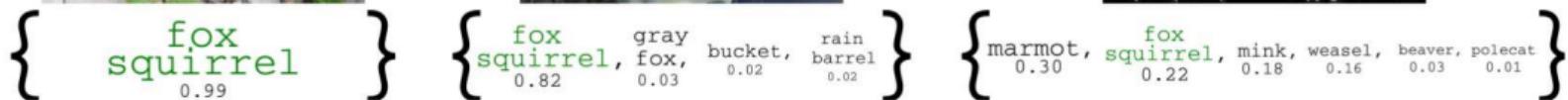
Right for the Right Reason ?

Attribution Maps

Counterfactual Explanation

Confounder

Conformal Prediction



for a test image I , find a set of labels $C(I)$ s.t.:

$$P(y \in C(I)) \geq 1 - \alpha$$

Conformal Prediction: Solving the problem

- ▶ heuristic uncertainty measure for the model ex: $\rightarrow \hat{f}(x)$
- ▶ Define a score function $s(x, y)$ (lower=better) ex: $\rightarrow s(x, y) = 1 - \hat{f}_y(x)$
- ▶ Define the set C_q :

$$C_q(x) = \{y \text{ s.t } s(x, y) < q\}$$

- ▶ Notice that

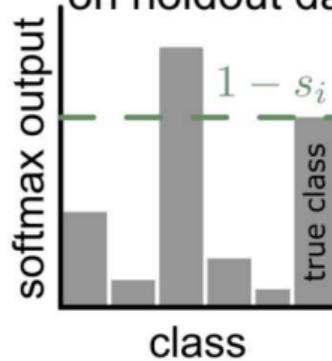
$$P(y \in C_q(x)) = P(s(x, y) < q) = F_S(q)$$

is the cumulative distribution function of $s(X, Y)$

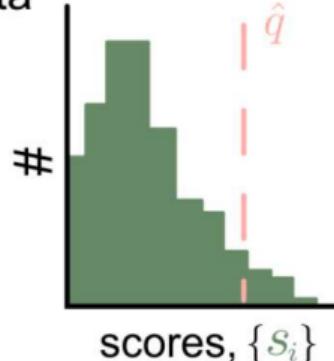
$$\Rightarrow \hat{q} \approx (1 - \alpha)\text{-quantile}$$

Conformal Prediction: Solving the problem

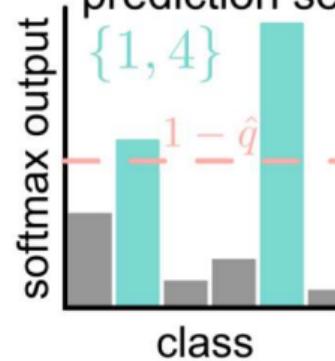
(1) compute scores
on holdout data



(2) get quantile



(3) construct
prediction set



Calibration data:

- conformal score:

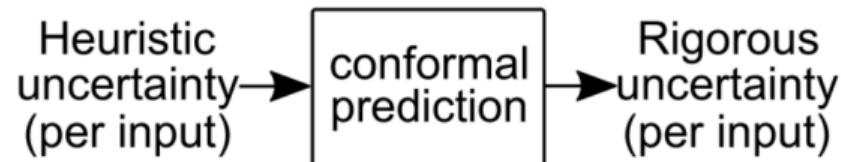
$$s_i = 1 - f_{y_i}(x_i)$$

- $q \approx (1 - \alpha)$ -quantile of scores

$x \in \text{TEST}$ (y unknown):

$$C(x) = \left\{ y \mid 1 - \hat{f}_y(x) \leq q \right\}$$

Conformal Prediction



► Limits

- bad model / bad heuristic uncertainties / small calib set
→ return all the labels
- sensitive to domain shift between calib / test

Angelopoulos & Bates, "A gentle introduction to conformal prediction and distribution-free uncertainty quantification." arXiv 2021

Content

Uncertainties Estimation

Calibration

Uncertainty Estimation for deep networks

Conformal Prediction

Out of Distribution Detection (OOD)

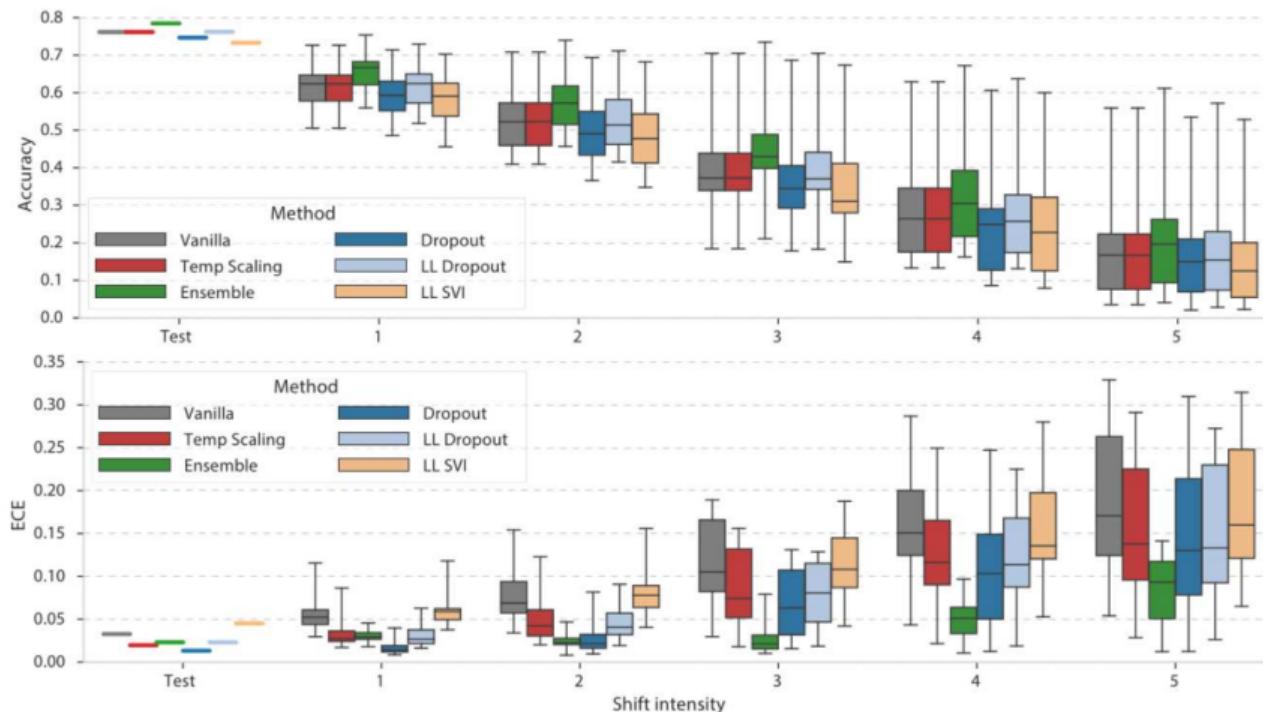
Right for the Right Reason ?

Attribution Maps

Counterfactual Explanation

Confounder

Uncertainties estimation is not robust to domain shift



Ovadia, et al. 2019 "Can you trust your model's uncertainty? evaluating predictive uncertainty under dataset shift."

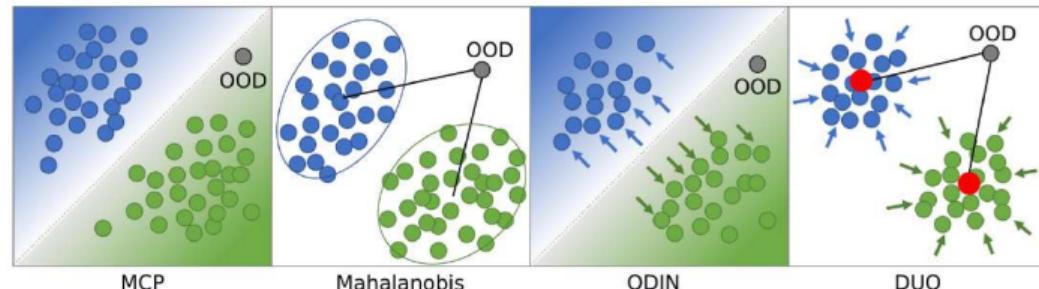
Some confidence and distance based OOD

Maximum Class Probability

ID : $\max_c p_c > \text{thresh}$

MCP with better calibration:

- MCP Dropout
- MCP Deep ensemble



ODIN

- \tilde{x} : adv attack to increase confidence
- adjust temperature (G-ODIN:large)
- score : MCP (\tilde{x})

Mahalanobis

DUQ

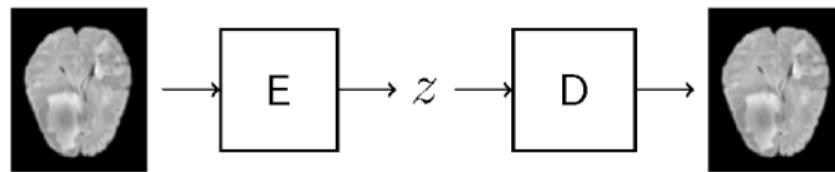
- feature learning
- contrastive learning
- score : RBF kernel distance to centroid

