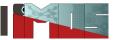




 École polytechnique fédérale



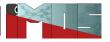
Which questions to follow to select an ML model



- 1. What is the nature of the problem?
- Is it a classification, regression, clustering, or recommendation problem?
- Are you dealing with supervised, unsupervised, or reinforcement learning?
- 2. What is the size and quality of the dataset?
- How many samples and features are available?
- Is the data clean, or does it contain a lot of noise and missing values?
- Is the data balanced, or are there significant class imbalances?
- 3. Is the training dataset representative of the expected operating / application conditions
- Is there a high diversity of operating conditions expected?
- 4. What are the characteristics of the data?
- Is the data structured or unstructured (e.g., text, images, audio)?
- Are the features numerical, categorical, or a mix of both?
- Do the features have a temporal or spatial component?



Which questions to follow to select an ML model



5. What are the performance requirements?

- Do you prioritize accuracy, interpretability, or computational efficiency?
- Is the problem domain one where explainability is crucial (e.g., healthcare, finance)?
- What are the acceptable trade-offs between bias and variance?

6. What is the computational complexity and resources available?

- How much time and computational power do you have for training?
- Do you have access to specialized hardware (e.g., GPUs)?

What are the specific goals and constraints?

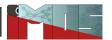
- What is the end goal of applying the ML model (e.g., prediction, anomaly detection, optimization)?
- Are there specific business or application constraints to consider (e.g., real-time) processing, deployment environment)?

What is the level of domain knowledge available?

- How well do you understand the domain and the problem?
- Is there domain expertise to help with feature engineering and model interpretation?



Which questions to follow to select an ML model



9. What is the expected model lifecycle?

- How often will the model need to be updated or retrained?
- Is the model expected to handle concept drift (changes in the underlying) data distribution over time)?

10. What is the potential for model interpretability?

- Do stakeholders need to understand how the model makes decisions?
- Is there a regulatory requirement for transparency?

11. What are the available benchmarks or baselines?

- Are there existing solutions or benchmarks that can provide a performance reference?
- How do different algorithms perform on similar problems in literature or industry standards?

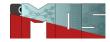


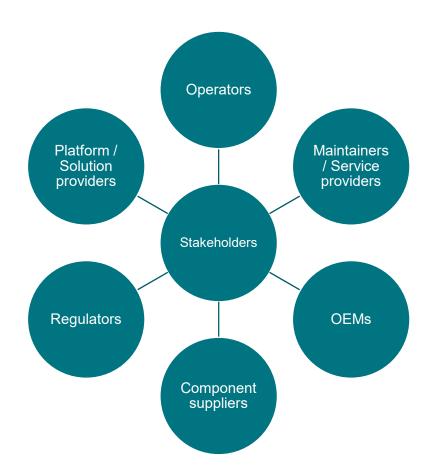


Data sharing/ Federated learning

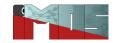


Different stakeholders in PHM





The importance of data for ML



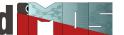
The biggest obstacle to using advanced data analysis isn't skill base or technology; it's plain old access to the data.

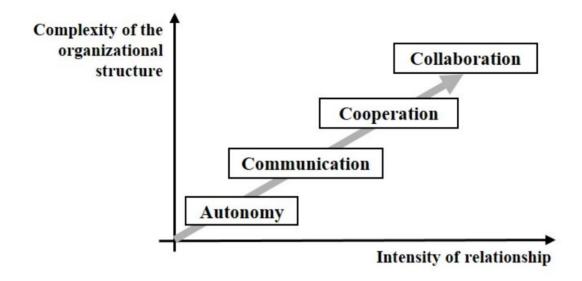
-Edd Wilder-James, Harvard Business Review



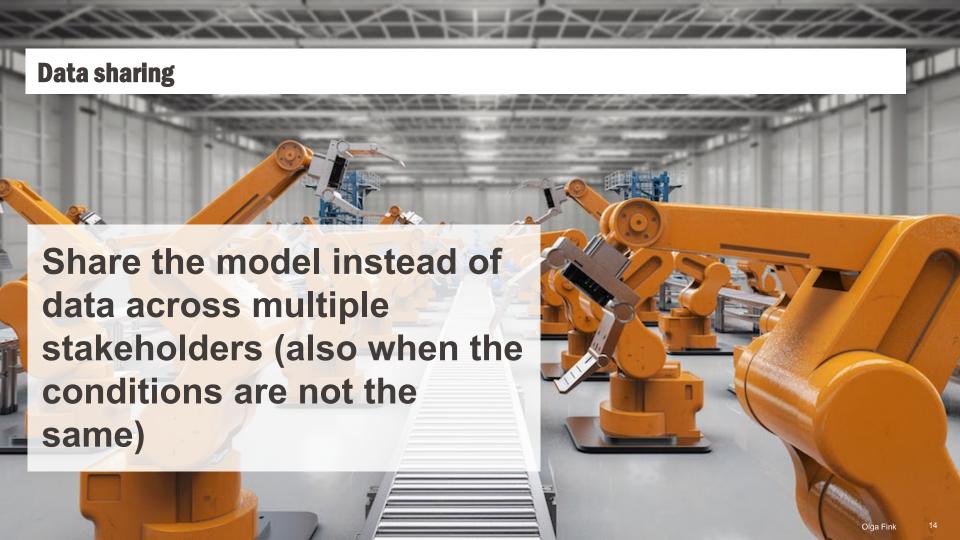


New forms of sharing and collaboration required



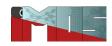








A new paradigm – Federated Learning



a synchronous update scheme that proceeds in rounds of communication

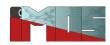
McMahan, H. Brendan, Eider Moore, Daniel Ramage, and Seth Hampson. "Communication-efficient learning of deep networks from decentralized data." *AISTATS*, 2017.

04 11 24

Source: Min Du, 2019



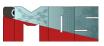
Federated learning



 Federated learning is a machine learning setting where multiple entities (clients) collaborate in solving a machine learning problem, under the coordination of a central server or service provider. Each client's raw data is stored locally and not exchanged or transferred; instead, focused updates intended for immediate aggregation are used to achieve the learning objective.



Federated learning characteristics (1/2)



1. Data Privacy and Security

- **Data Locality**: Raw data remains on local devices (e.g., smartphones, edge devices, or distributed servers) instead of being sent to a central server, enhancing privacy and security.
- **Secure Aggregation**: Only model updates (e.g., gradients) are sent to a central server, and these updates can be further encrypted or anonymized, reducing the risk of exposing sensitive information.
- **Privacy-Preserving Techniques**: Techniques like differential privacy or homomorphic encryption are often integrated to ensure that individual data points cannot be inferred from the shared updates.

2. Decentralized Training

- Local Training: Each device trains a model locally on its own data, which means training occurs independently
 across a distributed network.
- Global Model Aggregation: After local training, model parameters are sent to a central server that aggregates these updates to form a new global model, which is then sent back to each local device to continue training.

3. Handling Data Heterogeneity

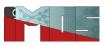
- **Non-IID Data**: In federated learning, data on different devices is often non-independent and identically distributed (non-IID), meaning data distributions may vary significantly between devices.
- Variable Data Quality and Quantity: Some devices may have abundant data, while others have minimal or noisy data, requiring federated learning algorithms to account for varying data quality and volume.

4. Communication Efficiency

- Minimizing Data Transfer: To limit network usage, federated learning focuses on minimizing the frequency and size
 of data transferred between the server and devices.
- **Compression Techniques**: Methods like quantization, sparsification, and pruning are used to reduce the size of model updates, which is essential for devices with limited bandwidth.



Federated learning characteristics (2/2)



5. Device Heterogeneity and Scalability

- Diverse Device Capabilities: Federated learning systems are designed to work across a wide variety of devices, from smartphones to IoT devices, each with different computational capabilities, network stability, and battery life.
- Scalability: Federated learning is scalable, as it can handle thousands or even millions of devices simultaneously, with participants joining or leaving the training process dynamically.

6. Model Personalization

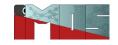
• **Personalized Models**: Federated learning allows each device to maintain its own personalized version of the model, especially if local data significantly deviates from the global model's training data. This feature can improve the model's performance on individual devices.

7. Fault Tolerance

• Robust to Device Failures: Federated learning systems can continue training even if some devices drop out or are temporarily unavailable, making them resilient to intermittent connectivity and device churn.



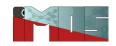
Federated Learning vs. Peer-to-peer learning

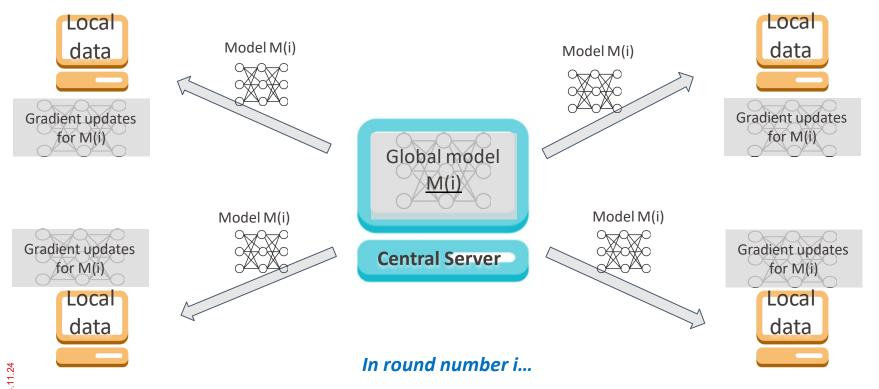


	Federated learning	Fully decentralized (peer-to-peer) learning
Orchestration	A central orchestration server or service organizes the training, but never sees raw data.	No centralized orchestration.
Wide-area communication	Typically a hub-and-spoke topology, with the hub representing a coordinating service provider (typically without data) and the spokes connecting to clients.	Peer-to-peer topology, with a possibly dynamic connectivity graph.



