



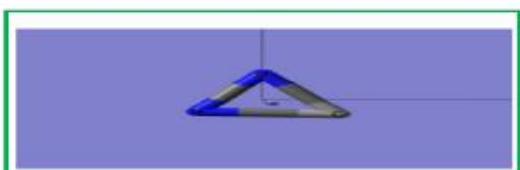
Journal of Applicable Chemistry

2023, 12 (5): 704-778

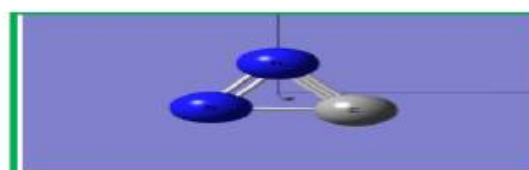
(International Peer Reviewed Journal)



New Chemistry News



New News of Chem (NNC)



ChemNewsNew (CNN)

CNN-56--Fit (Figure Image TableScript...) BasesPart

4. xAI (Bfit) 2022-2023

AI (1950-2023)

Information Source	sciencedirect.com ; ACS.org ;	
S. Narasinga Rao M D Associate Professor, Dept. of General Medicine, Government medical college, government general hospital, Srikakulam, AP, India	K. Somasekhara Rao, Ph D Dept. of Chemistry, Acharya Nagarjuna Univ., Dr. M.R.Appa Rao Campus, Nuzvid-521 201, India	R. Sambasiva Rao, Ph D Dept. of Chemistry, Andhra University, Visakhapatnam 530 003, India
srnaveen007@gmail.com (+91 9848136704)	sr_kaza1947@yahoo.com (+91 98 48 94 26 18)	rsr.chem@gmail.com (+91 99 85 86 01 82)

Conspectus: The start of first order logic based Symbolic Expert Systems dates back to 1960s as tools in the sub-goals of Artificial Intelligence domain. Dendral, Mycin, Xcon/Xsel were earliest the then large expert system products intended for use in organic structure elucidation in chemistry, medical diagnosis and to assist in the ordering of DEC's VAX computer systems. Xcon/Xsel use to automatically select the computer system components based on the customer's requirements. The limitations of ESs of those days were knowledge base was human extracted and antecedents were deterministic. Further they were implemented for complex real-life scenarios. The I/O transformation was transparent and thus, explanation is straight forward to all levels of users/stake holders.

Linear models, regression trees and fuzzy-logic systems are popularly known now as machine

learning tools. Here, data flow/model structure are transparent. The vivid crystal-clear explanation renders them to be called as white-box approaches.

Neural networks of first generation viz. Adaline (Adaptive Linear Neuron), Madaline (Multiple Adaptive Linear Neuron), SLP/MLP, Fuzzy-ARTMAP, SOM, RecNN compute with floating point data and incorporating probability scores. It led to generation of robust models of high accuracy in data-driven mode. But the black-tinge increases with complexity of model. This is the major stumble block to utilize these methods in Medicine, Defence, communication and industry. It necessitated the need to develop transparent models and explanations in multiple modes (visual, If-then-Else rules, numerical derived parameters, 2D-plots, Scripts etc.). With concerted efforts of DARPA, NSF and other agencies, a new trans-discipline called **eXplainable Artificial Intelligence (xAI)** emerged. The evolution of xAI is at a jet speed and serves different categories of stakeholders viz. Human Experts (Hes) with specific domain knowledge, common-man using xAI-imbedded/assisted services, experts in other field but not in discipline imbedding xAI, investors, managers and policy makers involved in approvals/sanctions. The future ventures in Science/ Technology/ products or Tools will be largely based on Trust-worthy/ Responsible/ Safe/ethical AI.

Keywords: Modelling; Artificial Intelligence; Symbolic expert systems; Second Generation AI – Neural Networks; Machine Learning; Deep Level neural architectures; Deep Learning; convolution Networks; ALEX; Capsule Nets; explainable Artificial Intelligence (xAI); Applications; DARPA Stipulations; NSF; European union;

Artificial Intelligence [1950-2023 ...]

xAI.

I(T)O.xAI

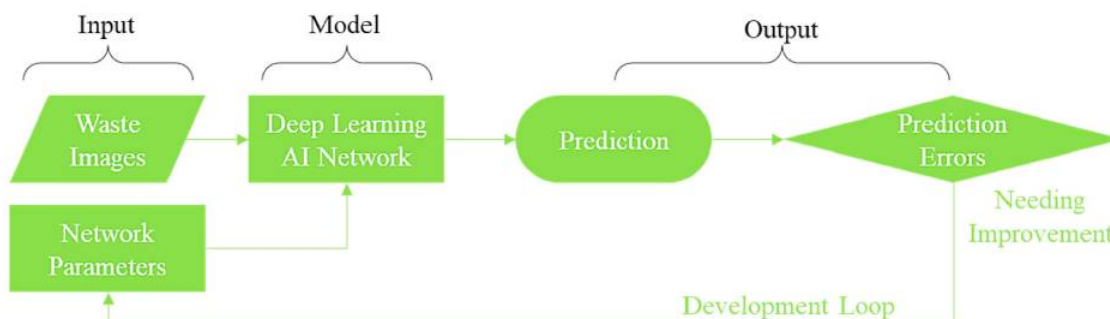
2022-074

Input operated by Transformer giving Output explained by xAI methods

xAI.