KNN

- contrastive
- KNN

Berger, et al. UNSURE MICCAI 2021, Confidence-based out-of-distribution detection: a comparative study and analysis
 Hendrycks & Gimpel, ICLR 2017, A baseline for detecting misclassified and out-of-distribution examples in neural networks
 Liang et al ICLR 2018: Enhancing the reliability of out-of-distribution image detection in neural networks
 Van Amersfoort et al, ICML 2020 : Uncertainty estimation using a single deep deterministic neural network
 Sun, et al, ICML 2022, Out-of-distribution detection with deep nearest neighbors

Reconstruction-based OOD

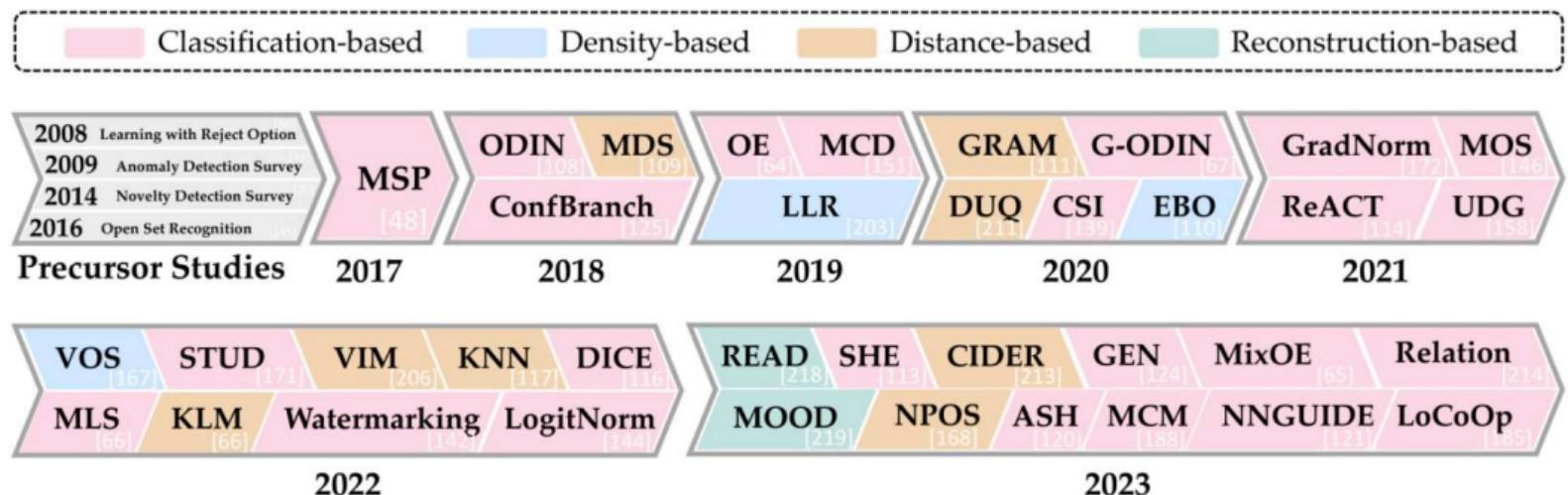


Assumption: better reco for ID than OOD

score: $\|\hat{x} - x\|$

$$\|\hat{x} - x\| + \lambda(z - \mu)\Sigma^{-1}(z - \mu)$$

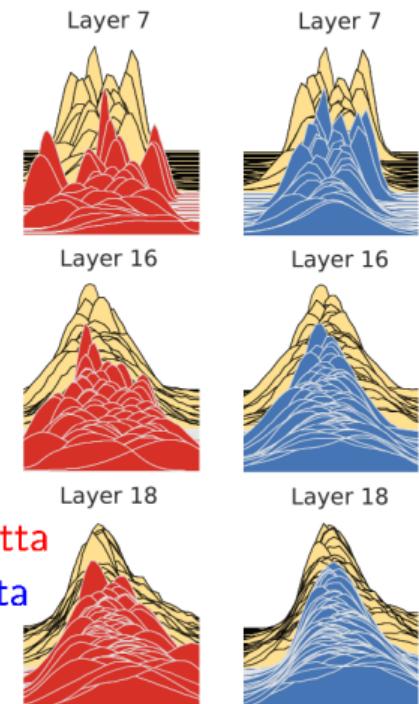
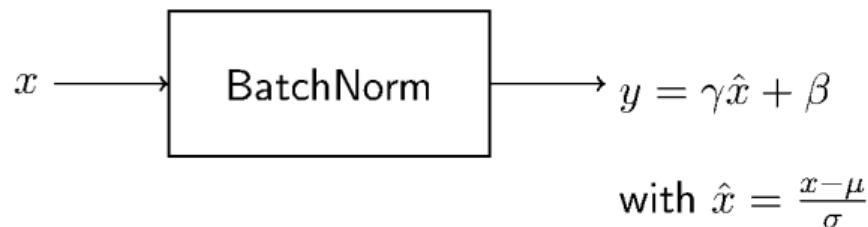
OOD



Yang et al., arXiv, Generalized out-of-distribution detection: A survey

Zhang et al., arxiv, OpenOOD v1.5: Enhanced Benchmark for Out-of-Distribution Detection

Test Time Adaptation with BatchNorm



TTA flavors:

- adjust μ/σ on each batch at test
- adjust on full test domain
- fine tune γ/β
with self sup. / consistency loss

Content

Uncertainties Estimation

Calibration

Uncertainty Estimation for deep networks

Conformal Prediction

Out of Distribution Detection (OOD)

Right for the Right Reason ?

Attribution Maps

Counterfactual Explanation

Confounder

Content

Uncertainties Estimation

Calibration

Uncertainty Estimation for deep networks

Conformal Prediction

Out of Distribution Detection (OOD)

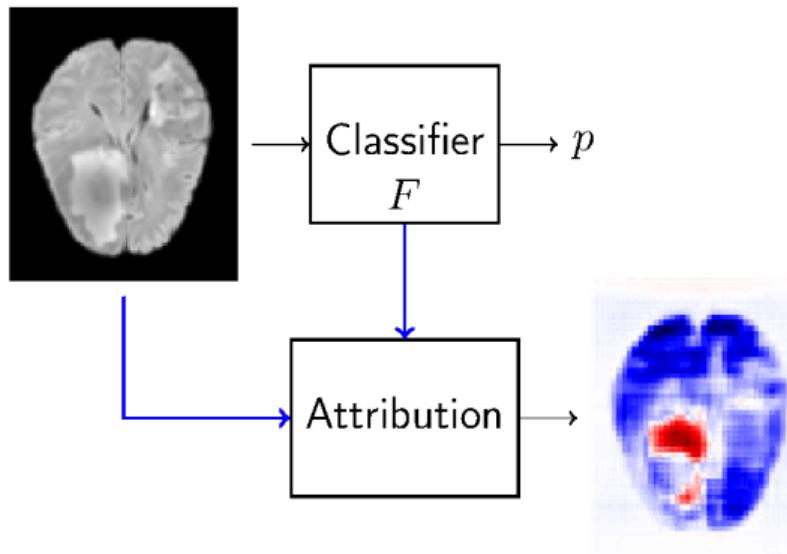
Right for the Right Reason ?

Attribution Maps

Counterfactual Explanation

Confounder

Attribution maps: WHERE, not WHAT



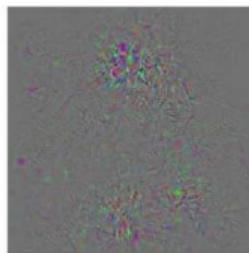
Post hoc: apply it once the network has been trained

blue: push toward negative
red : push toward positive

Lots of methods, several approaches



image



gradients

guided
backprop.

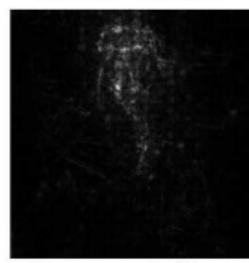
Grad-CAM

guided
Grad-CAM,

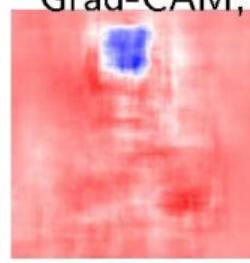
Score-CAM



FullGrad

Integrated
Gradients

SmoothGrad



occlusion

Perturbation

//

Gradient

//

Activation

Occlusion

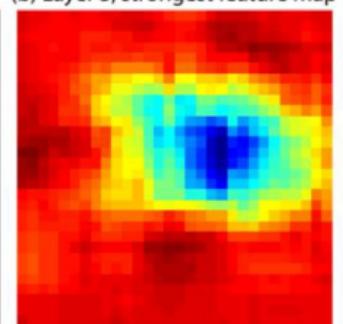


(a) Input Image

(b) Layer 5, strongest feature map

$$A(I, B)(x) = F(I) - F(I^x)$$

I^x : baseline value at pixel x
 otherwise I



- I^x might be OOD
- untractable \Rightarrow patch, superpixel...