Federated learning – overview

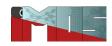


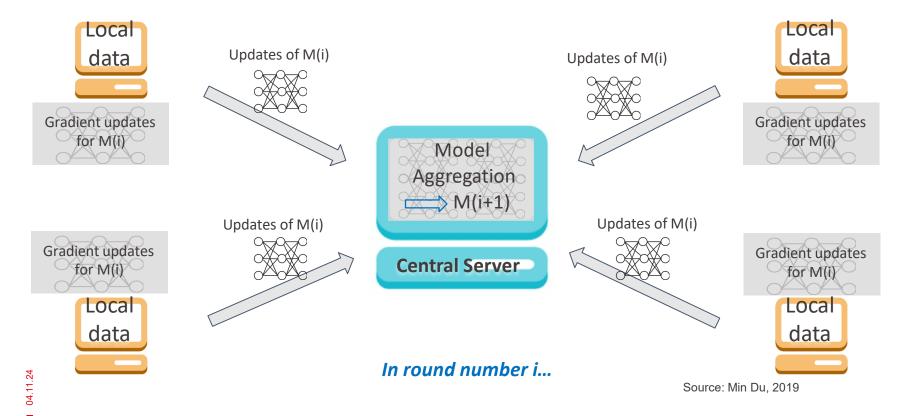


Source: Min Du, 2019



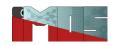
Federated learning – overview

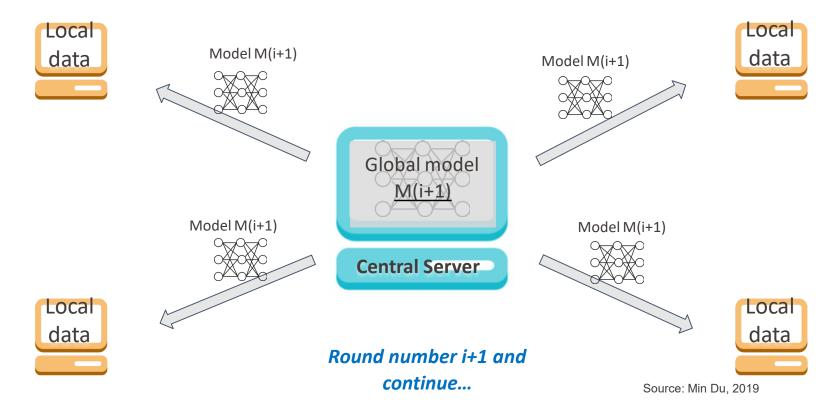






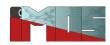
Federated learning – overview







Federated learning – detail



- Recall in traditional deep learning model training
 - For a training dataset containing n samples (x_i, y_i) , $1 \le i \le n$, the training objective is:

$$\min_{w \in \mathbb{R}^d} f(w) \qquad \text{where} \quad f(w) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n f_i(w)$$

$$f(w) = l(x_i, y_i, w)$$
 is the loss of the prediction on example (x_i, y_i)

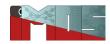
Deep learning optimization relies on SGD and its variants, through <u>mini-batches</u>

$$w_{t+1} \leftarrow w_t - \eta \nabla f(w_t; x_k, y_k)$$

Source: Min Du, 2019



Federated learning – detail

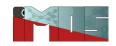


- In federated learning
 - Suppose n training samples are distributed to K clients, where $P_{\rm k}$ is the set of indices of data points on client k, and $n_{\rm k}=|P_{\rm k}|$
 - For training objective: $\min_{w \in \mathbb{R}^d} f(w)$

$$f(w) = \sum_{k=1}^{K} \frac{n_k}{n} F_k(w)$$
 where $F_k(w) \stackrel{\text{def}}{=} \frac{1}{n_k} \sum_{i \in P_k} f_i(w)$



A baseline - FederatedSGD (FedSGD)

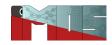


- A randomly selected client that has $n_{\rm k}$ training data samples in federated learning pprox A randomly selected sample in traditional deep learning
- Federated SGD (FedSGD): a single step of gradient descent is done per round
- Recall in federated learning, a C-fraction of clients are selected at each round.
 - C=1: full-batch (non-stochastic) gradient descent
 - <u>C<1</u>: stochastic gradient descent (SGD)





FederatedSGD (FedSGD) / FederatedAveraging (FedAvg)



Learning rate: η ; total #samples: n; total #clients: K; #samples on a client k: n_k ; clients fraction C=1

- In a round t:
 - The central server broadcasts current model w_I to each client; each client k computes gradient: $g_k = \nabla F_k(w_t)$, on its local data.
 - Approach 1: Each client k submits g_k ; the central server aggregates the gradients to generate a new model:

•
$$w_{t+1} \leftarrow w_t - \eta \nabla f(w_t) = w_t - \eta \sum_{k=1}^K \frac{n_k}{n} g_k$$
.

Recall
$$f(w) = \sum_{k=1}^{K} \frac{n_k}{n} F_k(w)$$

- Approach 2: Each client k computes: $w_{t+1}^k \leftarrow w_t \eta g_k$; the central server performs aggregation:
 - $W_{t+1} \leftarrow \sum_{k=1}^{K} \frac{n_k}{n} W_{t+1}^k$

For multiple times ⇒ FederatedAveraging (FedAvg)

Source: Min Du, 2019



FedSGD vs FedAvg



1. FederatedSGD (FedSGD)

- Process: In FedSGD, each device performs a single step of stochastic gradient descent (SGD) on its local data and then sends the computed gradients (i.e., updates) to the central server. The server then aggregates these gradients from all devices and applies them to the global model.
- Communication: FedSGD involves frequent communication between the devices and the server, as each gradient update requires synchronization with the server.
- Advantages: FedSGD ensures that each update to the global model reflects the most recent local gradient information from each device. This can sometimes lead to faster convergence if the communication overhead is manageable.
- **Drawbacks**: The high frequency of communication in FedSGD is a major drawback, especially in federated learning scenarios where devices might have limited bandwidth. Frequent communication also consumes more energy on devices, which is a concern for battery-powered devices like mobile phones.

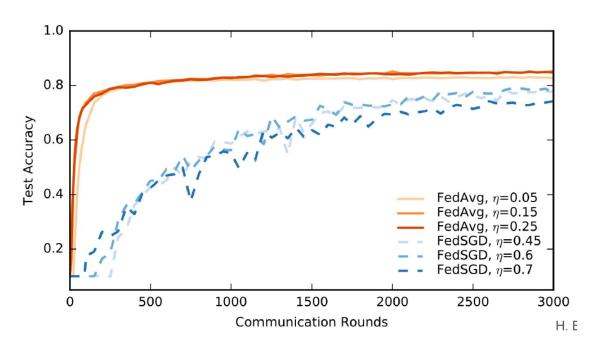
2. FederatedAveraging (FedAvg)

- Process: In FedAvg, each device performs multiple local updates (i.e., several steps of SGD) on its local data before sending the updated model parameters to the central server. The server then averages the model parameters from all participating devices to update the global model.
- **Communication**: By performing multiple local updates before communicating, FedAvg significantly reduces the frequency of communication between devices and the server.
- Advantages: FedAvg is much more communication-efficient than FedSGD since each device communicates with the server only after several local updates. This makes it more suitable for federated learning, where communication costs are high. Additionally, local updates often lead to faster convergence in practical applications, as they allow each device to adapt the model more effectively to its own data before syncing.
- **Drawbacks**: FedAvg might converge more slowly or be less stable in scenarios where data on each device is highly heterogeneous (non-IID). In such cases, performing multiple local updates can lead to models that diverge from each other, making aggregation less effective.



CIFAR-10 Convolutional model





Updates to reach 82% SGD 31,000 FedSGD 6,600 FedAvg 630

decrease in communication (updates) vs SGD

(IID and balanced data)

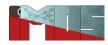




Example



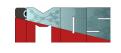
Federated learning: limitations



- Data heterogeneity in industrial settings:
- 1. Domain shift: Clients do not share a similar data distribution.
- Label heterogeneity: Clients' datasets have a different number and types of faults.



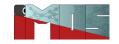
Federated learning adapted to the specificities of the unit

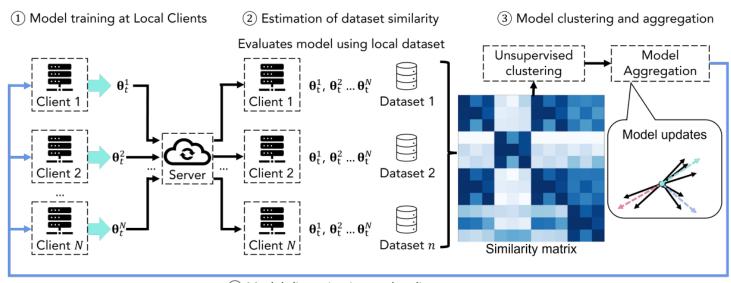


- Distance awareness: the model's ability to quantify the distance of a testing example from the training data.
- For personalized federated learning, during model aggregation → a local client would assign higher weights to the model that was trained using similar training data.
- Spectral-normalized Neural Gaussian Process (SNGP) used to quantify the prediction uncertainty



Federated Learning with Uncertainty-Based Client Clustering for Fleet-Wide Fault Diagnosis

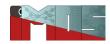


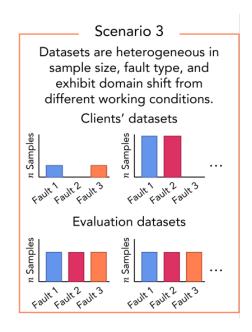


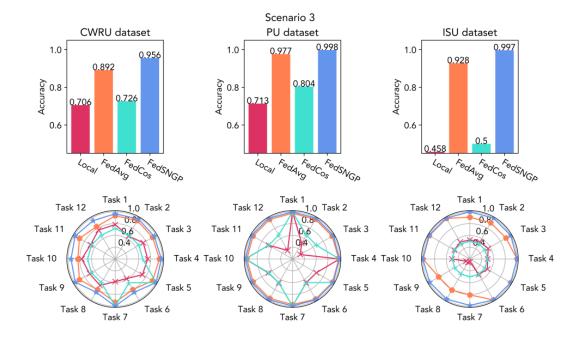
4 Model dissemination to the clients.