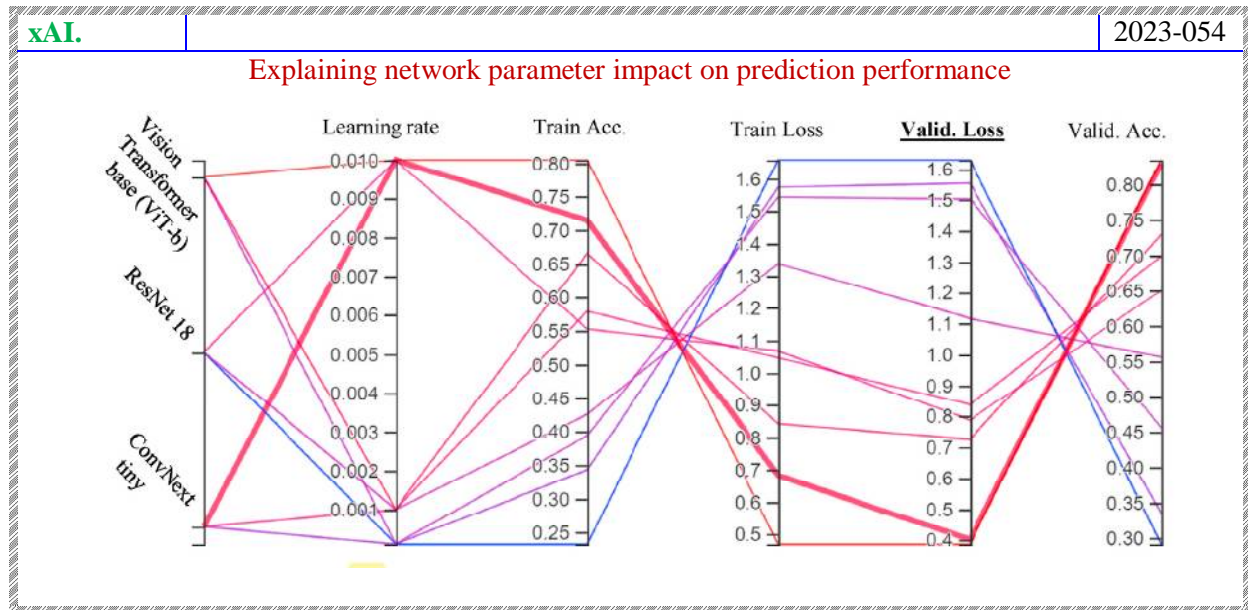
2023-054

Triadic structure of AI networks: input, model and output



Explanation

Input, output, parameters, Network flow



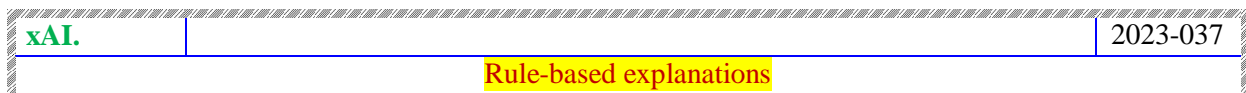
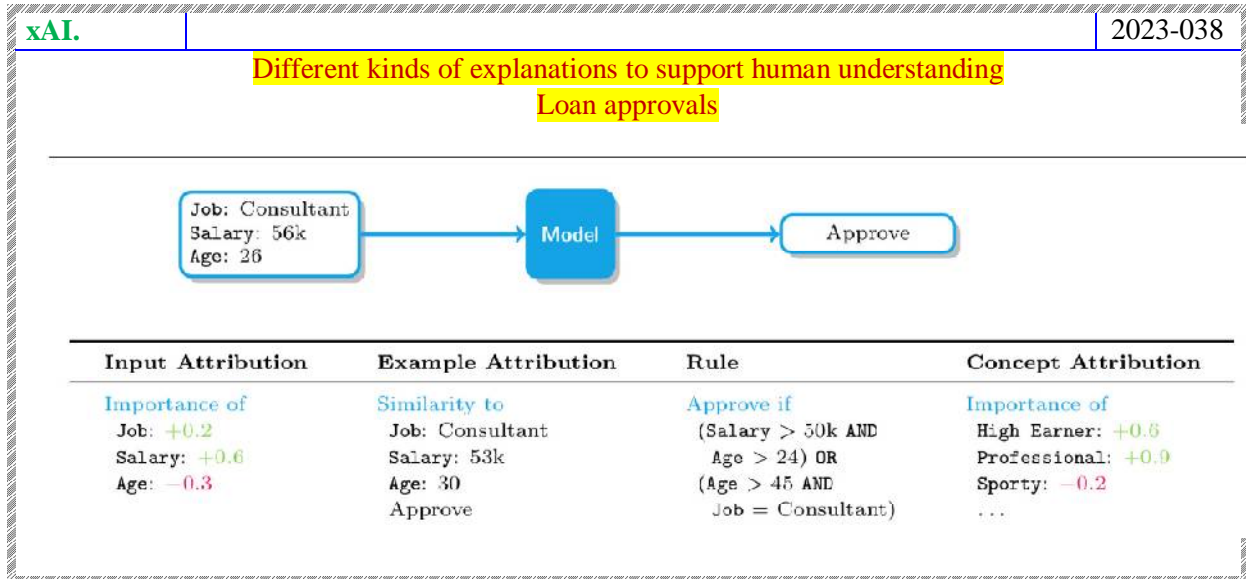
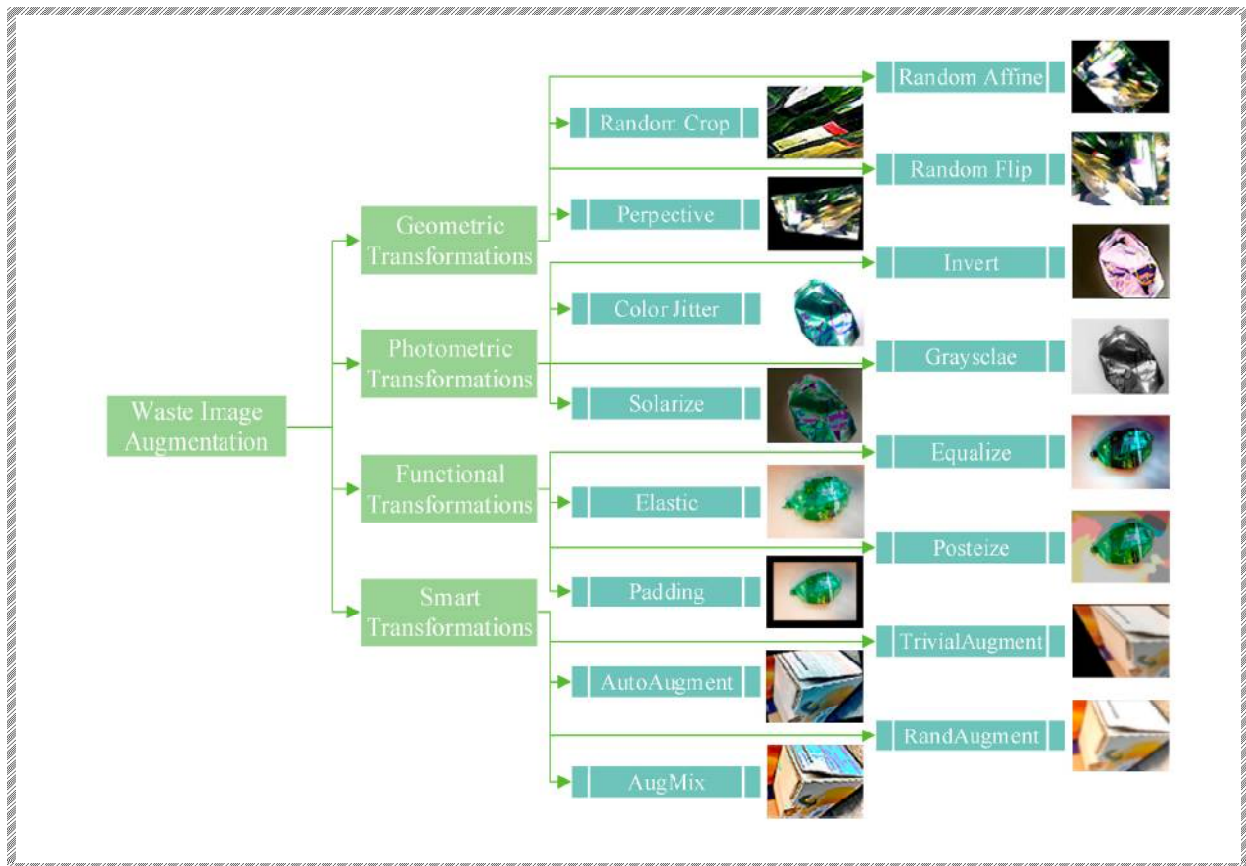
Input & Output