Saliency maps: the gradient as an attribution maps

Saliency Maps / Gradient:

$$A(I) = \nabla_I F$$

Input x Gradient:

$$A(I) = (I - B)\nabla_I F$$

$$y = w_1 f_1 + w_2 f_1 + w_3 f_3$$

G: w_1 (local)
 IxG: $w_1 f_1$ (global)



saliency map

Simonyan et al, Deep inside convolutional networks: Visualising image classification models and saliency maps. ICLR 2014
 Shrikumar et al, Learning important features through propagating activation differences. ICML 2017

Path Methods: Integrated Gradient

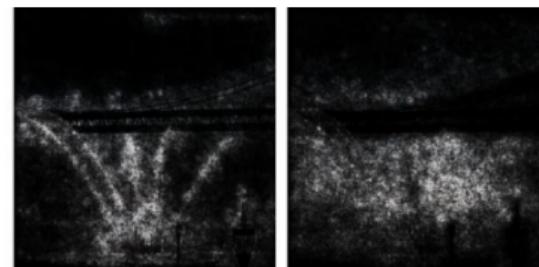
$$A(I, B) = (I - B) \int_0^1 \frac{\partial F}{\partial I} (B + \alpha(I - B)) d\alpha$$

satisfy: Linearity w.r. model, Sensitivity, Implementation invariance

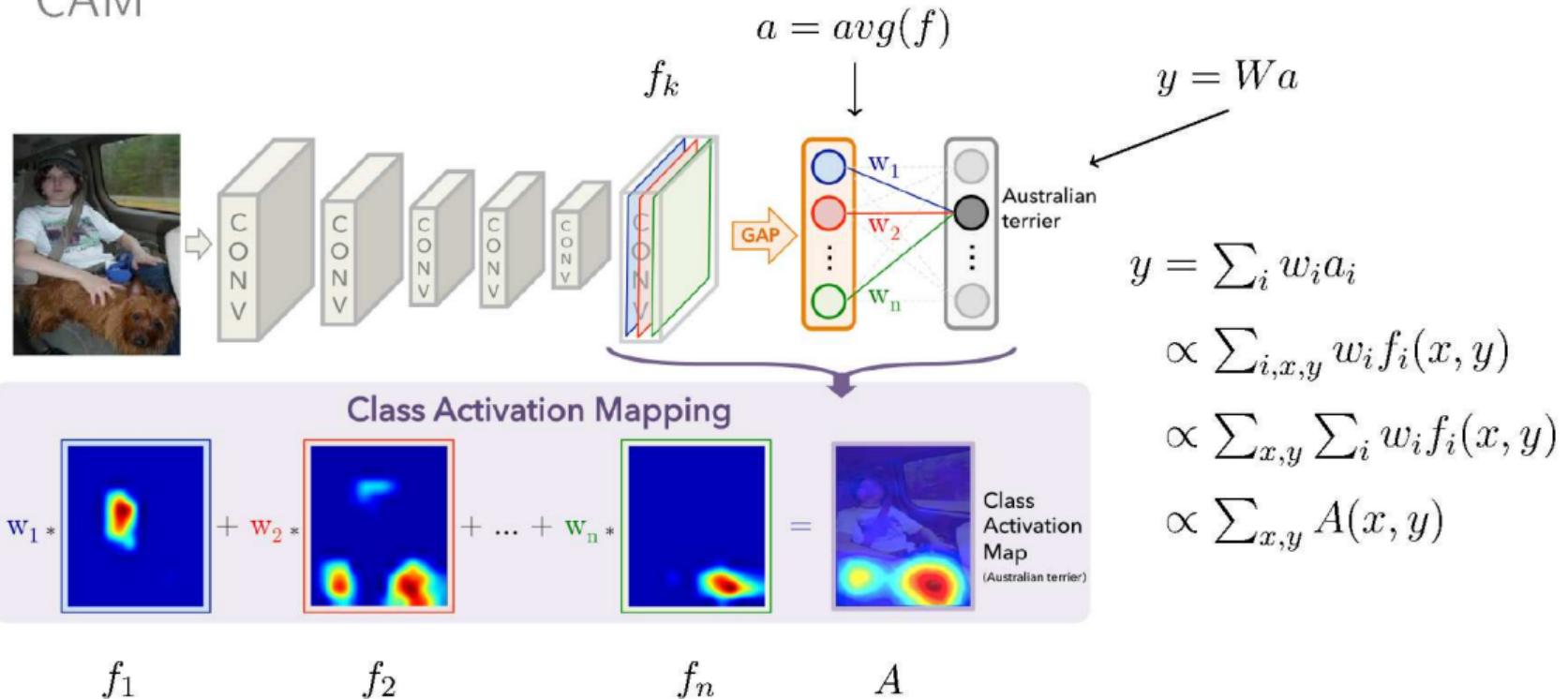
Completeness ($\sum_x A(I)_x = F(I) - F(B)$)



Top label: fireboat
Score: 0.999961



CAM



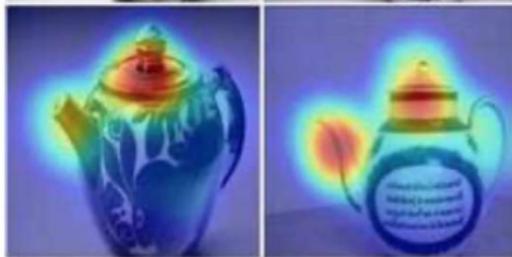
Mushroom



Penguin

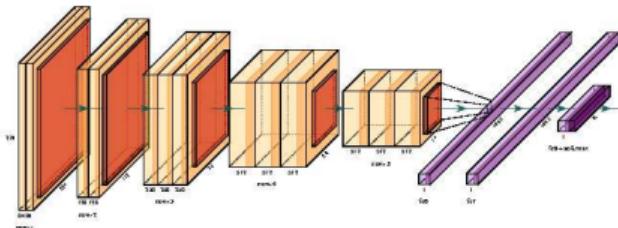


Teapot



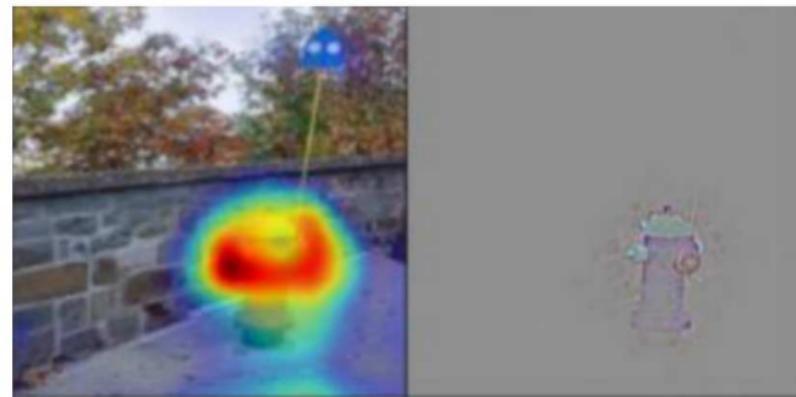
Grad-CAM

Can be used on any layer
(usually last conv)



$$w_k = \text{avg}_X \left(\frac{\partial y}{\partial f_k} \right)$$

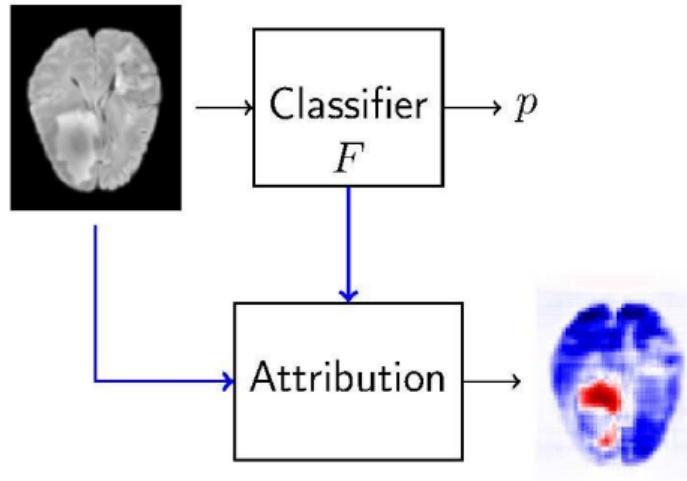
$$A(I) = \text{relu} \left(\sum_k w_k f_k(X) \right)$$



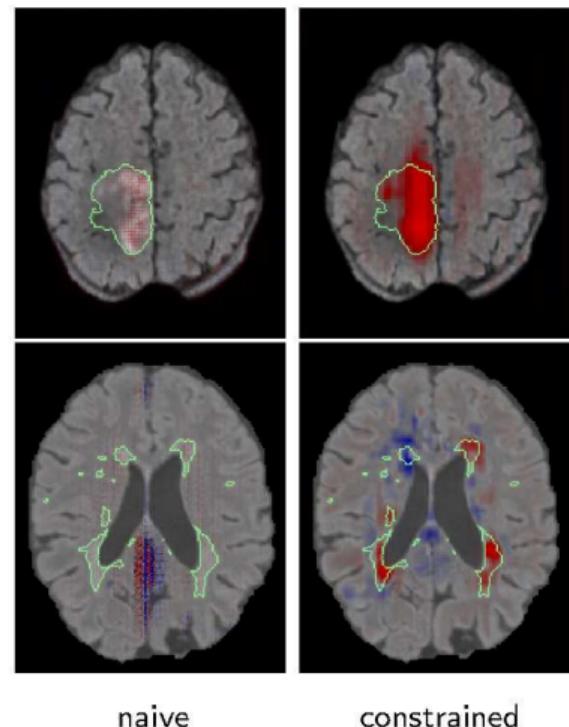
GradCAM

Guided GradCAM

Attribution maps constrained training



Constraint:
 I is Healthy $\Rightarrow A(I)$ is blue



Content

Uncertainties Estimation

Calibration

Uncertainty Estimation for deep networks

Conformal Prediction

Out of Distribution Detection (OOD)

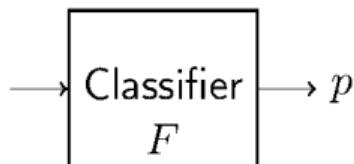
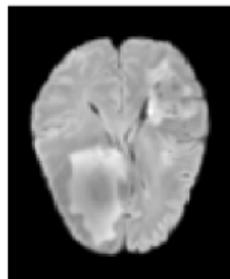
Right for the Right Reason ?