Results on FedSNGP







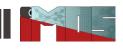




Explainability / Interpretability



Most of the time, the algorihtms work really well but sometimes...

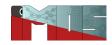




A refrigerator filled with lots of food and drinks



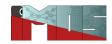
Tesla Autopilot Misidentified On-Road Horse-Drawn Carriage

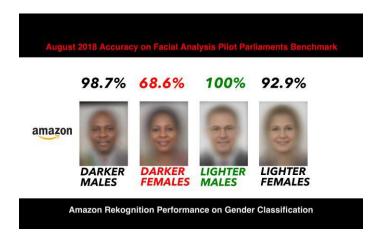






Bias in algorithms





https://medium.com/@Joy.Buolamwini/response-racial-and-gender-bias-in-amazon-rekognition-commercial-ai-system-for-analyzing-faces-a289222eeced

Machine Learning can amplify bias.



- . Data set: 67% of people cooking are women
- Algorithm predicts: 84% of people cooking are women

https://www.infoq.com/presentations/unconscious-bias-machine-learning/

37

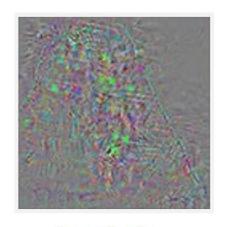


Adversarial Examples





Original image
Temple (97%)



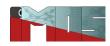
Perturbations



Adversarial example
Ostrich (98%)



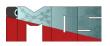
Current limitations of the ML algorithms



- Interpretability (Explainability, Transparency, Understanding, Trust)
- Physical consistency
- Complex and uncertain data
- Limited labels
- Bias
- (Computational demand)



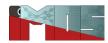
What is interpretability?



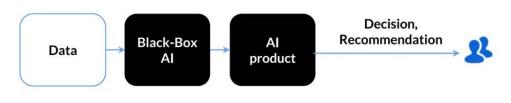
 Ability to explain or to present a model in understandable terms to humans (Doshi-Velez 2017)



Why explainable Al?



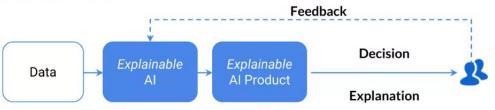
Black Box Al



Confusion with Today's Al Black Box

- Why did you do that?
- Why did you not do that?
- When do you succeed or fail?
- How do I correct an error?

Explainable Al

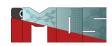


Clear & Transparent Predictions

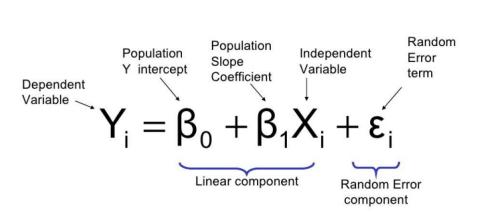
- I understand why
- I understand why not
- I know why you succeed or fail
- I understand, so I trust you

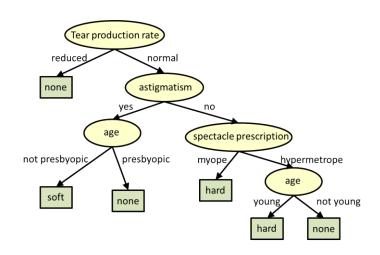


Simple explainability



 In pre-deep learning models, some models are considered "interpretable"





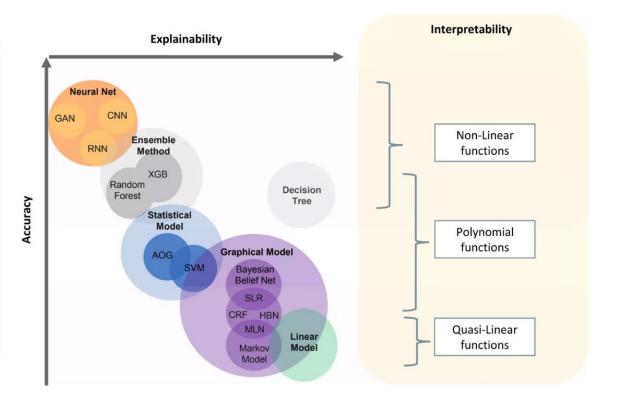


Accuracy vs. Explainability



Learning

- · Challenges:
 - Supervised
 - · Unsupervised learning
- · Approach:
 - Representation Learning
 - · Stochastic selection
- · Output:
 - Correlation
 - · No causation





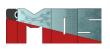
Different ways to achieve interpretability



- Build interpretability into the model (e.g. by fusing physical models and machine learning or learning the underlying physics explicitly)
- Post-hoc approach to interpretability → trying to explain given models and their output



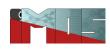
Some properties of Interpretations



- Faithfulness how to provide explanations that accurately represent the true reasoning behind the model's final decision.
- Plausibility Is the explanation correct or something we can believe is true, given our current knowledge of the problem?
- Understandable Can I put it in terms that end user without in-depth knowledge of the system can understand?
- Stability Do similar instances have similar interpretations?



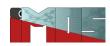
Evaluating Interpretability



- Application level evaluation Put the model in practice and have the end users interact with explanations to see if they are useful.
- Human evaluation Set up a Mechanical Turk task and ask nonexperts to judge the explanations
- Functional evaluation Design metrics that directly test properties of your explanation.



Global vs Local

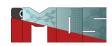


- Do we explain individual prediction?
- Example :
- Heatmaps
- Rationales

- Do we explain entire model?
- Example :
- Linear Regression
- Decision Trees



Inherent vs Post-hoc



Is the explainability built into the model?

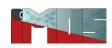
- Examples:
- Linear Regression
- Decision Trees

Is the model black-box and we use external method to try to understand it?

- Examples:
- Heatmaps (Some forms)



Model based vs Model Agnostic



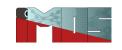
Can it explain only few classes of models?

- Examples:
- Decision Trees
- Attention
- Gradients (Differentiable Models only)

- Can it explain any model ?
- Examples:
- LIME Locally Interpretable Model Agnostic Explanations
- SHAP Shapley Values



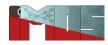
Different approaches to explain the behavior of **ML** approaches post-hoc



- Explaining with Surrogates
- Explaining with local perturbations
- Propagation-Based Approaches (Leveraging Structure)
- Meta-explanations
- Attribution-based methods
- Counterfactual explanation
- Interaction explanations
- Attention Mechanisms



Explainable AI (1/2)



Feature Importance and Attribution Methods

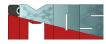
- These methods identify the most influential features or components in making a prediction.
- SHAP (SHapley Additive exPlanations): Based on cooperative game theory, SHAP assigns each feature a "Shapley value" that indicates its contribution to a particular prediction. SHAP values provide consistent and fair attribution by treating each feature as a "player" in a game.
- LIME (Local Interpretable Model-agnostic Explanations): LIME approximates complex models by locally fitting an interpretable model (like linear regression) around a specific instance. It perturbs input features to see how predictions change, helping identify which features influence the prediction.
- Integrated Gradients: A technique for neural networks that calculates feature importance by integrating the gradients of a model's output with respect to its inputs, starting from a baseline (e.g., zero). Integrated gradients quantify each input's contribution to a model's prediction.
- **Permutation Feature Importance**: Measures feature importance by shuffling each feature in turn and observing the impact on model accuracy. If a feature's importance is high, shuffling its values will significantly degrade model performance.

Visualization Techniques

- Visualization is key to making model insights understandable for humans, especially in image and text processing.
- Saliency Maps: Often used in image processing, saliency maps highlight the parts of an image that have the most influence on a model's prediction. This can help in understanding what parts of an image a convolutional neural network (CNN) is focusing on.
- Partial Dependence Plots (PDP): Show the relationship between one or two features and the predicted outcome, averaging out the effect of other features. PDPs help understand how the model's predictions change as the value of a feature changes.
- Layer-wise Relevance Propagation (LRP): Assigns relevance scores to each neuron in the network, which can be aggregated to show which pixels or words are most influential in a decision. This technique is especially useful for deep neural networks.
- t-SNE and UMAP: Dimensionality reduction techniques that help visualize high-dimensional data in two or three dimensions, making it easier to observe clusters or patterns in the data that the model may have learned.



Explainable AI (2/2)



Surrogate Models

Surrogate Models: These are simpler, interpretable models (like linear or decision tree models) trained to approximate a complex model's
predictions. Surrogate models can provide a global understanding of a complex model's behavior by approximating its decision boundaries.

Counterfactual Explanations

Counterfactual explanations help answer "what if" questions by indicating how a model's prediction would change if the input features were
altered. For example, in loan applications, a counterfactual explanation might state, "If the applicant's income were \$5,000 higher, the loan
would be approved." This type of explanation is valuable for understanding decision boundaries and fairness.