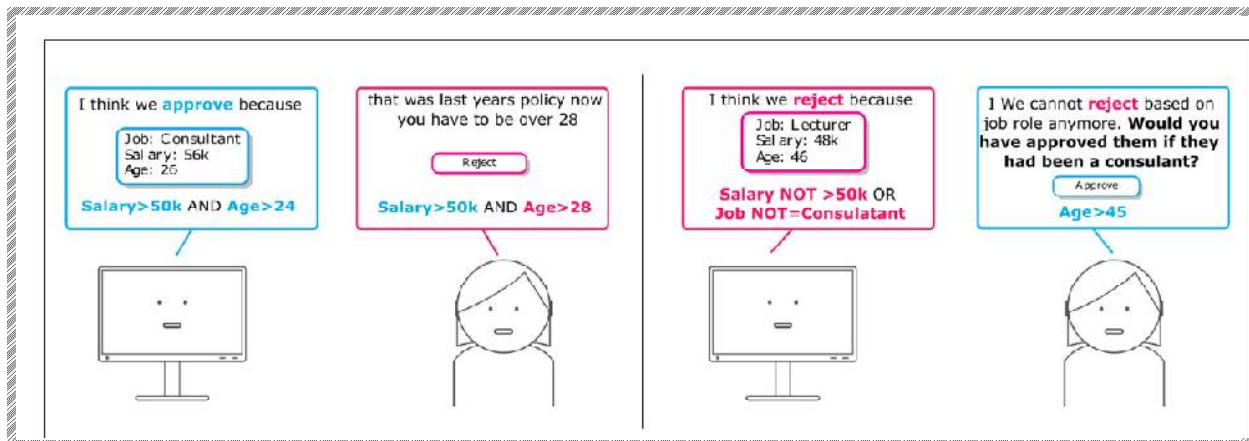
Augmentation.Input

xAI. 2023-054

Input augmentation

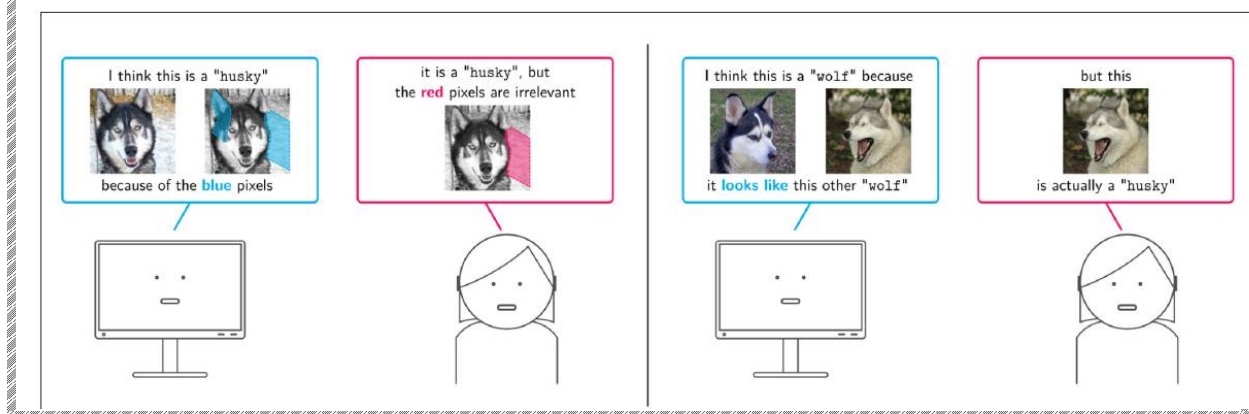
Input augmentations to maximize AI network generalizability and minimize overfitting





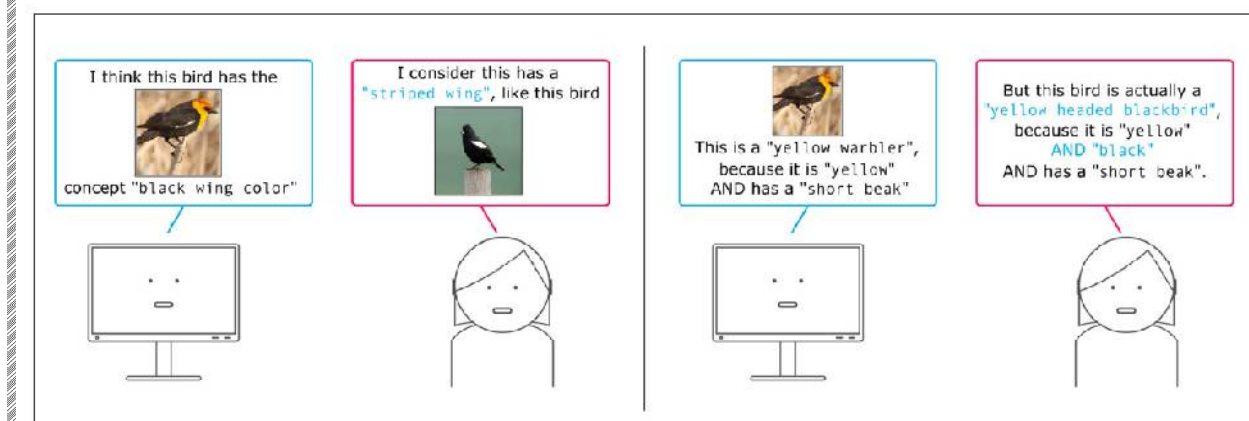
xAI. 2023-038

Predictions by highlighting relevant input Variables Machine justifies its predictions in terms of Training examples



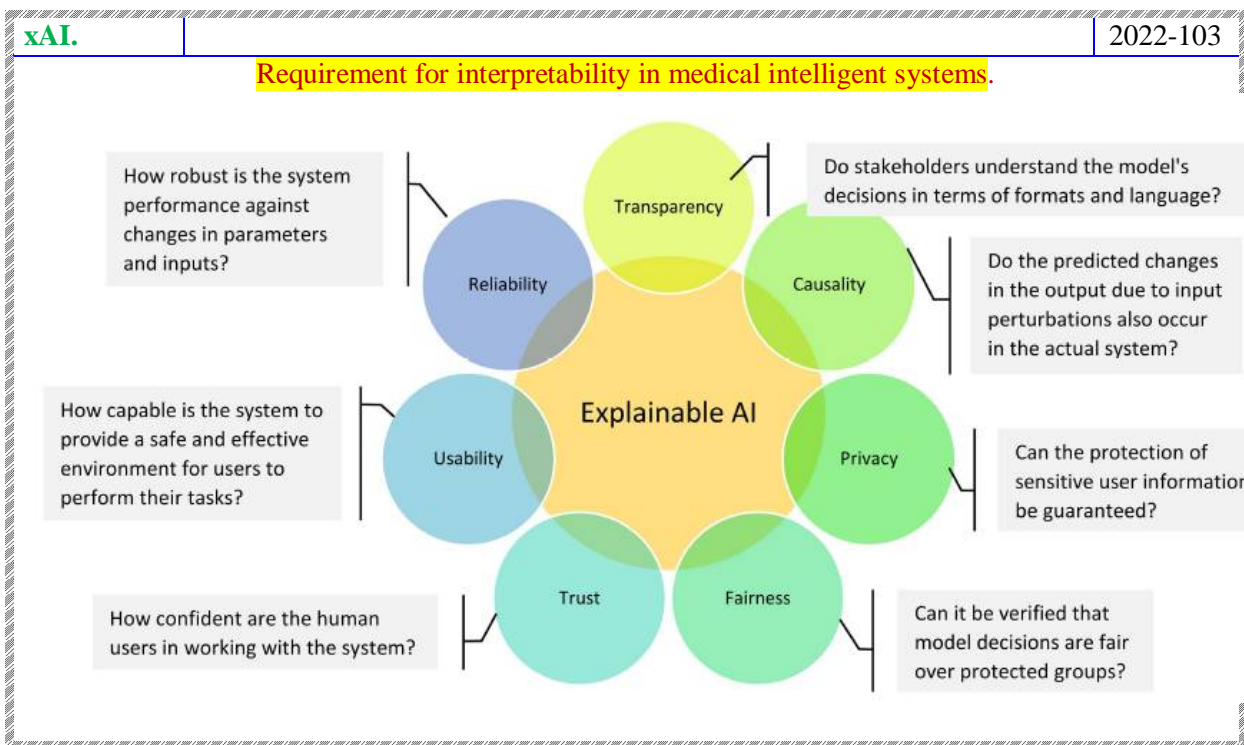
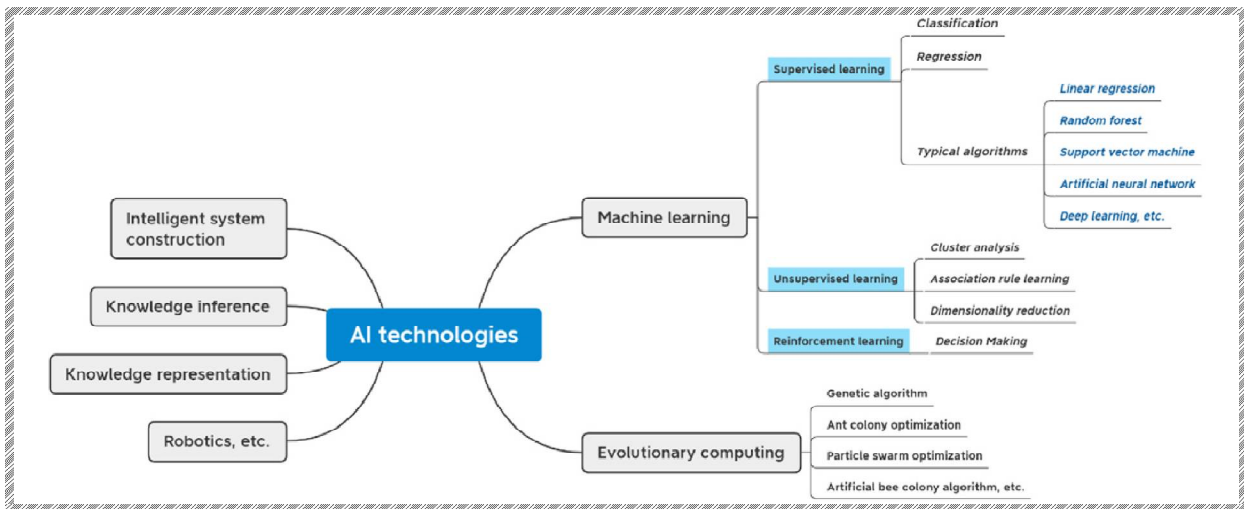
xAI. 2023-038