Attribution Maps

Counterfactual Explanation

Confounder

Counterfactual Explanation



Minimal changes in the input
data to switch class ?

Naive:

$$\min_{F(x')=1-y} \|x' - x\|$$

Naive = Adversarial attack



x
 “panda”
 57.7% confidence

$$+ .007 \times$$



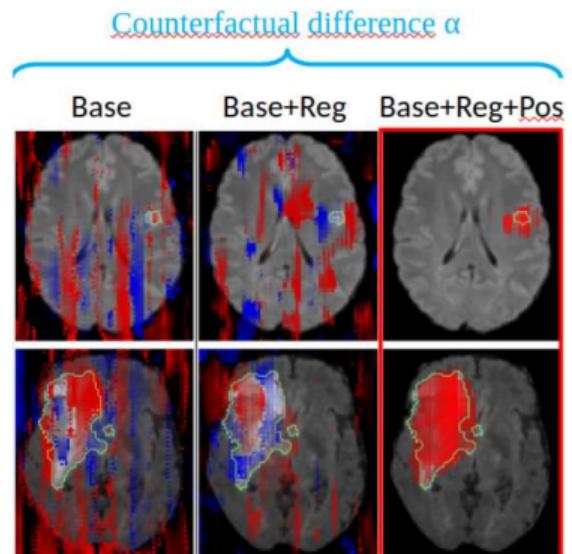
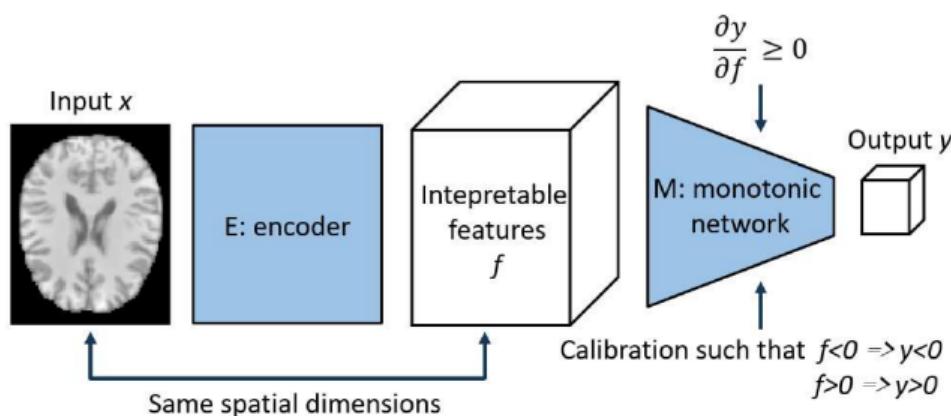
$\text{sign}(\nabla_x J(\theta, x, y))$
 “nematode”
 8.2% confidence

$$=$$



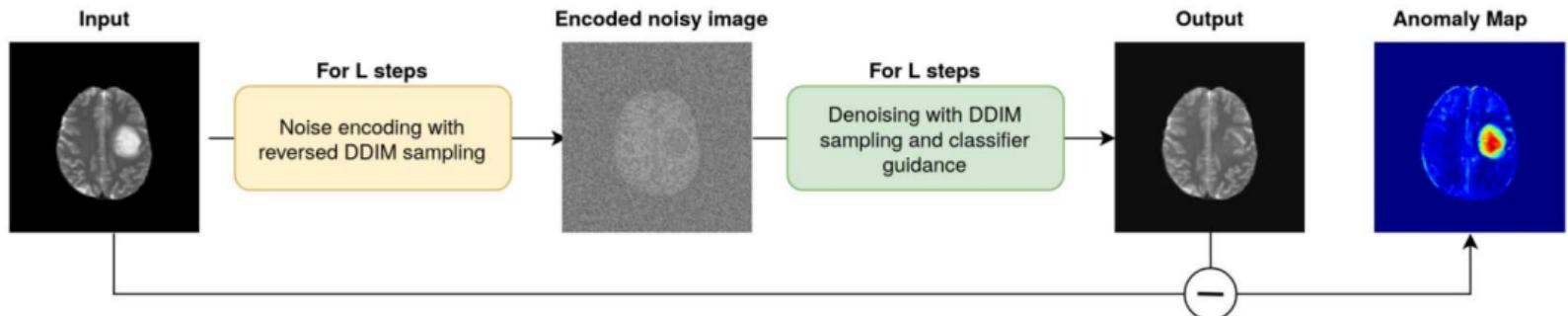
$x +$
 $\epsilon \text{sign}(\nabla_x J(\theta, x, y))$
 “gibbon”
 99.3 % confidence

Counterfactual explanation with a constrained architecture



$$\min_{\alpha} M(f - \alpha) + \lambda \|\alpha\|_1$$

Using Generative Models



counterfactual example s.t.:

- ▶ minimal change from input
- ▶ classif: opposite class
- ▶ output of a generative model

Content

Uncertainties Estimation

Calibration

Uncertainty Estimation for deep networks

Conformal Prediction

Out of Distribution Detection (OOD)

Right for the Right Reason ?

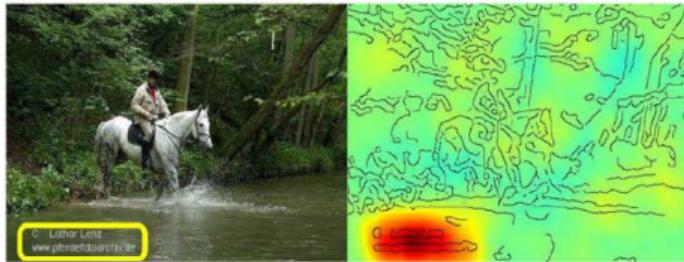
Attribution Maps

Counterfactual Explanation

Confounder

Counfounder effect: Right for the wrong reason

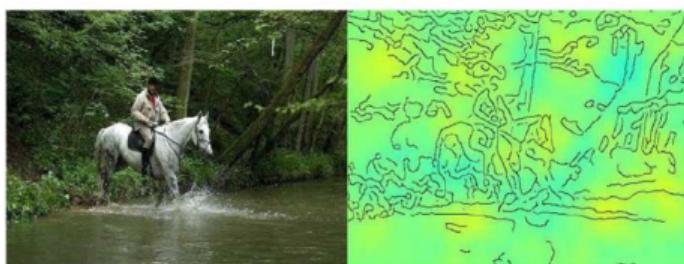
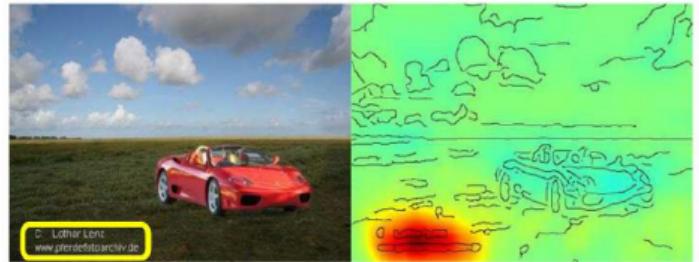
Horse-picture from Pascal VOC data set



Source tag
present

↓
Classified
as horse

Artificial picture of a car



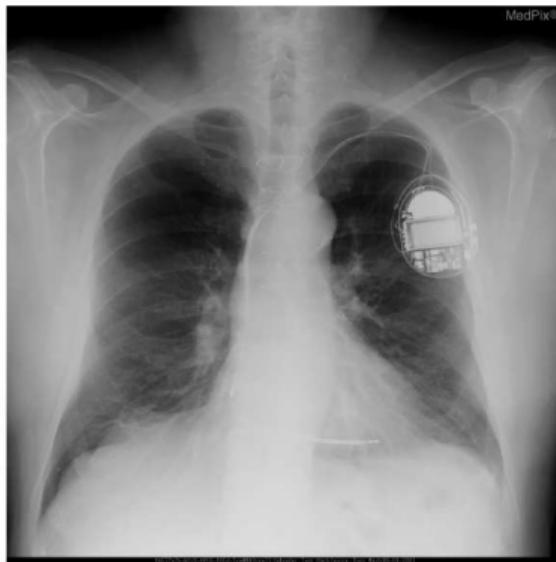
No source
tag present

↓
Not classified
as horse



Lapuschkin et al. "Unmasking Clever Hans predictors and assessing what machines really learn." Nature communications 2019