Causal Inference and Causal Models

Causal inference goes beyond correlation by trying to understand causal relationships between variables. Causal models can help identify
which factors genuinely influence predictions rather than merely correlate with the outcome. These methods are particularly important in
applications where knowing causality is crucial, such as medicine or social sciences.

Rule-Based Explanations and Decision Rules

Some models, such as decision trees and rule-based models, are naturally interpretable. These models use a series of if-then rules that are
easy for humans to follow. For more complex models, rule extraction techniques can derive decision rules that approximate the model's decision
process, making it more interpretable.

Prototype and Example-Based Explanations

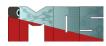
- **Prototypes**: These are representative examples from the training set that encapsulate common patterns or features. Showing a prototype can help users understand typical cases within a class.
- Influence Functions: These methods determine which training examples are most responsible for a specific prediction. They help identify cases in the training data that the model relied on heavily, which can be especially useful for diagnosing model biases.

Bayesian and Probabilistic Approaches

Probabilistic models, such as Bayesian networks, can inherently quantify uncertainty in their predictions. By providing a probabilistic
interpretation of model decisions, these methods allow users to understand the confidence level associated with each prediction, helping to
address questions of model reliability.

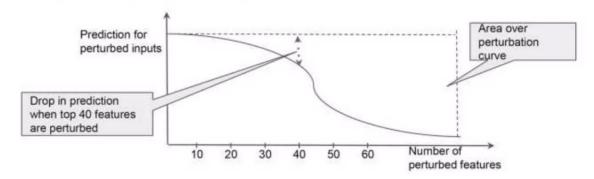


Perturbation-based approaches



Perturb top-k features by attribution and observe change in prediction

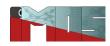
- Higher the change, better the method
- Perturbation may amount to replacing the feature with a random value
- Samek et al. formalize this using a metric: Area over perturbation curve
 - Plot the prediction for input with top-k features perturbed as a function of k
 - Take the area over this curve



Source: KDD 2019 Tutorial



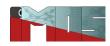
Attribution-based methods



- Attribute a model's prediction on an input to features of the input
- Attribiution methods
 - Ablations (drop each feature and attribute the change in prediction to that feature)
 - Gradient-based methods (attribution to a feautre is feature value times) gradient)
 - Score backpropation based methods
 - Guided BackProp: Only consider ReLUs that are on (linear regime), and which contribute positively
 - LRP: Use first-order Taylor decomposition to linearize activation function
 - **DeepLift**: Distribute activation difference relative a reference point in proportion to edge weights



Saliency Based Methods

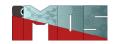


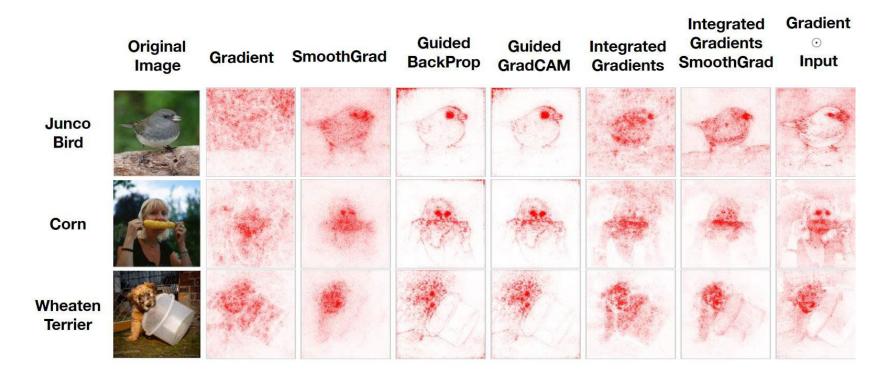
- Heatmap based visualization
- Need differentiable model in most cases
- Normally involve gradient





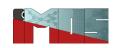
Saliency Based Methods







Saliency Example - Gradients

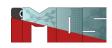


$$f(x): R^d \to R$$

$$E(f)(x) = \frac{df(x)}{dx}$$



Saliency Example – Leave-one-out



$$f(x): \mathbb{R}^d \to \mathbb{R}$$

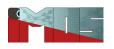
$$E(f)(x)_i = f(x) - f(x \setminus i)$$

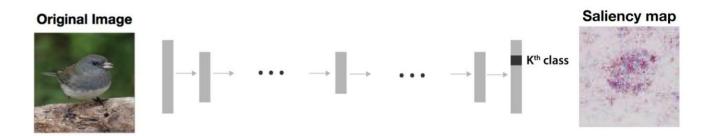
How to remove?

- 1. Zero out pixels in image
- 2. Remove word from the text
- 3. Replace the value with population mean in tabular data



Sanity check: When predictions change, do the explanations change?







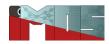


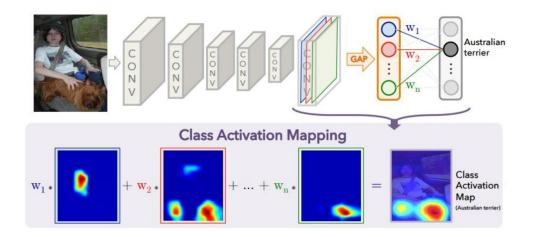


Source: Julius Adebayo



Class Activation Mapping

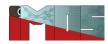




- GAP: Global average pooling
- We can identify the importance of the image regions by projecting back the weights of the output layer on the convolutional feature maps obtained from the last Convolution Layer.

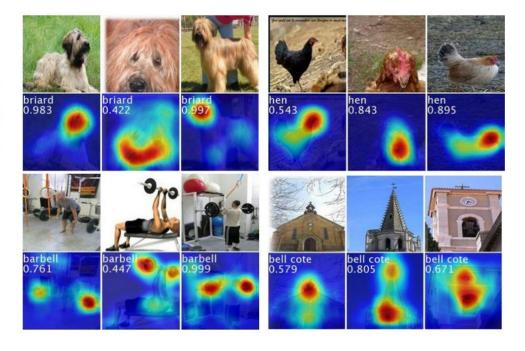


Class Activation Mapping: examples



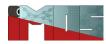
Brushing teeth

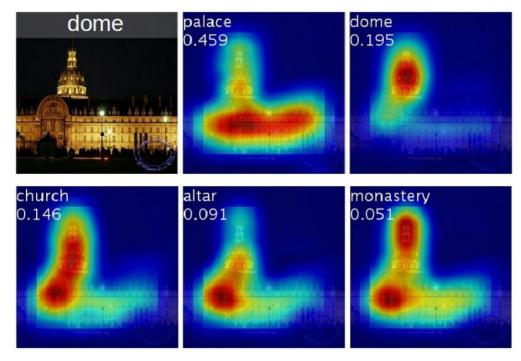
Cutting trees





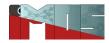
Class Activation Mapping

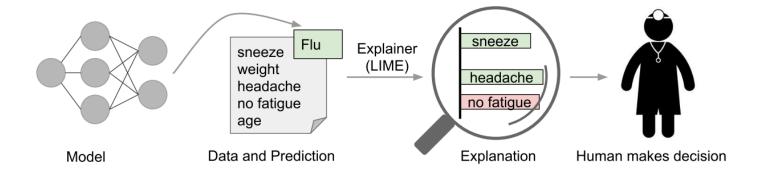






Local Interpretable Model-Agnostic Explanations (LIME): basic idea

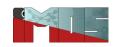


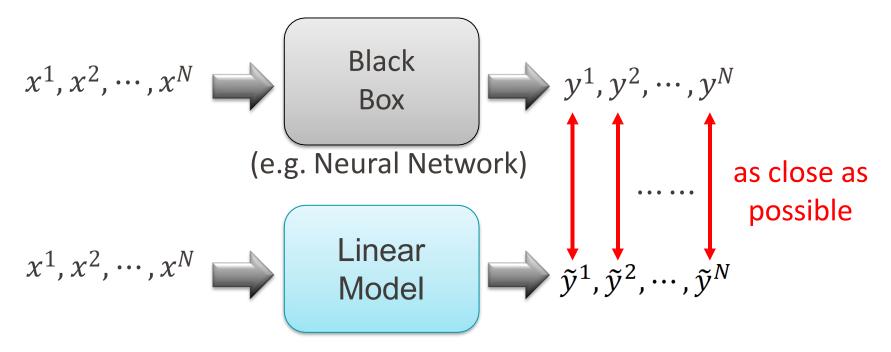






LIME - locally interpretable model agnostic

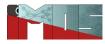


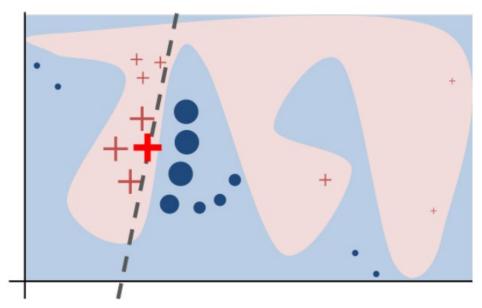


Can't do it globally of course, but locally? Main Idea behind LIME



LIME: Toy example of the basic concept



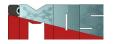


04.11.24

M. T. Ribeiro, S. Singh, and C. Guestrin, "Why should I trust you?' Explaining the predictions of any classifier," in *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2016, vol. 13-17-Augu, pp. 1135–1144.