Concept learning a model Known basic concepts a model



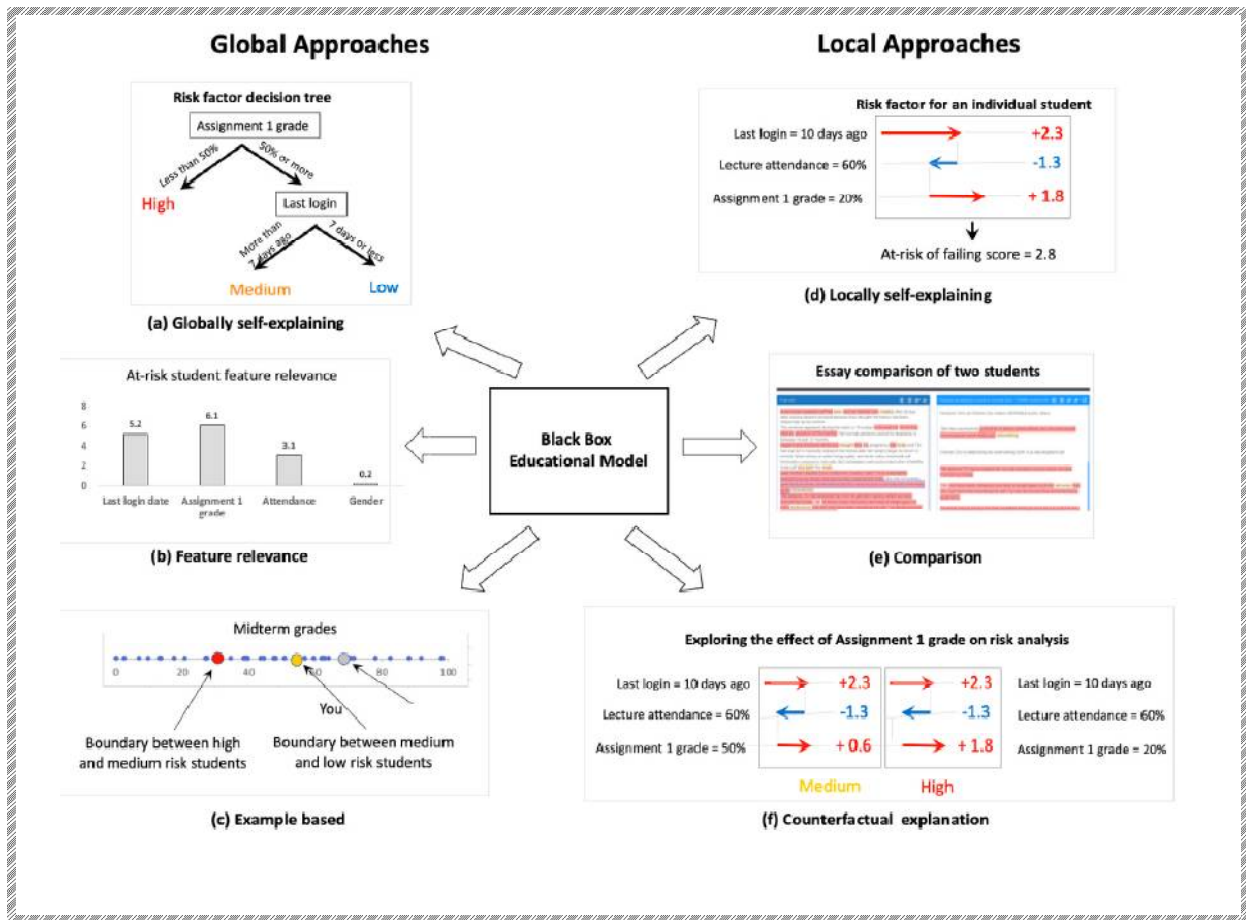
xAI. 2023-039

Framework of AI technologies



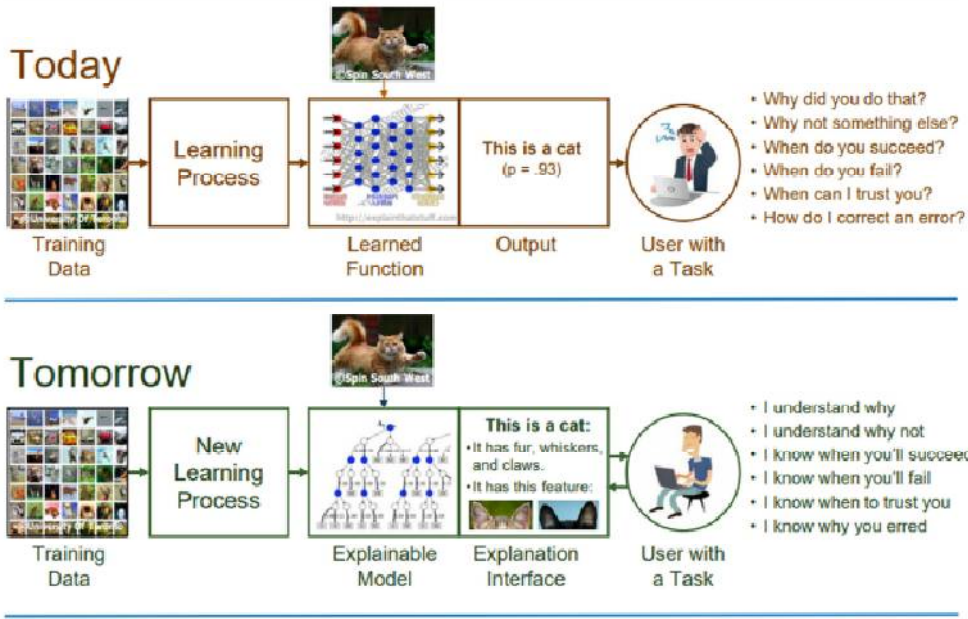
xAI. | 2022-172

Explainability approaches

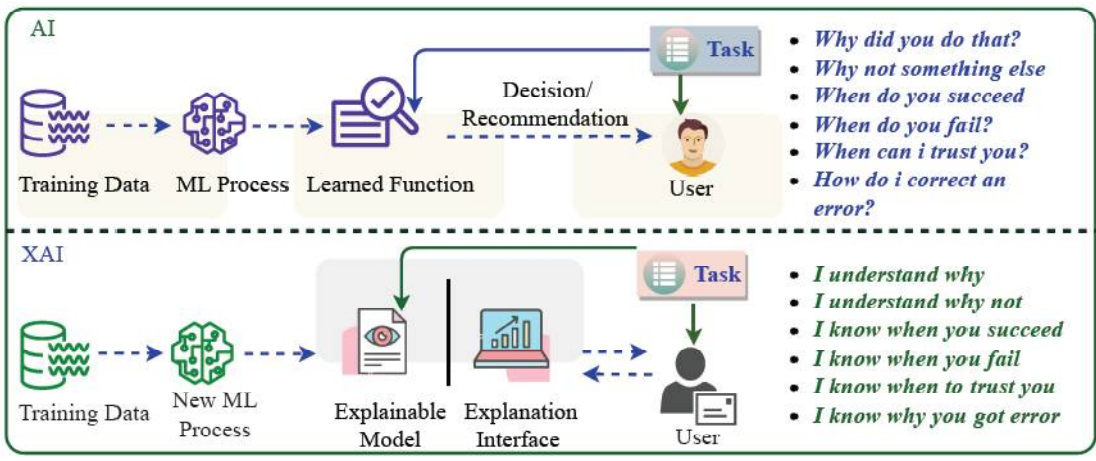


DARPA

xAI.	Expected effect through application of XAI	2022-159
-------------	---	----------