Obvious Confounder in Medical Images

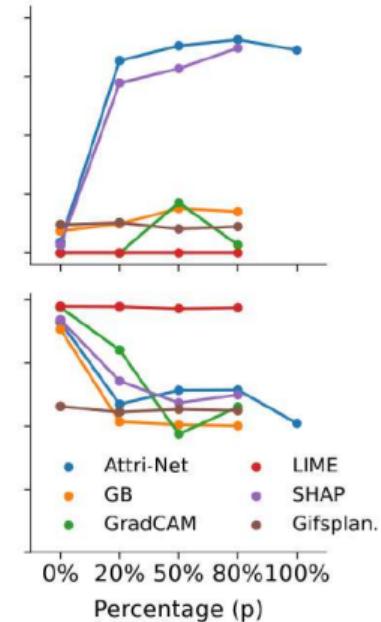
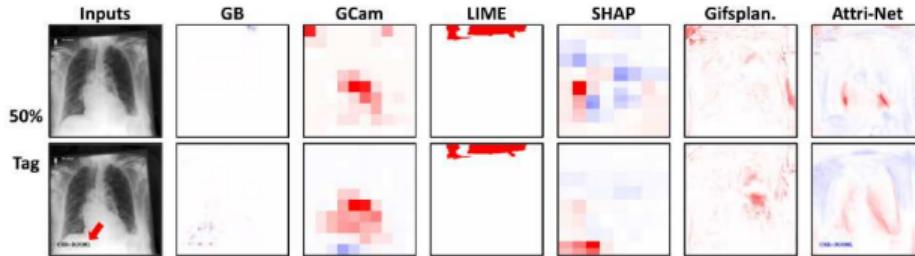
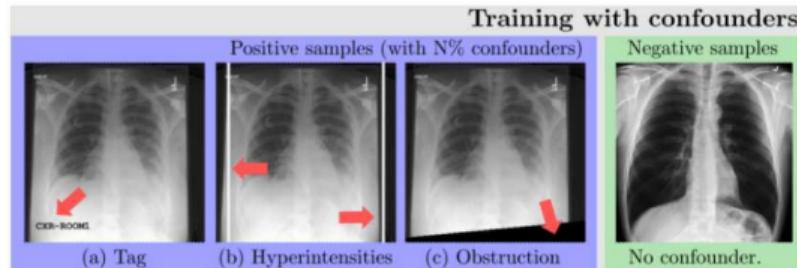


pacemaker on chest radiograph



ruler on skin lesion images

Attribution maps comparison with artificial confounders



The brain shape as a confounder

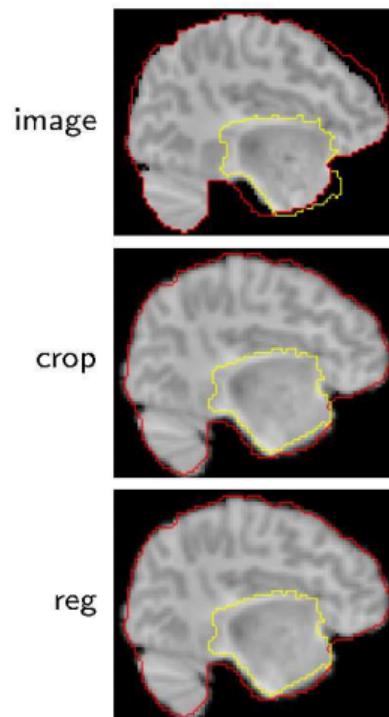


	Classification task		TNR	TPR	BA
Healthy vs Pathological	IXI	vs	BraTS	1.00	1.00 ✓
	IXI	vs	OFSEP	0.54	0.70 ✓
	IXI	vs	HCP	1.00	1.00 ✓
Healthy vs Healthy	IXI	vs	IBC/kirby/MPI	1.00	0.86 ✓
	HCP	vs	IBC/kirby/MPI	0.99	0.95 ✓
Healthy vs Pathological	ADNI-CN	vs	ADNI-AD	0.35	0.74
				0.55 X	Intra dataset

Age matching

Intra dataset

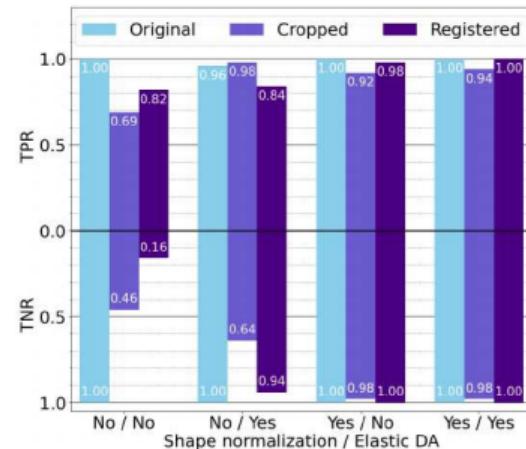
The brain shape as a confounder: assessment and solution



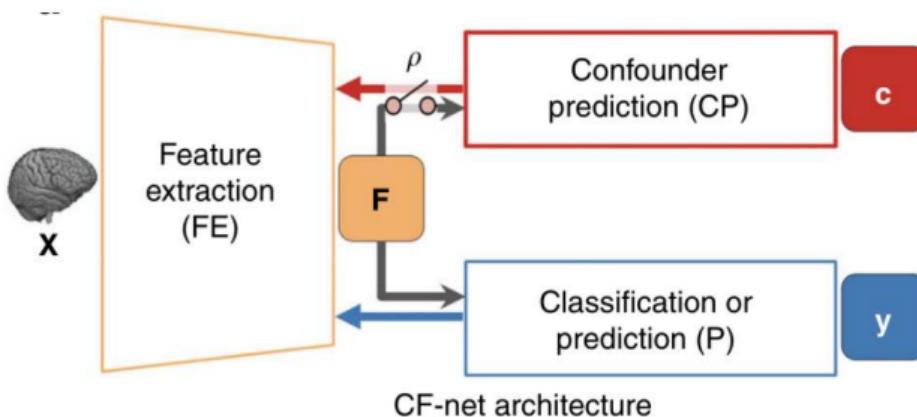
During test:

- Change an image from class A s.t. its brain shape is in class B
- Check if the performance decrease

Solution:
brain shape
normalization

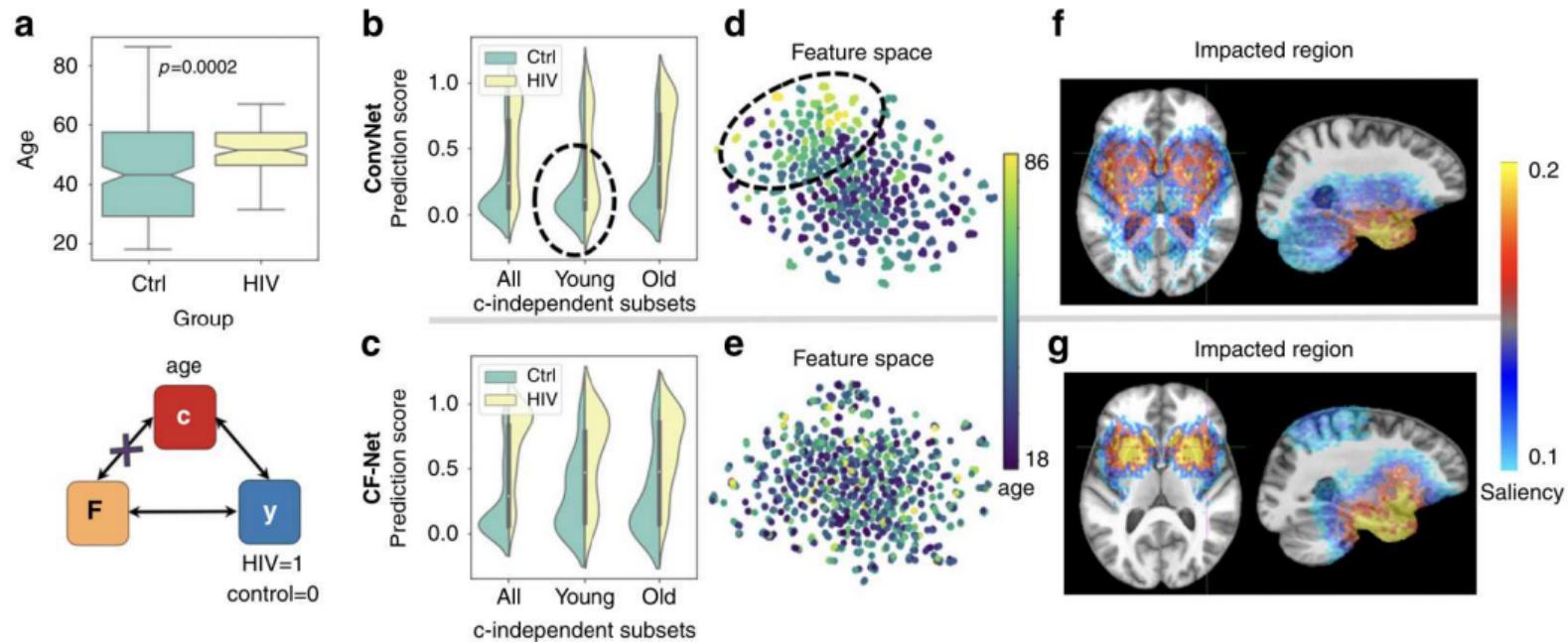


Confounder Free Training



- update FE+P
estimate y
- update CP **on control only**
estimate c
avoid confounding of y on CP
- adversarial update FE:
to fool CP
 \Rightarrow F cannot predict c

Confounder Free Training: results



Conclusion

Thanks for your attention !!