LIME: divide images into interpretable components (contiguous superpixels)





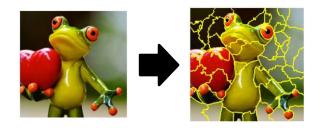
Original Image

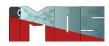


Interpretable Components

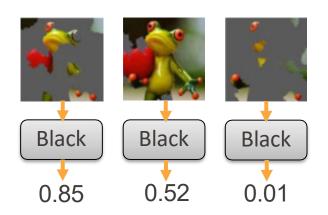


LIME — Image





- 1. Given a data point you want to explain
- 2. Sample at the nearby Each image is represented as a set of superpixels (segments).

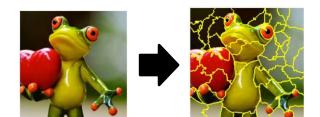


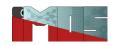
- Randomly delete some
- segments.

Compute the probability of "frog" by black box

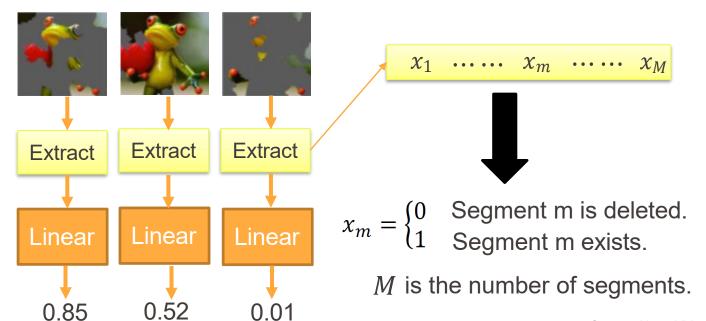


LIME — Image

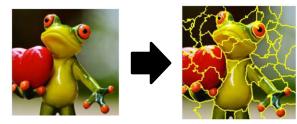


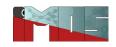


• 3. Fit with linear (or interpretable) model

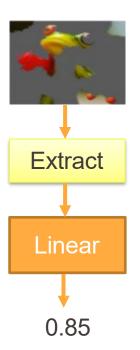


LIME — Image





4. Interpret the model you learned



$$y = w_1 x_1 + \cdots + w_m x_m + \cdots + w_M x_M$$

$$x_m = \begin{cases} 0 & \text{Segment m is deleted.} \\ 1 & \text{Segment m exists.} \end{cases}$$

M is the number of segments.

If
$$w_m \approx 0$$

If $W_m \approx 0$ segment m is not related to "frog"

If
$$w_m$$
 is positive

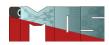
If w_m is positive \blacksquare segment m indicates the image is "frog"

If
$$w_m$$
 is negative

If w_m is negative \blacksquare segment m indicates the image is not "frog"



LIME: underlying algorithm



Algorithm 1 Sparse Linear Explanations using LIME

Require: Classifier f, Number of samples N

Require: Instance x, and its interpretable version x'

Require: Similarity kernel π_x , Length of explanation K

$$\mathcal{Z} \leftarrow \{\}$$

for $i \in \{1, 2, 3, ..., N\}$ do
 $z'_i \leftarrow sample_around(x')$
 $\mathcal{Z} \leftarrow \mathcal{Z} \cup \langle z'_i, f(z_i), \pi_x(z_i) \rangle$

end for

 $w \leftarrow \text{K-Lasso}(\mathcal{Z}, K) \triangleright \text{with } z_i' \text{ as features, } f(z) \text{ as target}$ return w

Match interpretable model to black box

Control complexity of the model

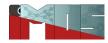
$$\xi(x) = \underset{g \in G}{\operatorname{argmin}} \ \mathcal{L}(f, g, \pi_x) + \Omega(g)$$

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

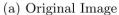
70



LIME: Example











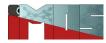
(b) Explaining Electric guitar (c) Explaining Acoustic guitar



(d) Explaining Labrador



LIME: bringing trust («Husky vs Wolf)





(a) Husky classified as wolf

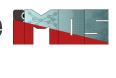


(b) Explanation

	Before	After
Trusted the bad model	10 out of 27	3 out of 27
Snow as a potential feature	12 out of 27	25 out of 27



SHAP (Shapley values Additive exPlanations) – additive feature attribution method to explain the output of any ML model

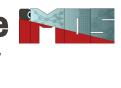


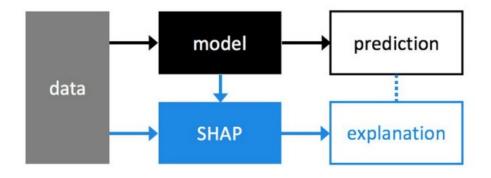
Classic result in game theory on distributing gain in a coalition game

- **Coalition Games**
 - Players collaborating to generate some gain (think: revenue)
 - Set function v(S) determining the gain for any subset S of players
- **Shapley Values** are a fair way to attribute the total gain to the players based on their contributions
 - Concept: Marginal contribution of a player to a subset of other players (v(S U {i}) v(S))
 - Shapley value for a player is a specific weighted aggregation of its marginal over all possible subsets of other players



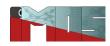
SHAP (Shapley values Additive exPlanations) – additive feature attribution method to explain the output of any **ML** model





The Shapley value is a mathematical concept in game theory that was introduced by Lloyd Shapley in 1951.

Shapley values



$$\phi_i = \sum_{S \subseteq F \setminus \{i\}} \frac{|S|!(|F| - |S| - 1)!}{|F|!} \left[f_{S \cup \{i\}}(x_{S \cup \{i\}}) - f_S(x_S) \right]$$

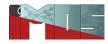
where S feature subset, F – set of all features, $\ S \subseteq F \setminus \{i\}$ all possible subsets

 $f_S{\cup}\{i\}$ is trained with that feature present and f_S with that feature withheld

predictions from the two models are then compared $f_{S \cup \{i\}}(x_{S \cup \{i\}}) - f_{S}(x_{S})$



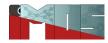
Individual SHAP value plot

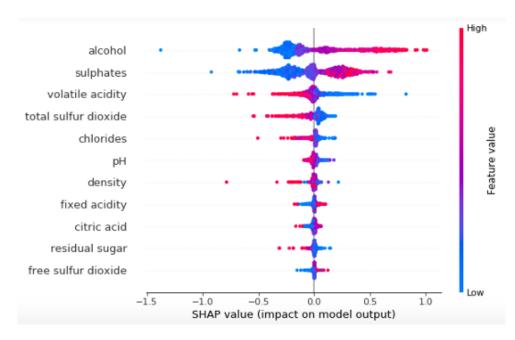






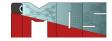
SHAP variables importance plot

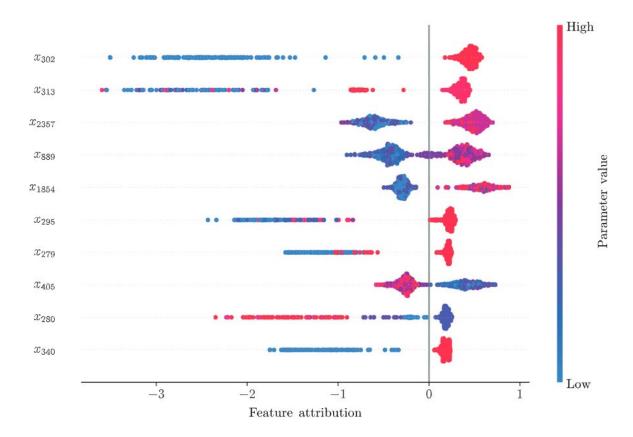






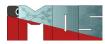
Drivers of the semiconductor production quality

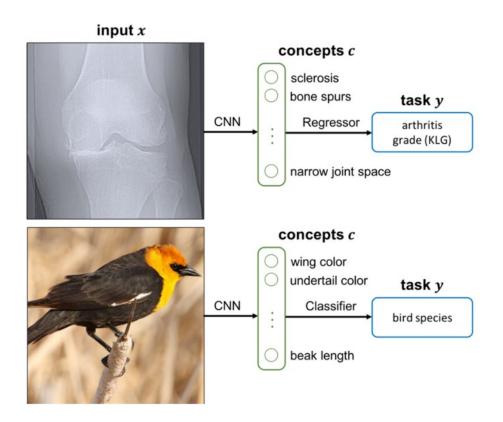






Concept bottleneck models









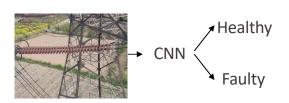
Example: Detecting defective insulators



Detecting defective insulators



- XAI on Power Grid Insulators:
 - Fault detection → binary classifier
 - When faulty, apply XAI to show which part of the image led to the decision,
 - highlighting the defect region



Interpretability



Fault localization



Evaluation of the prediction results

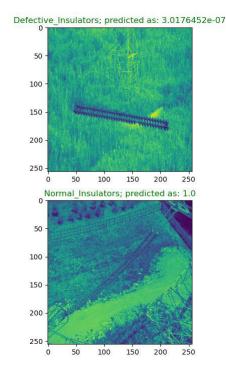


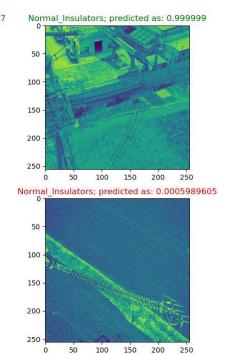
Test Accuracy = 96.4%

FPR = 3.6 %

FNR= 3.8 %

Test Samples Examples (defective \rightarrow 0; Normal \rightarrow 1)