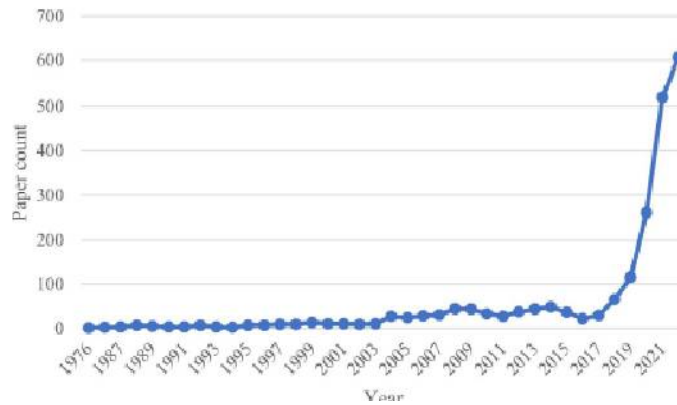


Differences in methodology of AI and XAI



xAI LiteratureSearch

Number of XAI publications added per year from 1976 to 2021

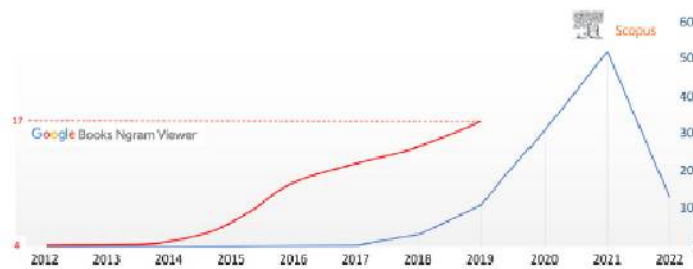


- ✓ Database: PubMed (accessed on 1 July 2022: <https://pubmed.ncbi.nlm.nih.gov>)
- ✓ Terms searched: (explainable AI OR explainable artificial intelligence) AND (medicine OR healthcare)

xAI.

2023-116

Last 10-year trends in academic publications in xAI

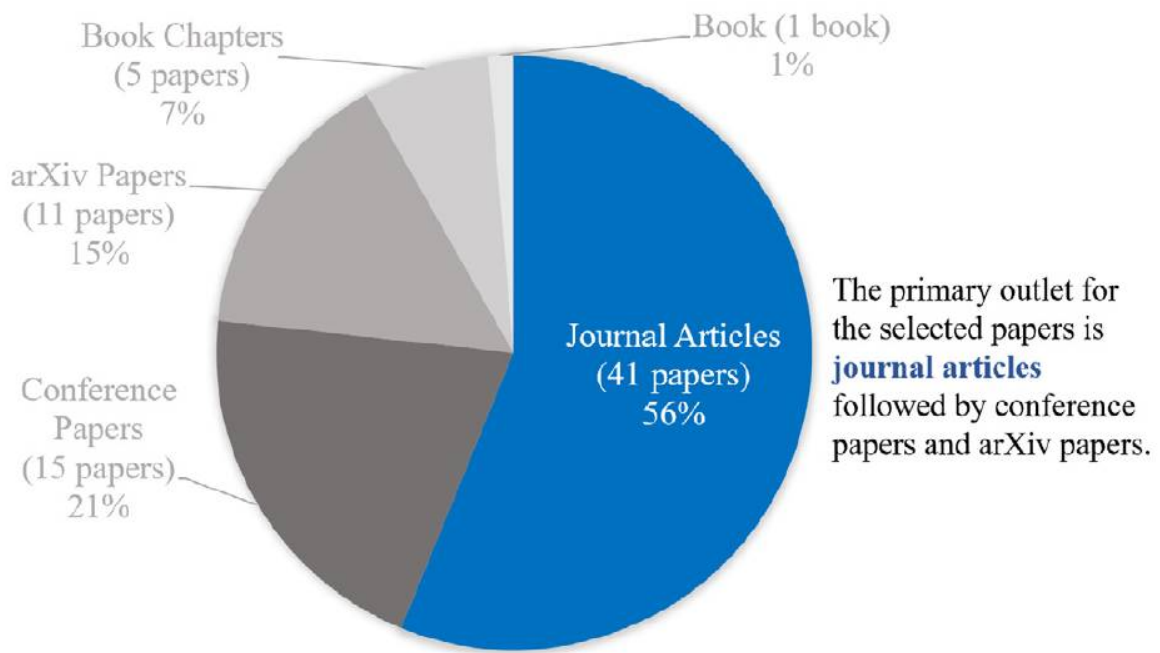


- ✓ Last 10-year trends in academic production (blue line)
- ✓ Articles indexed in: Elsevier Scopus database, and generic publishing (red line)

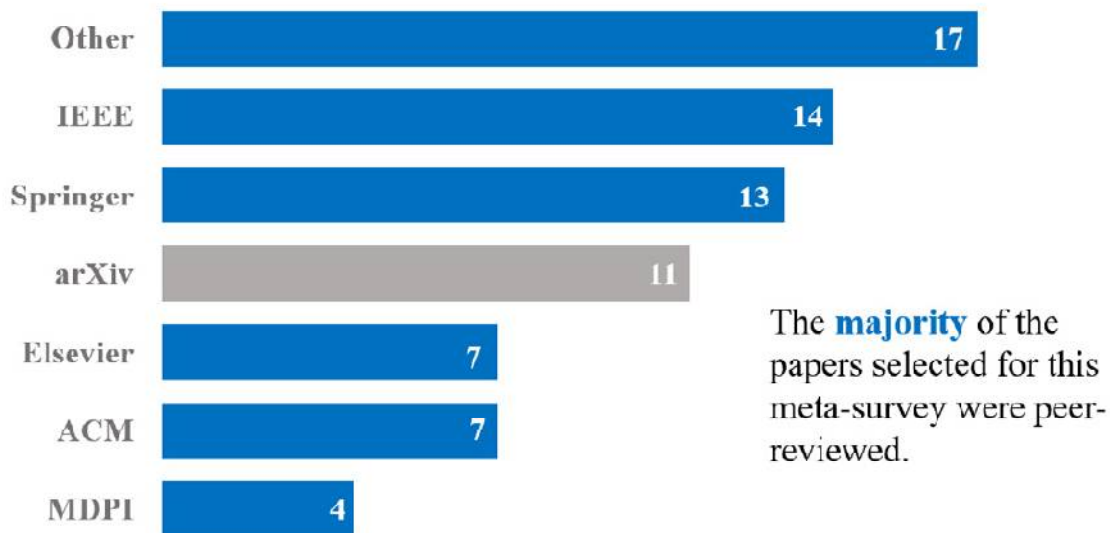
xAI.