Path Methods: Integrated Gradient

$$A(I, B) = (I - B) \int_0^1 \frac{\partial F}{\partial I} (B + \alpha(I - B)) d\alpha$$

Completeness	$\sum_x A(I)_x = F(I) - F(B)$
Sensitivity	F do not depend on pixel $x \implies A(I)(x) = 0$
Linearity	$A(aF_1 + bF_2) = aA(F_1) + bA(F_2)$
Implementation	$\forall I, F_1(I) = F_2(I) \rightarrow A_1 = A_2$
Invariance	

Th. Only path methods always satisfy these 4 properties.

Global vs Local Attributions

$$C = 1.05x_1 + 10x_2$$

$$x_1 = 1.000.000 \text{ and } x_2 = 10.000$$

Local:

$$R1 = \frac{\partial C}{\partial x_1} = 1.05$$

$$R2 = \frac{\partial C}{\partial x_2} = 10$$

How do you behave locally ?

Global:

$$R1 = x_1 \frac{\partial C}{\partial x_1} = 1.050.000$$

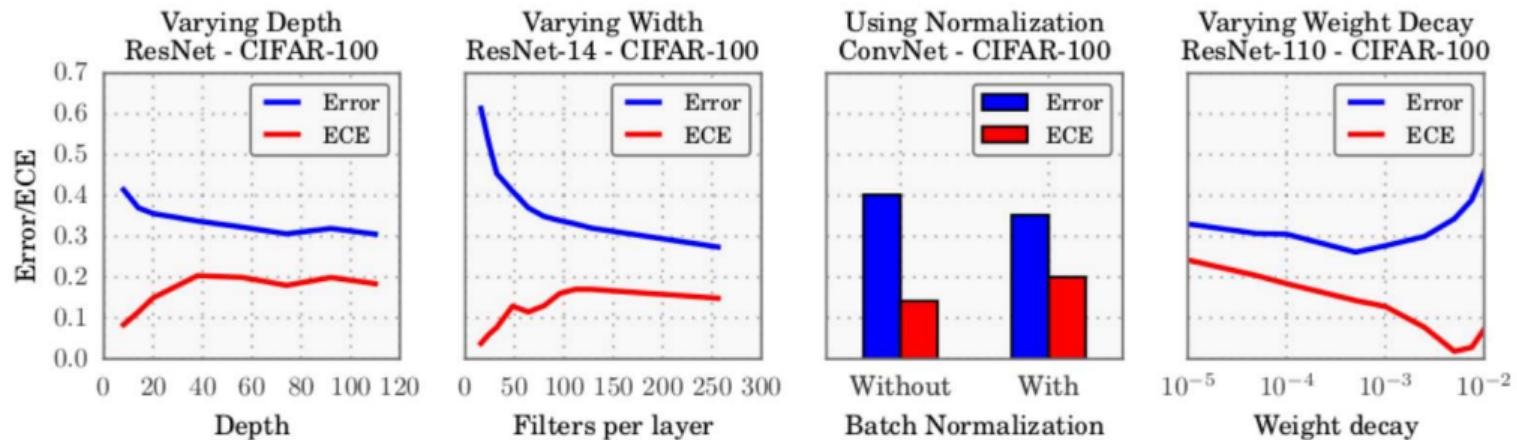
$$R2 = x_2 \frac{\partial C}{\partial x_2} = 10.000$$

How initial investments contributed ?

Retropagation rules: LRP, DeepLIFT, ...

Method	Attribution $R_i^c(x)$	Example of attributions on MNIST			
		ReLU	Tanh	Sigmoid	Softplus
Gradient * Input	$x_i \cdot \frac{\partial S_c(x)}{\partial x_i}$				
Integrated Gradient	$(x_i - \bar{x}_i) \cdot \int_{\alpha=0}^1 \frac{\partial S_c(\tilde{x})}{\partial (\tilde{x}_i)} \Big _{\tilde{x}=\bar{x}+\alpha(x-\bar{x})} d\alpha$				
<u>ϵ-LRP</u>	$x_i \cdot \frac{\partial^g S_c(x)}{\partial x_i}, \quad g = \frac{f(z)}{z}$				
<u>DeepLIFT</u>	$(x_i - \bar{x}_i) \cdot \frac{\partial^g S_c(x)}{\partial x_i}, \quad g = \frac{f(z) - f(\bar{z})}{z - \bar{z}}$				
Occlusion-1	$S_c(x) - S_c(x_{[x_i=0]})$				

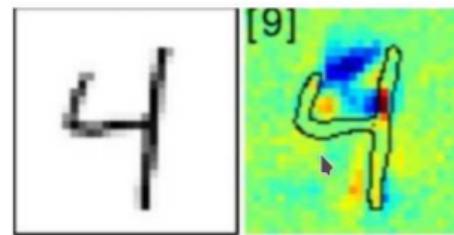
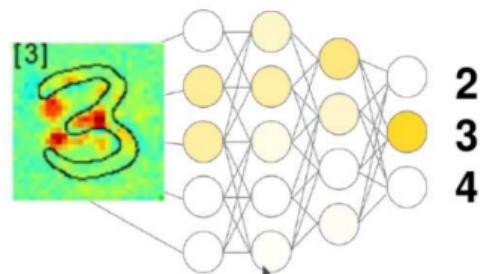
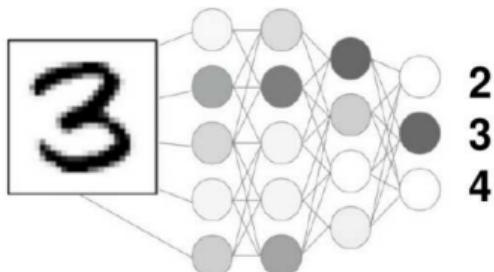
Network probability output



Guo et. al. 2017, On Calibration of Modern Neural Networks

About Interpreting Deep Learning...

Layer-wise Relevance Propagation (LRP)



Investigate bias on training dataset

- ▶ Age Predictions from face images
- ▶ with pretraining on ImageNet, laughing speaks against 60+
- ▶ model learned that old people do not laugh



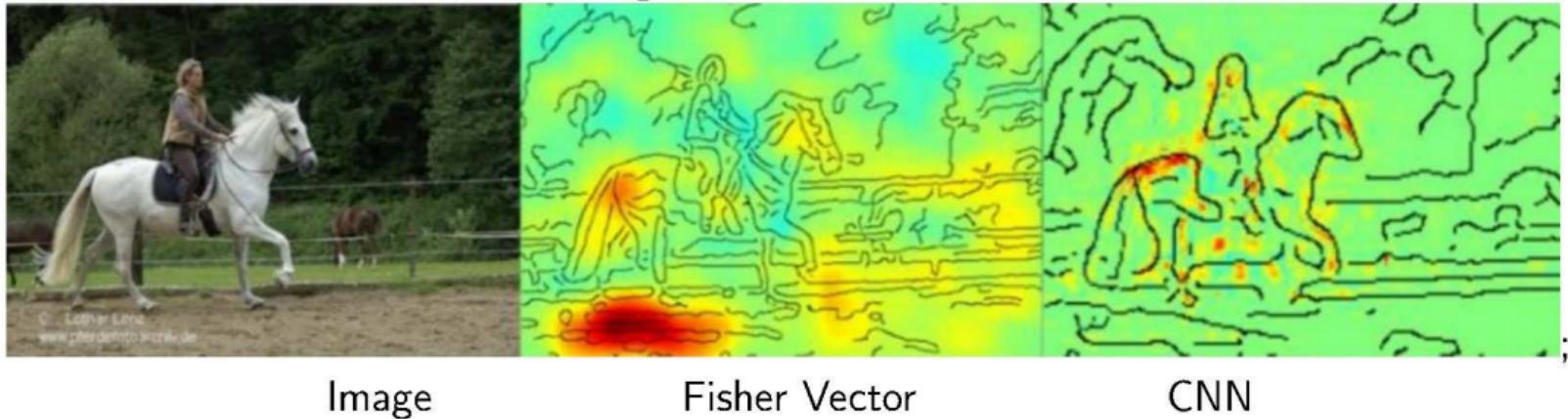
▶ 25-32 years old relevance image

▶ 60+ years old relevance image
(pre training on image net)

▶ 60+ years old relevance image
(pre training on IMDB-WIKI)