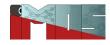


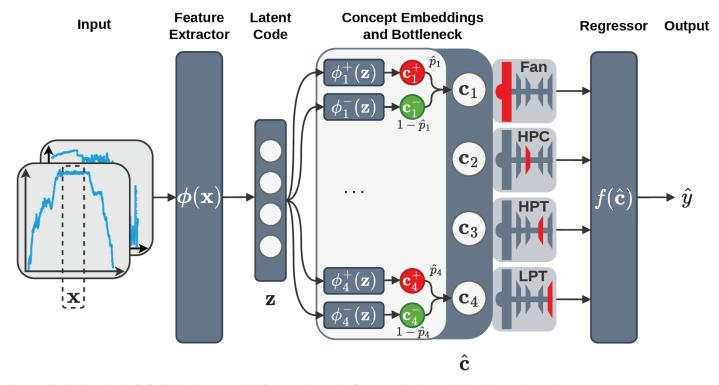
Example:

Integrated concept embeddings



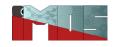
Concept embeddings Turbofan engines







RUL prediction performance RMSE and NASA

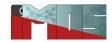


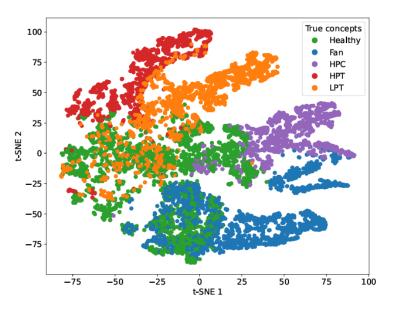
Method	RMSE	NASA score
CNN	6.79	0.851
CEM	6.15	0.752



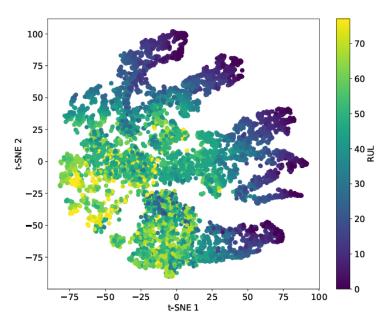


Latent space visualisation





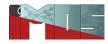
(a) Latent space colored by true concept.

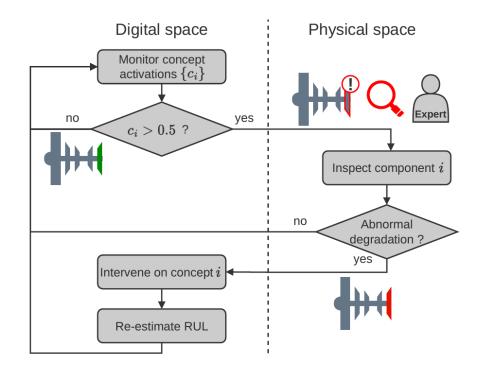


(b) Latent space colored by true RUL value.



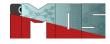
Concept intervention strategy

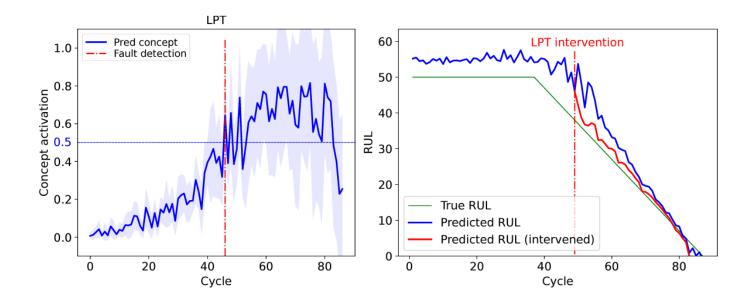






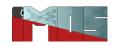
Concept interventions







RUL prediction performance RMSE and NASA



Method	RMSE	NASA score
CNN	6.79	0.851
CEM	6.15	0.752
CEM after intervention	5.96	0.692





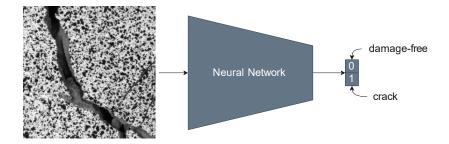


Example:

Using explanations for labeling

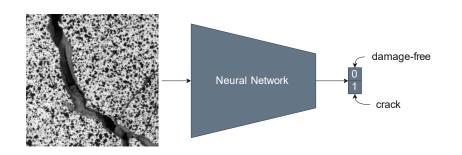
Data-driven approaches based on supervised deep learning have demonstrated excellent performance in detection of cracks in images, but they require **large annotated datasets** for training.

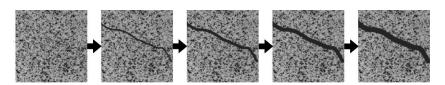
Classification task:



Data-driven approaches based on supervised deep learning have demonstrated excellent performance in detection of cracks in images, but they require large annotated datasets for training.

Classification task:

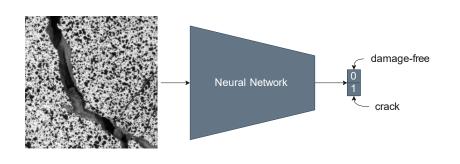


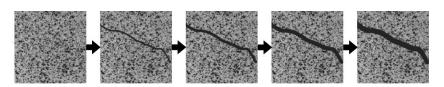


Severity quantification and monitoring is crucial for timely decision-making.

Data-driven approaches based on supervised deep learning have demonstrated excellent performance in detection of cracks in images, but they require **large annotated datasets** for training.

Classification task:





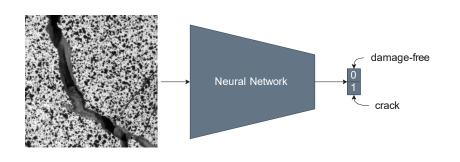
Severity quantification and monitoring is crucial for timely decision-making.

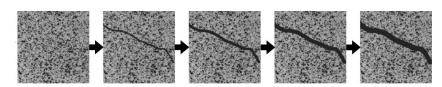


Data-driven approaches based on supervised deep learning have demonstrated excellent performance in

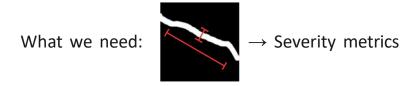
detection of cracks in images, but they require large annotated datasets for training.

Classification task:



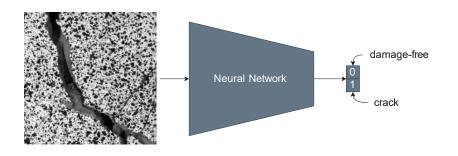


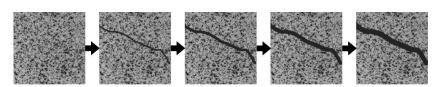
Severity quantification and monitoring is crucial for timely decision-making.



Data-driven approaches based on supervised deep learning have demonstrated excellent performance in detection of cracks in images, but they require large annotated datasets for training.

Classification task:





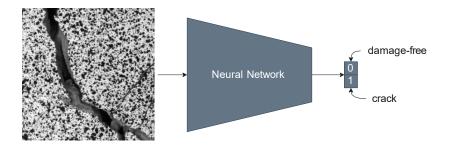
Severity quantification and monitoring is crucial for timely decision-making.



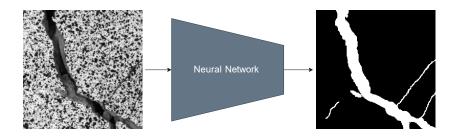
- Does not allow severity quantification
- Fast and easy image-level annotation (1 bit)

Data-driven approaches based on supervised deep learning have demonstrated excellent performance in detection of cracks in images, but they require large annotated datasets for training.

Classification task:



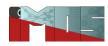
Semantic segmentation task:



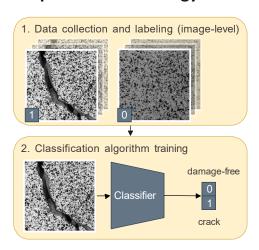
- Does not allow severity quantification
- Fast and easy image-level annotation (1 bit)

- Allows severity quantification and monitoring
- Tedious and costly pixel-level annotation $(256 \times 256 \rightarrow 2^{16} = 64 \text{ Kb})$

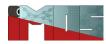




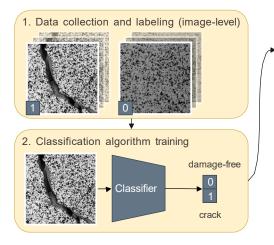
Proposed methodology:

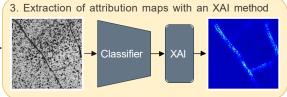




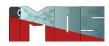


Proposed methodology:

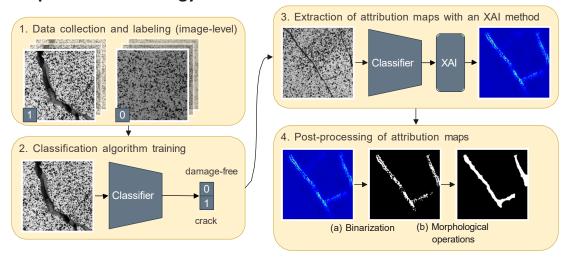




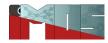




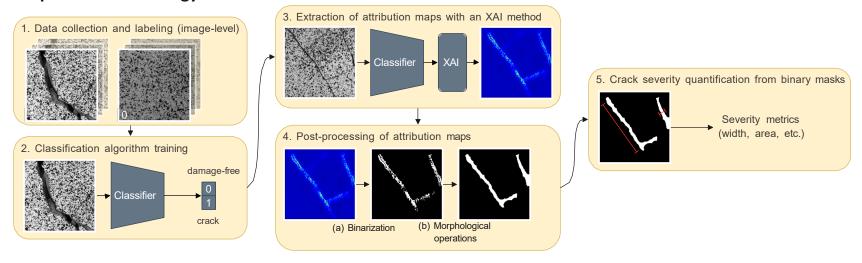
Proposed methodology:







Proposed methodology:

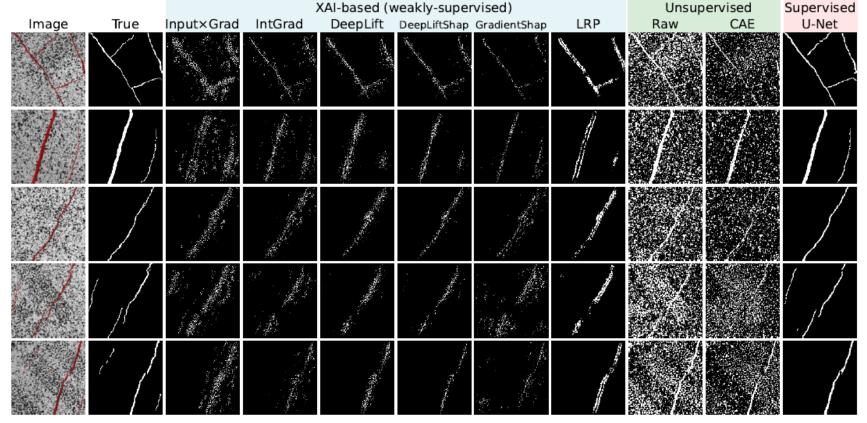


100



Visualization after binarization





04.11.24



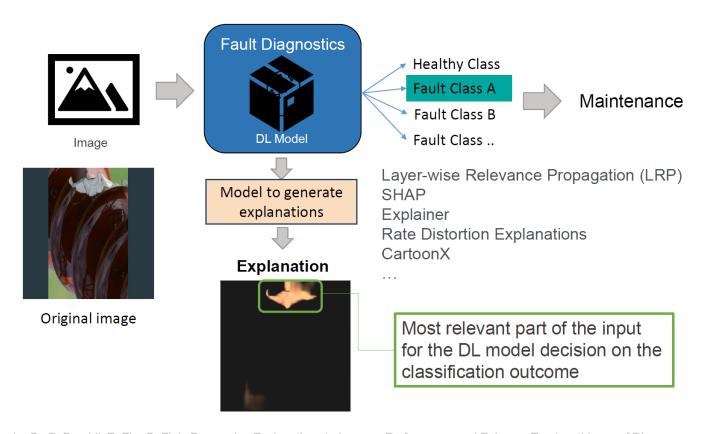


Automated explanation assessments



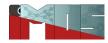
Processing explanations

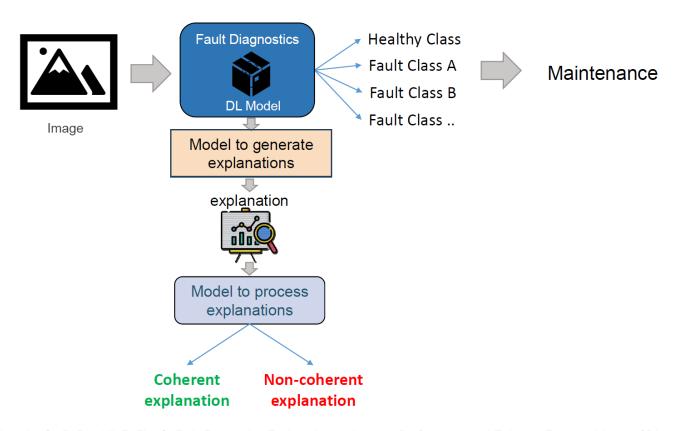






How can explanations be used?

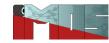




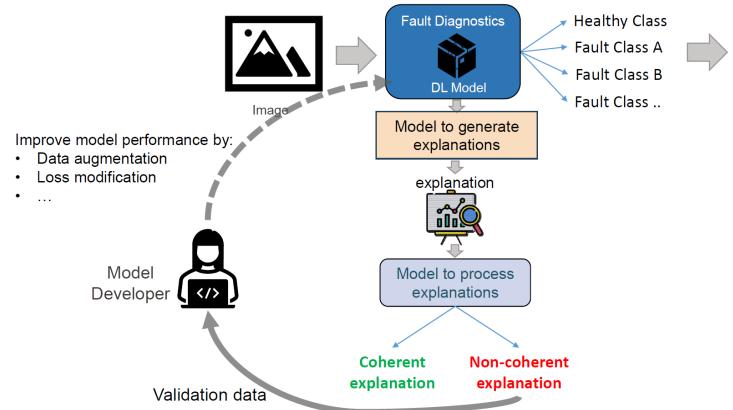
104



Model improvement

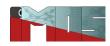


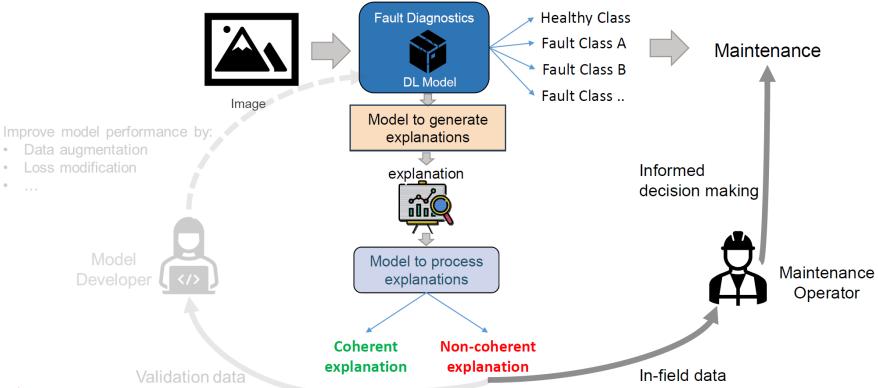
Maintenance





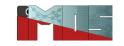
Informed decision making

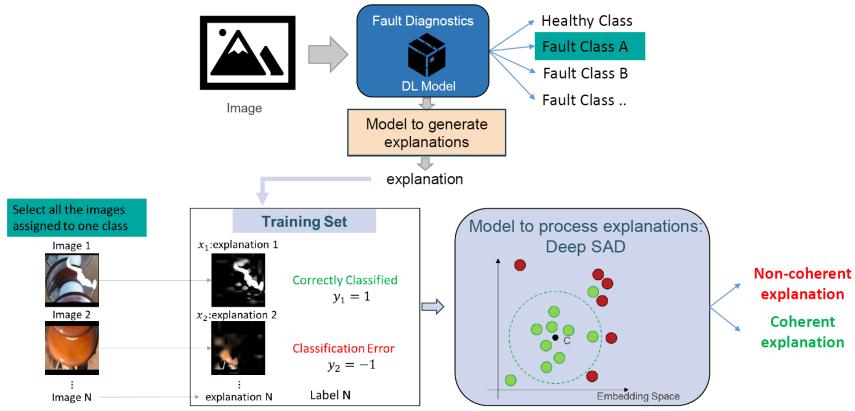






Proposed method for automating the evaluation

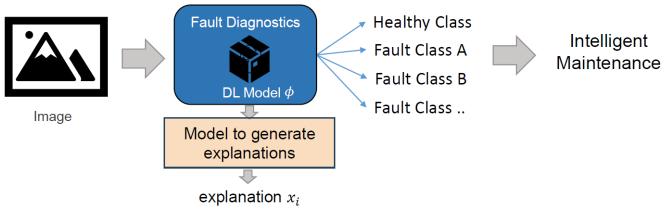




04.11.24

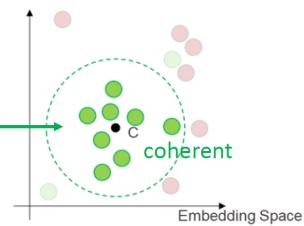
EPFL Proposed method: detecting unusual explanations

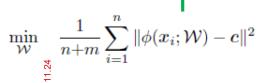




Embedding space definition: loss function

Correct classifications: minimize the distance from the centre

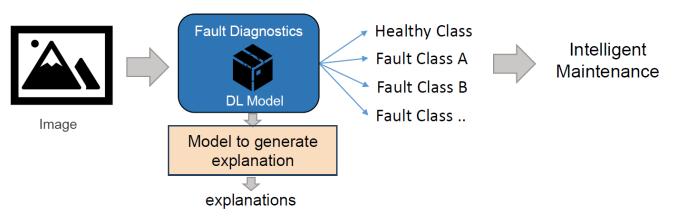




Floreale, G., P. Baraldi, E. Zio, O. Fink: Processing Explanations to Improve Performance and Enhance Trustworthiness of DL Models for Fault Diagnostics of Power Grid Components, in preparation

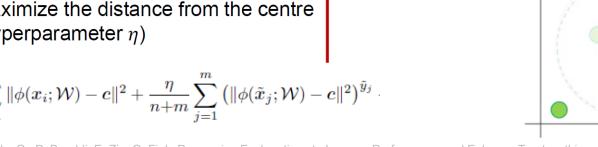
EPFL Proposed method: detecting unusual explanations





Embedding space definition: loss function

Classifications errors: maximize the distance from the centre (hyperparameter η)

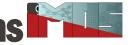


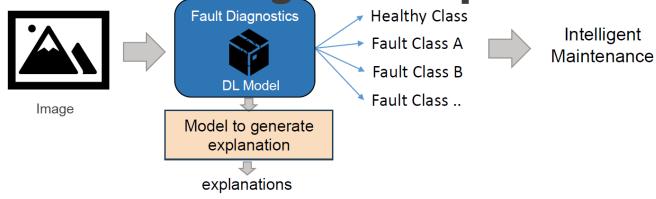
Floreale, G., P. Baraldi, E. Zio, O. Fink: Processing Explanations to Improve Performance and Enhance Trustworthiness of DL Models for Fault Diagnostics of Power Grid Components, in preparation