2023-123

Distribution of Selected Papers Per Publication Type



Distribution of Selected Papers Per Publisher



Recent research published in 2020 and 2021 has focused on most of these challenges and research directions

Challenges and research directions	2017	2018	2019	2020	2021	Total
Towards more formalism	1	4	4	11	11	31
Explanations and the nature of user experience and expertise	1	1	2	6	2	12
XAI for trustworthiness AI	1	1	0	4	5	11
Multidisciplinary research collaborations	0	1	1	3	5	10
Interpretability vs. performance trade-off	0	1	1	3	2	7
XAI for non-image, non-text, and heterogeneous data	0	0	0	3	3	6
Explainability methods composition	0	3	0	1	1	5
Causal explanations	0	1	0	3	1	5
Challenges in the existing XAI models/methods	0	1	1	1	2	5
Contrastive and counterfactual explanations	0	0	1	0	2	3
Communicating uncertainties	0	0	0	2	0	2
Time constraints	1	0	0	1	0	2
Natural language generation	0	0	0	1	0	1
Analyzing assumption-free black-box models, not assumption-based data models	0	0	1	0	0	1
Reproducibility	0	0	0	1	0	1
The economics of explanations	0	1	0	0	0	1

Recent research published in 2020 and 2021 has focused on most of these challenges and research directions

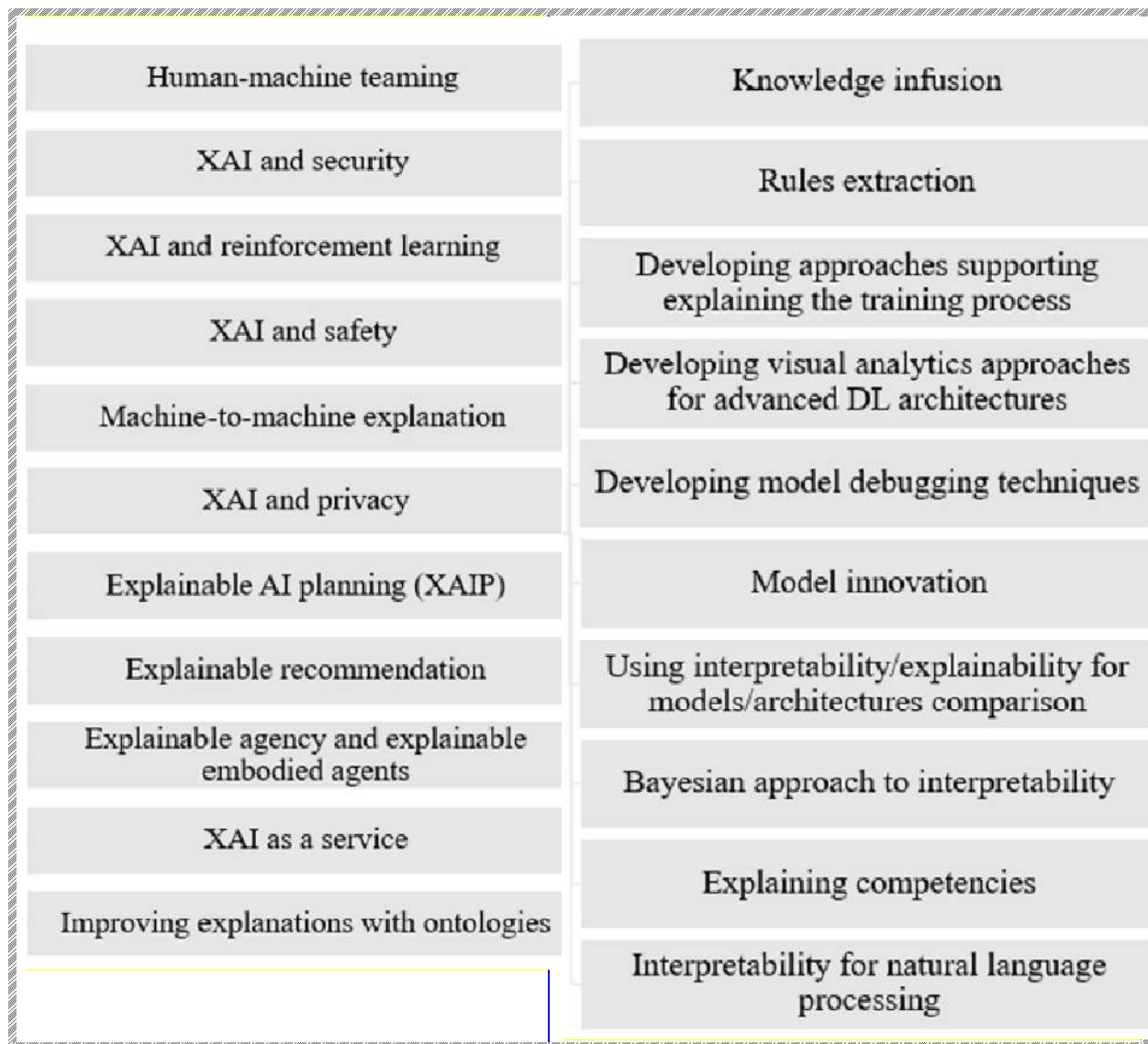
Phases	Challenges and research directions	2017	2018	2019	2020	2021	Total
Design	Communicating data quality	0	0	0	0	2	2
	Sparsity of analysis	0	0	0	1	0	1
Development	Knowledge infusion	0	2	0	4	0	6
	Rules extraction	0	1	0	1	2	4
	Developing approaches supporting explaining the training process	0	1	1	2	0	4
	Developing visual analytics approaches for advanced DL architectures	0	1	0	1	0	2
	Developing model debugging techniques	0	1	0	1	0	2
	Model innovation	0	0	0	0	2	2
	Using interpretability/explainability for models/architectures comparison	0	0	0	1	0	1
	Bayesian approach to interpretability	1	0	0	0	0	1
	Explaining competencies	0	0	1	0	0	1
	Interpretability for natural language processing	0	0	0	0	1	1
Deployment	Human-machine teaming	0	4	2	3	3	12
	XAI and security	0	1	0	3	3	7
	XAI and reinforcement learning	0	1	1	1	3	6
	XAI and safety	0	0	0	3	1	4
	Machine-to-machine explanation	0	2	1	0	0	3
	XAI and privacy	0	1	0	1	0	2
	Explainable AI planning (XAIP)	1	1	0	0	0	2
	Explainable recommendation	0	0	0	2	0	2
	Explainable agency and explainable embodied agents	0	0	1	0	1	2
	XAI as a service	0	0	1	0	1	2
	Improving explanations with ontologies	0	0	0	0	1	1

xAI.