Investigate bias on training dataset

horse images in PASCAL VOC 2007



Accuracy:	FV	80.45%
	CNN	81.60%

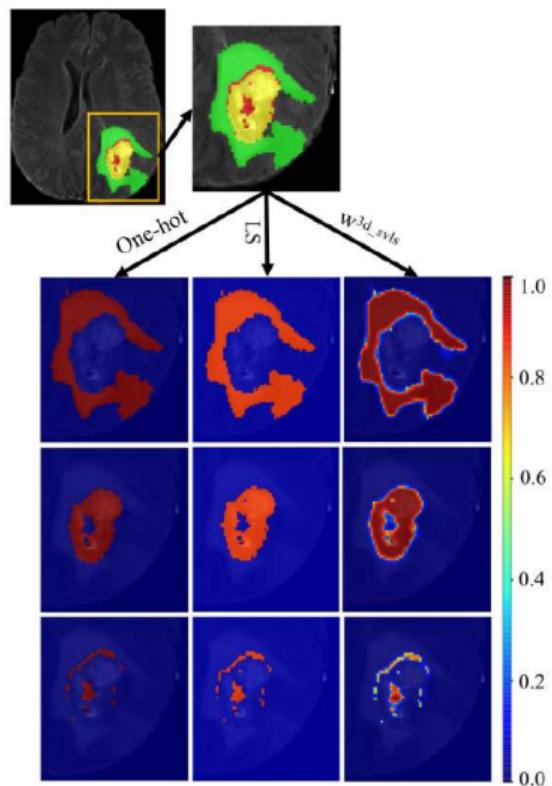
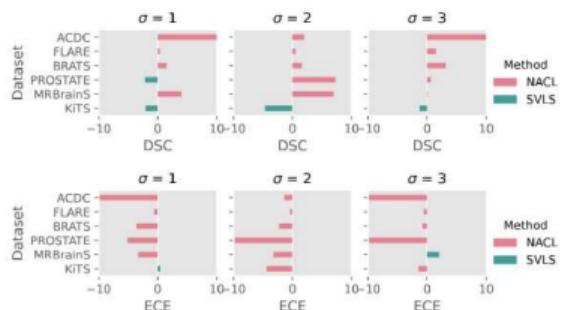
Spatially Varying Label Smoothing

SVLS

- gaussian smooth GT label maps
- use it in CE

NACL

$$\min CE + P(\|l - \tau\|)$$



Islam & Glocker, IPML 2021. Spatially varying label smoothing: Capturing uncertainty from expert annotations
 Murugesan ... Dolz , MEDIA 2024, Neighbor-Aware Calibration of Segmentation Networks with Penalty-Based Constraints

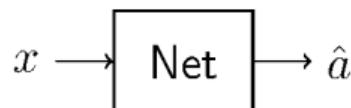
Regression



$$\text{MSE} = \|\hat{a} - a\|^2$$

Nix & Weigend, ICNN 1994. Estimating the mean and variance of the target probability distribution

Regression



$$\text{MSE} = \|\hat{a} - a\|^2$$

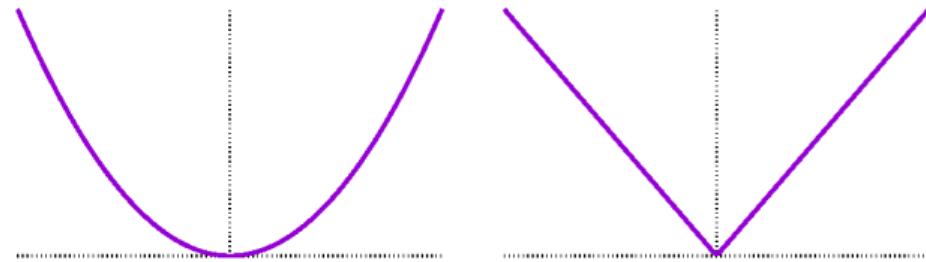
$$A(x) \sim \mathcal{N}(\mu(x), \sigma(x))$$

$$\text{likelihood: } \frac{1}{\sqrt{2\pi}\sigma} e^{\frac{-\|\mu-a\|^2}{2\sigma^2}}$$

homoscedastic ($\sigma = \text{cst}$): $\rightarrow \text{NLL} : \log(\sigma) + \frac{\|\mu-a\|^2}{2\sigma^2}$
 $\Rightarrow \text{MSE} = \text{NLL}$

Quantile Regression

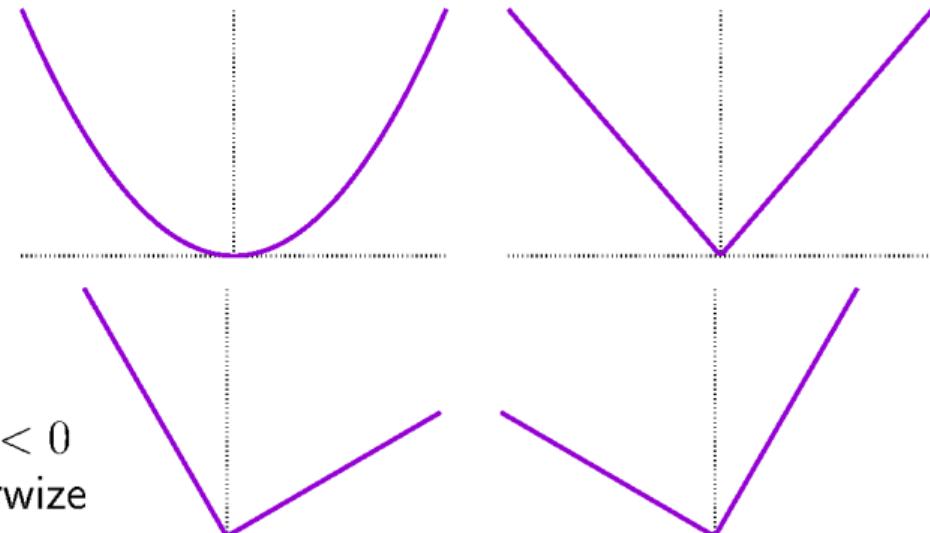
$$\min_x \sum_i r(x - x_i)$$



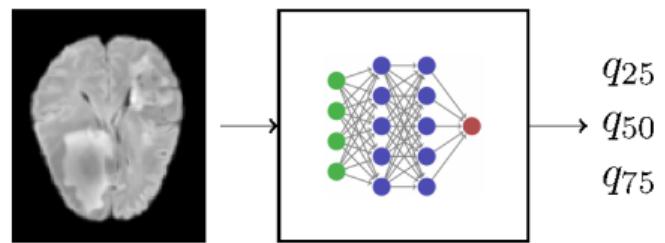
Quantile Regression

$$\min_x \sum_i r(x - x_i)$$

$$r_\tau(x) = \begin{cases} -(1-\tau)x & \text{if } x < 0 \\ x\tau & \text{otherwise} \end{cases}$$



Quantile Regression



$$\text{Loss: } r_{25}(q_{25} - y) + r_{50}(q_{50} - y) + r_{75}(q_{75} - y)$$

More on conformal prediction

Adaptive prediction set:

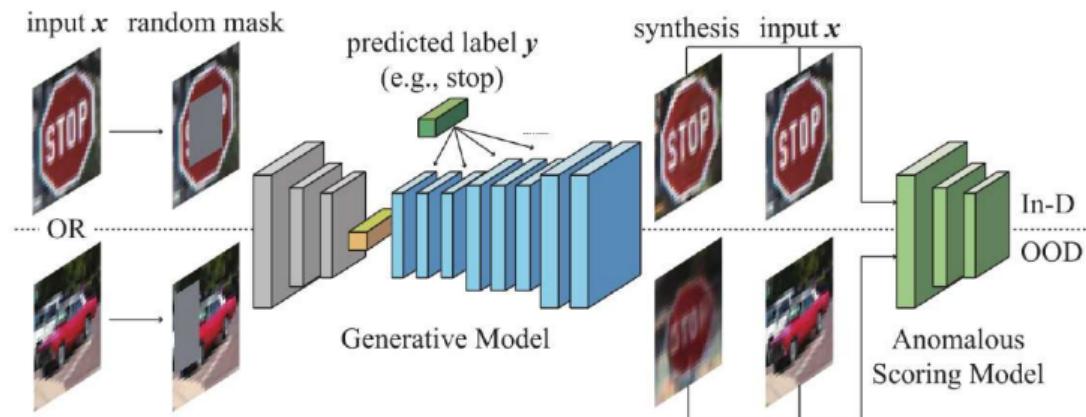
$$s(x, y) = 1 - \sum_{f_z(x) \geq f_y(x)} f_z(x)$$

Regression:

- quantile regression: $\rightarrow f_q(x), f_{1-q}(x)$ (with $q = \alpha/2$)
- score: $s(x, y) = \max\{f_q(x) - y, y - f_{1-q}(x)\}$
- ...