Embedding Space

Non coherent

Proposed method: detecting unusual explanations

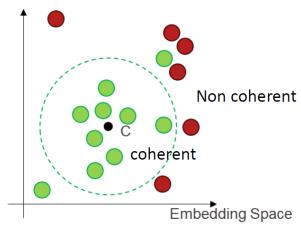




Embedding space definition: loss function

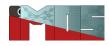
Regularization term to avoid overfitting (hyperparameter λ)

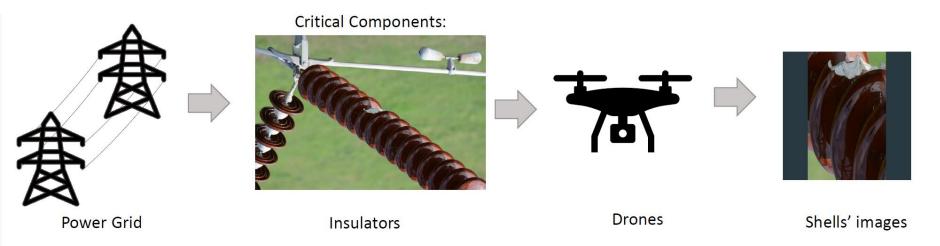
$$\min_{\mathcal{W}} \quad \frac{1}{n+m} \sum_{i=1}^{n} \|\phi(x_i; \mathcal{W}) - c\|^2 + \frac{\eta}{n+m} \sum_{j=1}^{m} \left(\|\phi(\tilde{x}_j; \mathcal{W}) - c\|^2 \right)^{\tilde{y}_j} + \frac{\lambda}{2} \sum_{\ell=1}^{L} \|W^\ell\|_F^2.$$





Case study: infrastructure monitoring

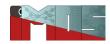


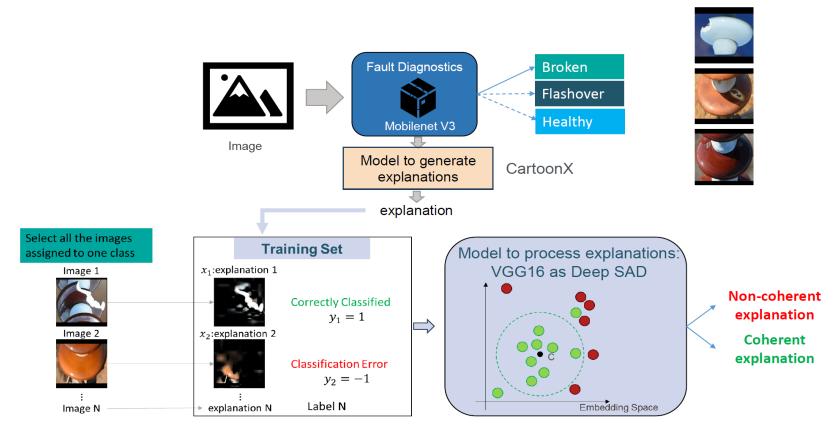


e.g. swiss power grid: 6700 km long → 12 000 pylons



Application of the methodology



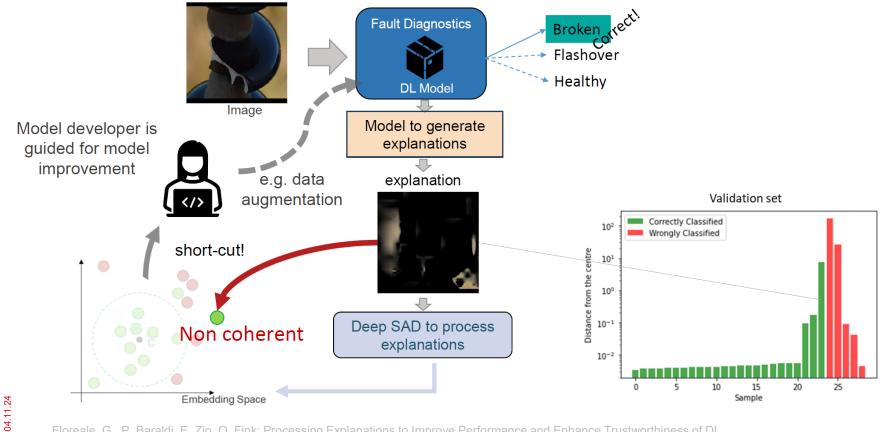


Floreale, G., P. Baraldi, E. Zio, O. Fink: Processing Explanations to Improve Performance and Enhance Trustworthiness of DL Models for Fault Diagnostics of Power Grid Components, in preparation



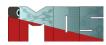
Application of the methodology



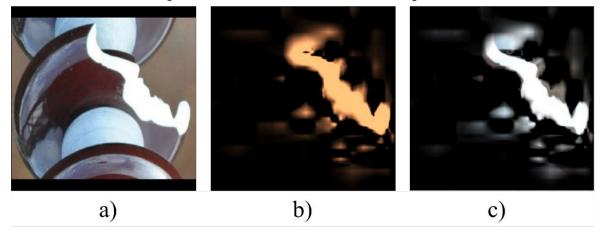




Examples

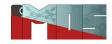


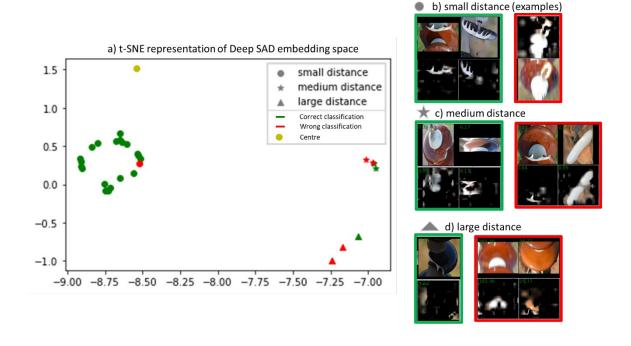
Example of one broken shell and its explanation





Coherent vs. Non-coherent explanations

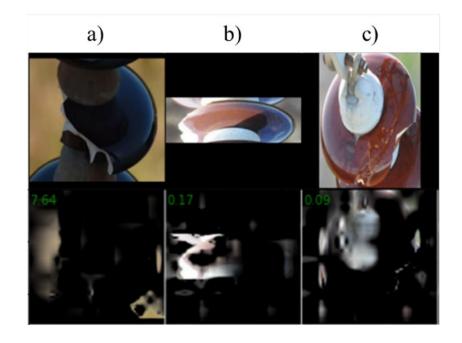






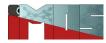
Examples

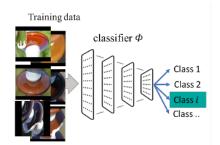




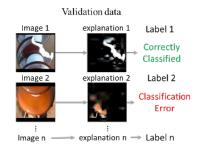


Summary

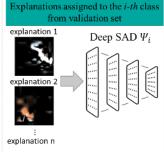




Step 1. Train the supervised classifier Φ

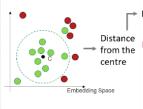


Step 2. Compute explanations and assign binary label for correct/incorrect classifications



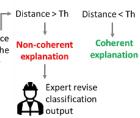
Step 3. Use labelled explanations to Step 4. Apply Ψ_i to test data train Deep SAD model Ψ_i for the i-explanations and map them th class

For each i-th class



Test data assigned to i-th class

in the embedding space



Step 5. Apply threshold to spot non-coherent explanations and ask expert to revise them

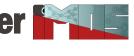


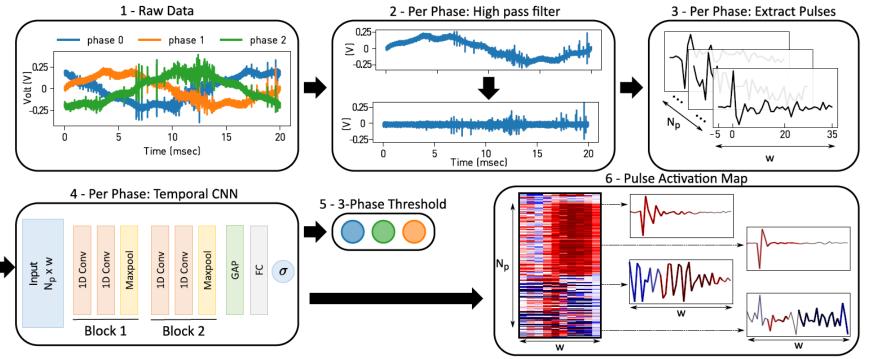


Example:Detection of Partial Discharge in Power Lines



Interpretable Detection of Partial Discharge in Power **Lines with Deep Learning** → **Framework**



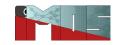


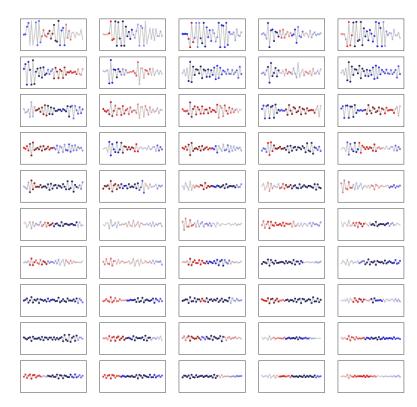
Michau, Gabriel, Chi-Ching Hsu, and Olga Fink. "Interpretable Detection of Partial Discharge in Power Lines with Deep Learning." Sensors 21.6 (2021): 2154.

Olga Fink



Pulse activations







Pulse Activation Maps (PAM)