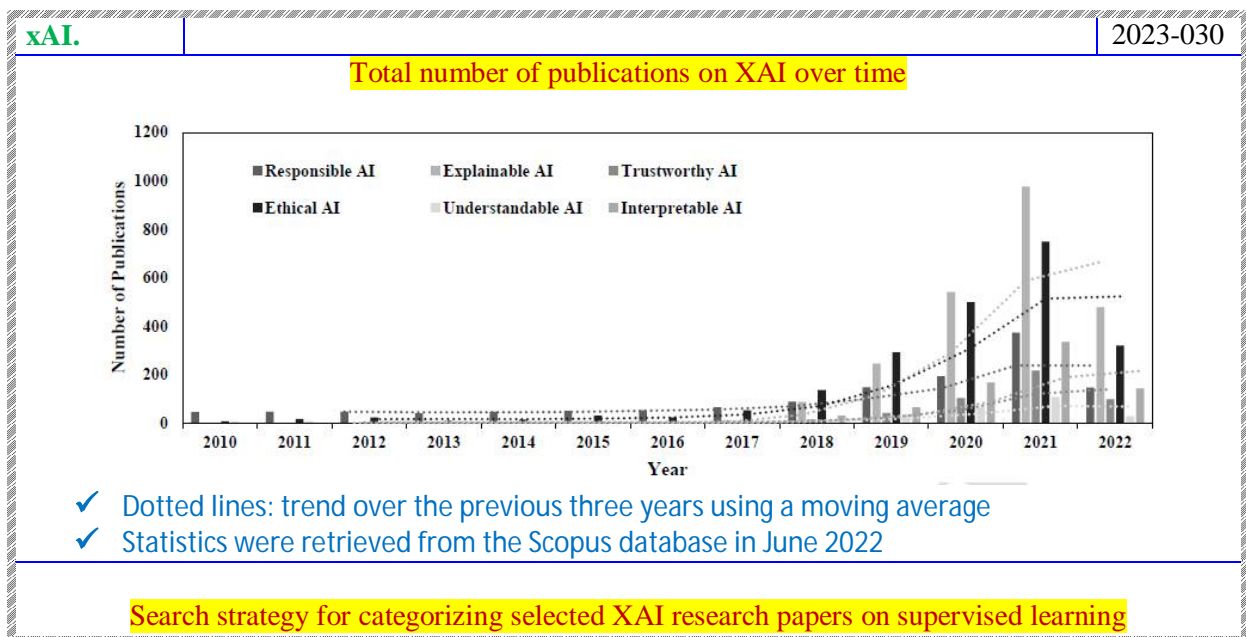
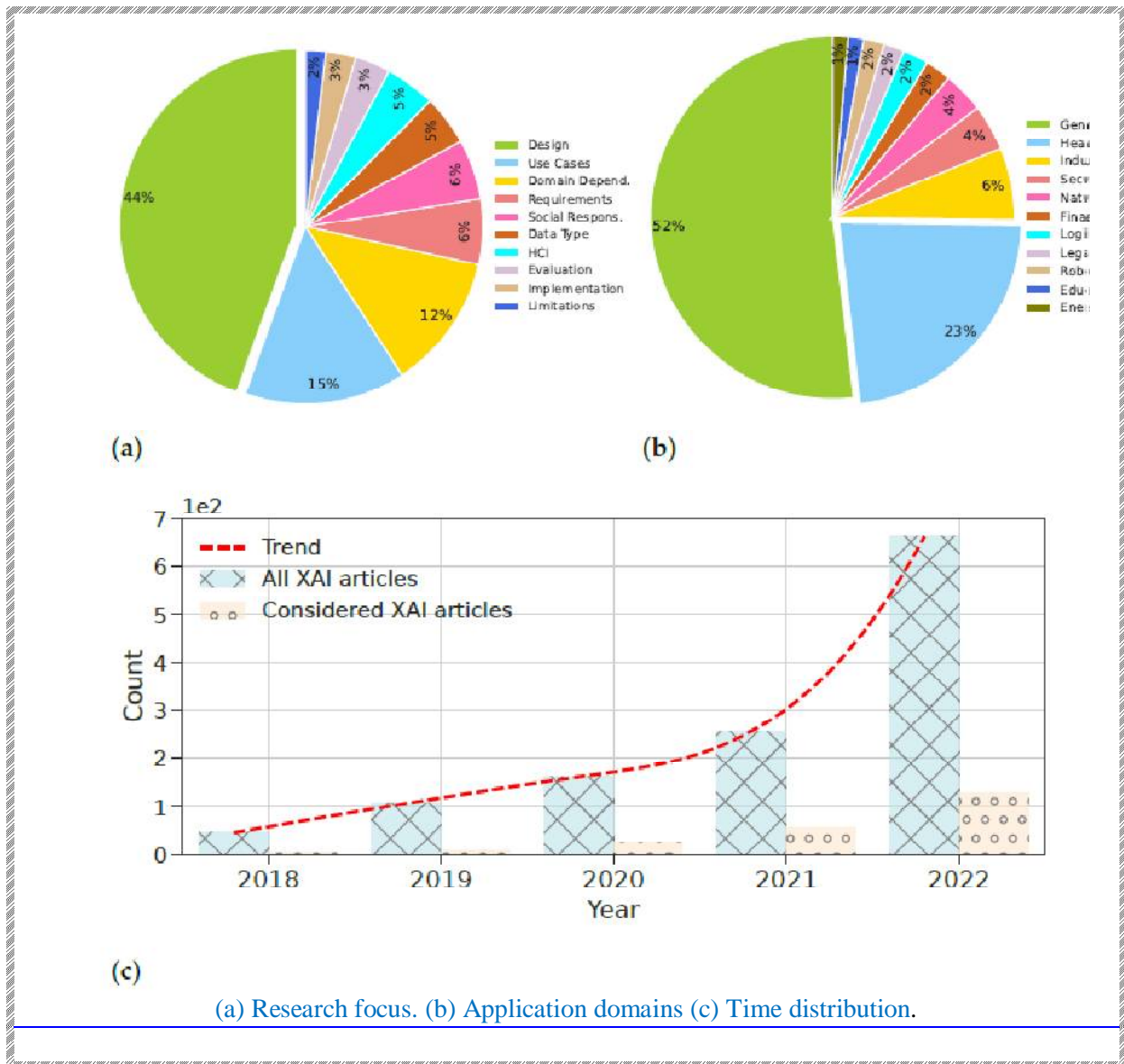
2023-123