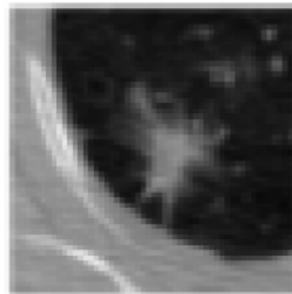
Angelopoulos & Bates, "A gentle introduction to conformal prediction and distribution-free uncertainty quantification." arXiv 2021

Reconstruction-based OOD



Denouden et al, arxiv 2018, Improving Reconstruction Autoencoder Out-of-distribution Detection with Mahalanobis Distance

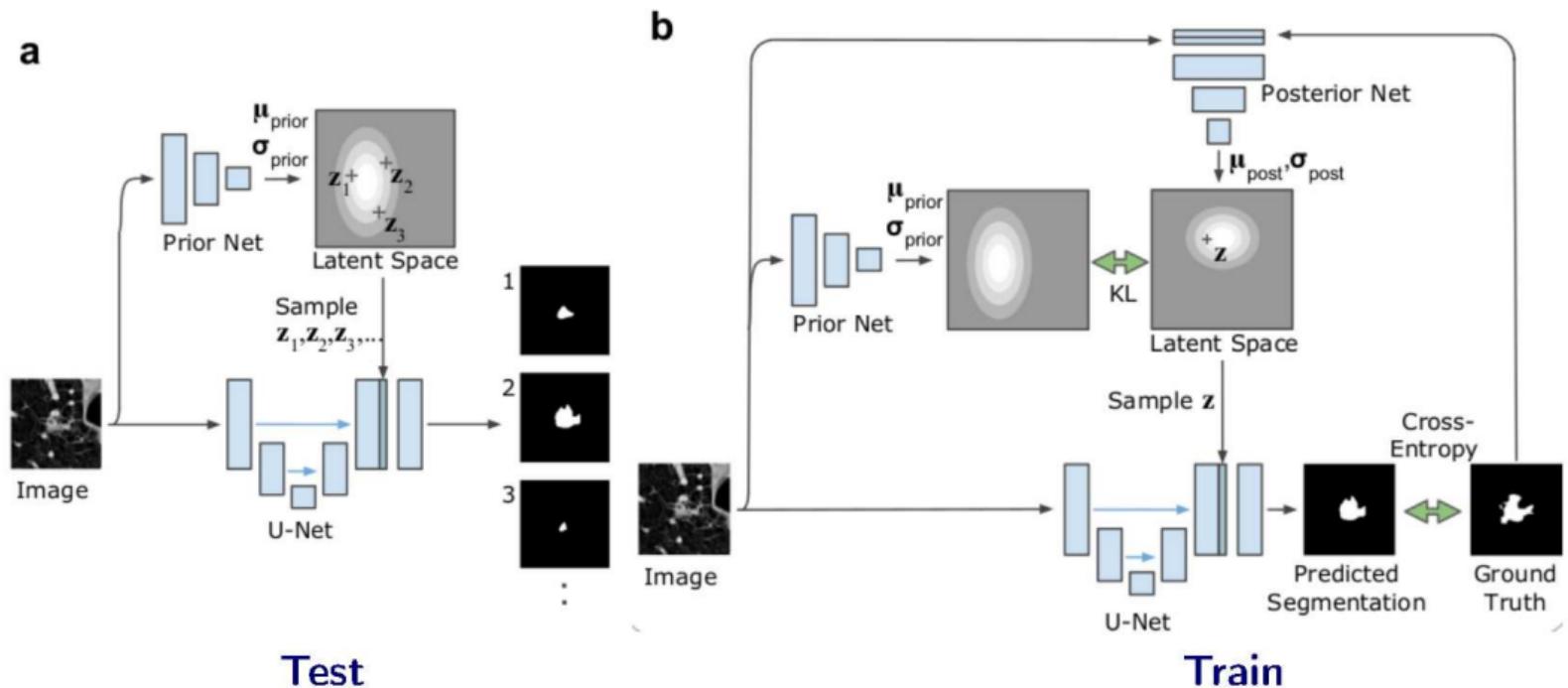
Uncertainty with respect to expert annotations



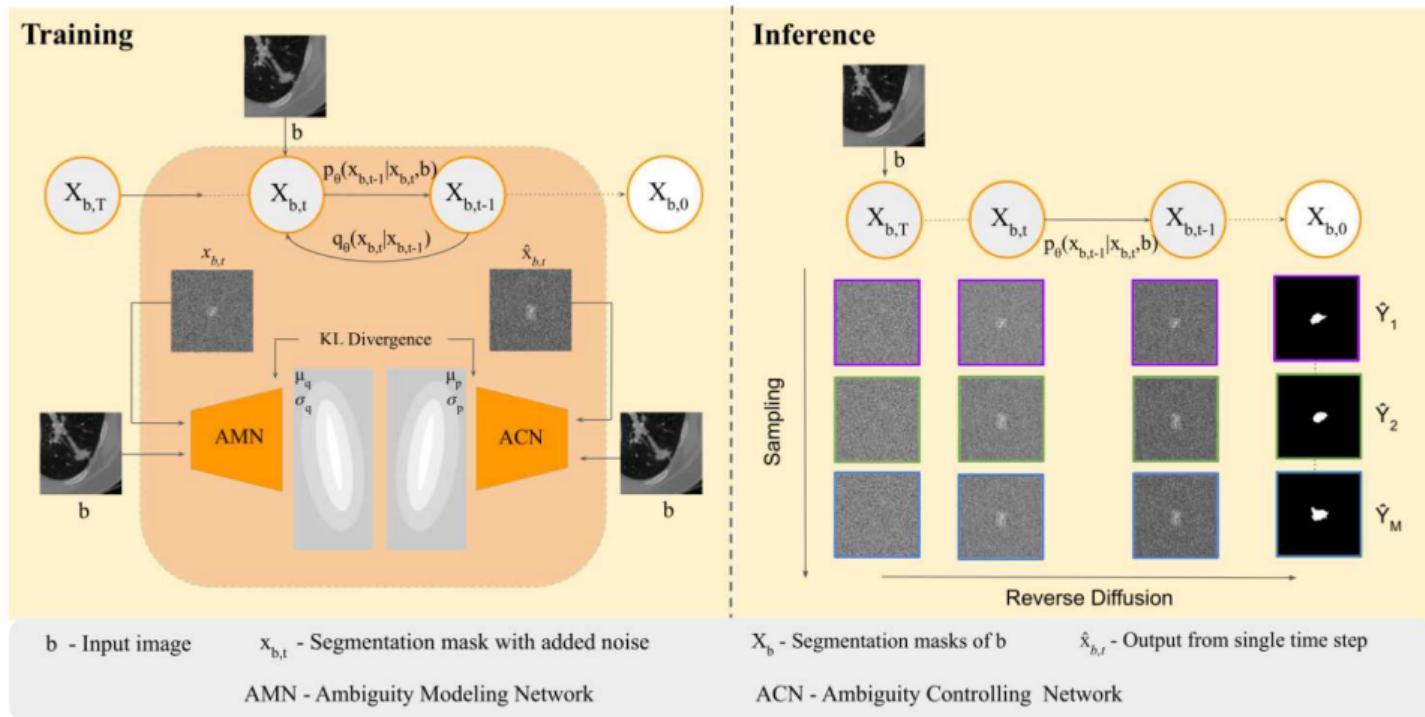
Lung CT images with segmentation made by 4 operators

Kohl et. al. NIPS 2018, A probabilistic u-net for segmentation of ambiguous images

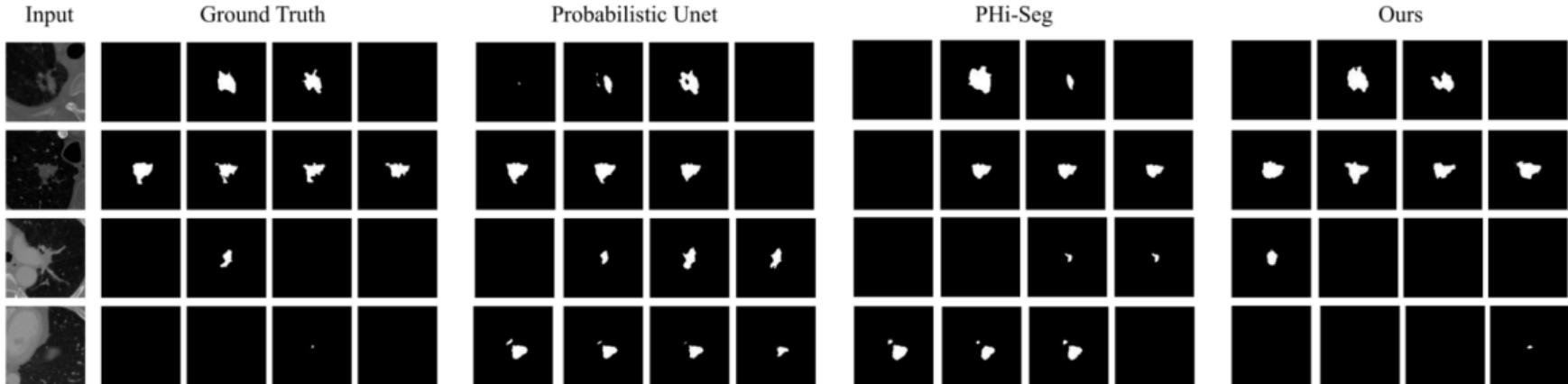
Uncertainty with respect to expert annotations



Generating segmentation distribution with diffusion models



Generating segmentation distribution with diffusion models



Method	GED (\downarrow)	CI (\uparrow)	D_{max} (\uparrow)
DDPM-det-Seg [59]	1.081	0.616	0.548
DDPM-Prob-Seg	0.417	0.683	0.689
CIMD (Ours)	0.321	0.759	0.915