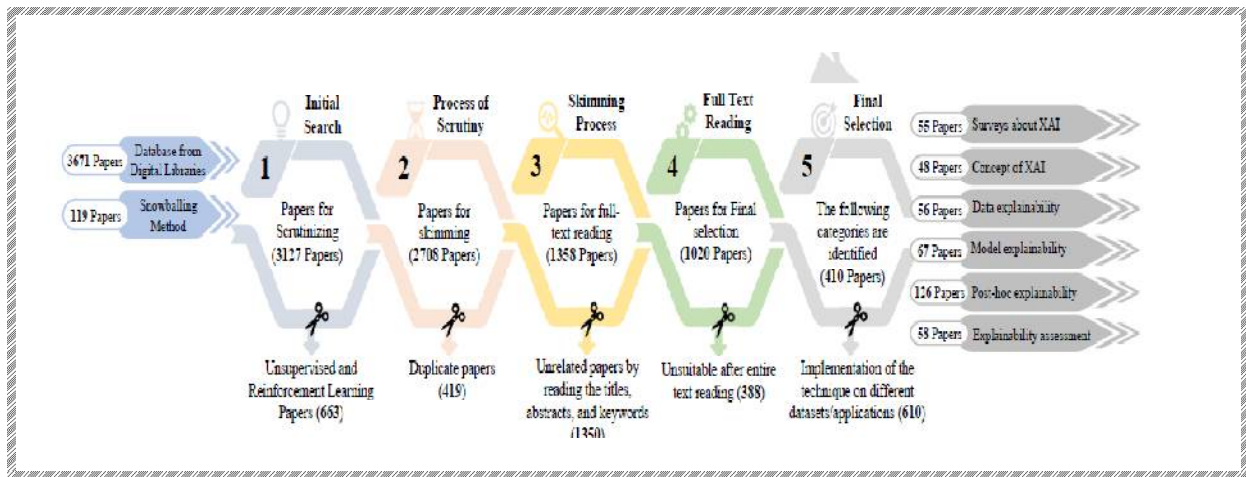
Deployment Phase

Development Phase



xAI.	2023-145
Research publications -- xAI	

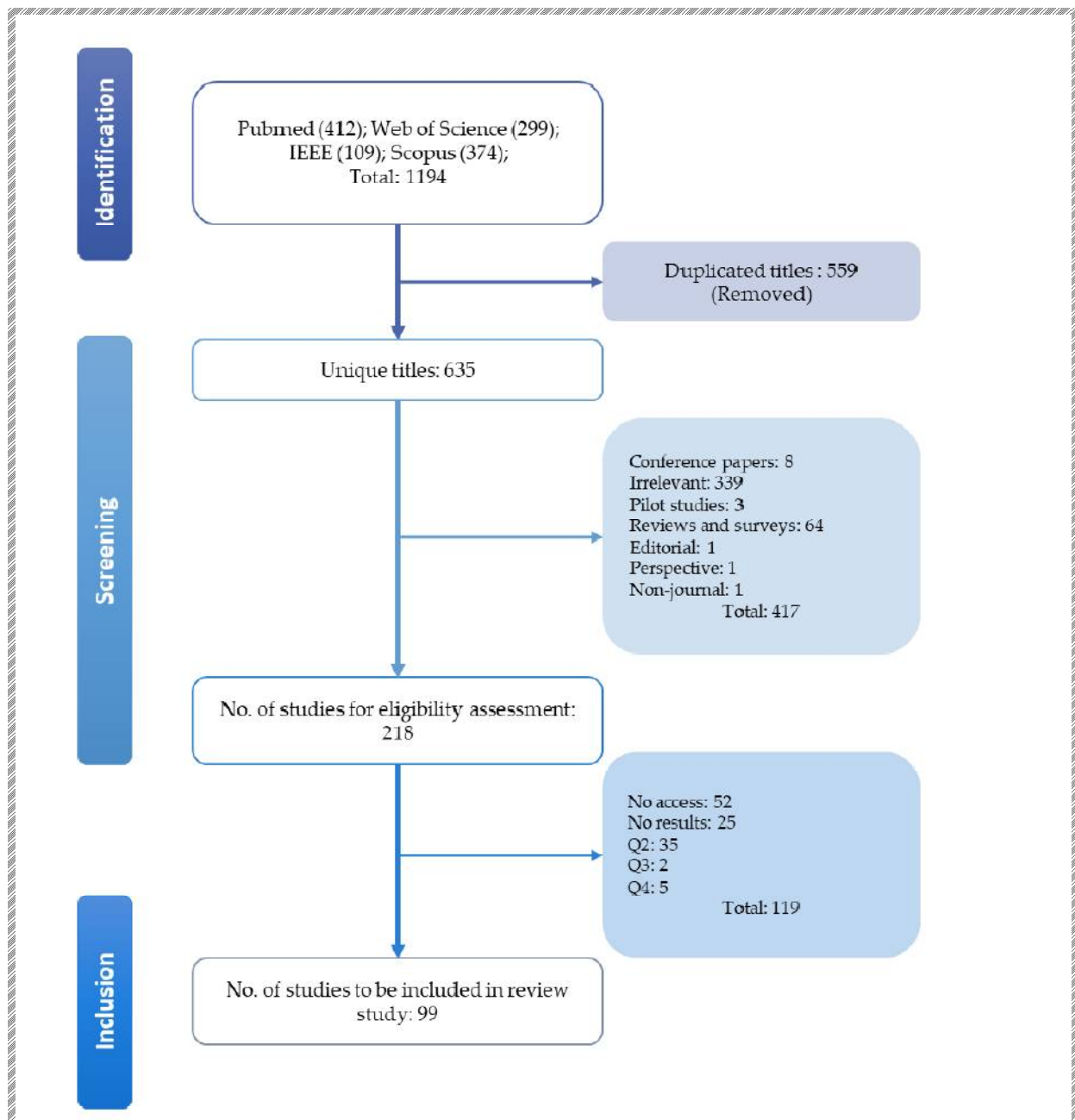




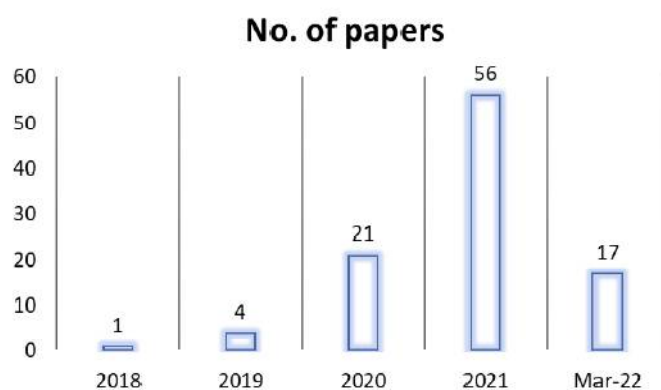
Keyword Search

xAI.		2023-150
Boolean search strings used for journal databases		
Database	Boolean search strings	
Scopus	(TITLE-ABS (explainability) OR TITLE-ABS (explainable)) AND (TITLE-ABS (machine learning) OR TITLE-ABS (deep learning)) AND (TITLE-ABS (medical) OR TITLE-ABS (clinical) OR TITLE-ABS (health) OR TITLE-ABS (healthcare) OR TITLE-ABS (biomedical)) AND (LIMIT-TO (DOCTYPE, "ar"))	
PubMed	(explainability[Title/Abstract] OR explainable[Title/Abstract]) AND (machine learning[Title/Abstract] OR deep learning [Title/Abstract]) AND (medical[Title/Abstract] OR biomedical[Title/Abstract] OR health[Title/Abstract] OR healthcare[Title/Abstract] OR clinical[Title/Abstract])	
Web of Science	(AB-(explainable) OR AB-(explainability)) AND (AB-(machine learning) OR AB-(deep learning)) AND (AB-(medical) OR AB-(biomedical) OR AB-(health) OR AB-(healthcare) OR AB-(clinical))	
IEFF	(("Abstract": "explainable" OR "Abstract": "explainability") AND ("Abstract": "machine learning" OR "Abstract": "deep learning"))	

xAI.		2023-150
PRISMA systematic filtration of journal articles		

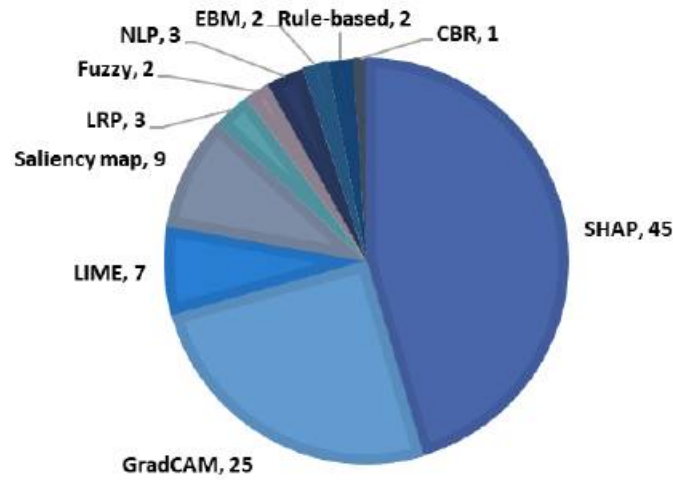


Bar chart representation of XAI studies from 2018 to March 2022

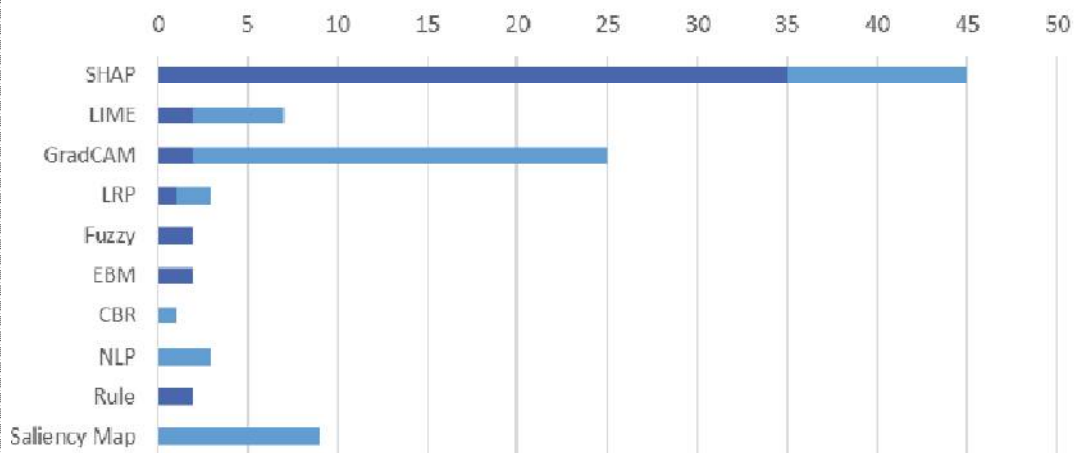


Pie chart of XAI techniques employed

IMPLEMENTED XAI METHOD IN HEALTHCARE APPLICATION



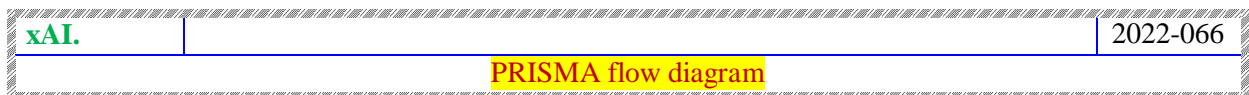
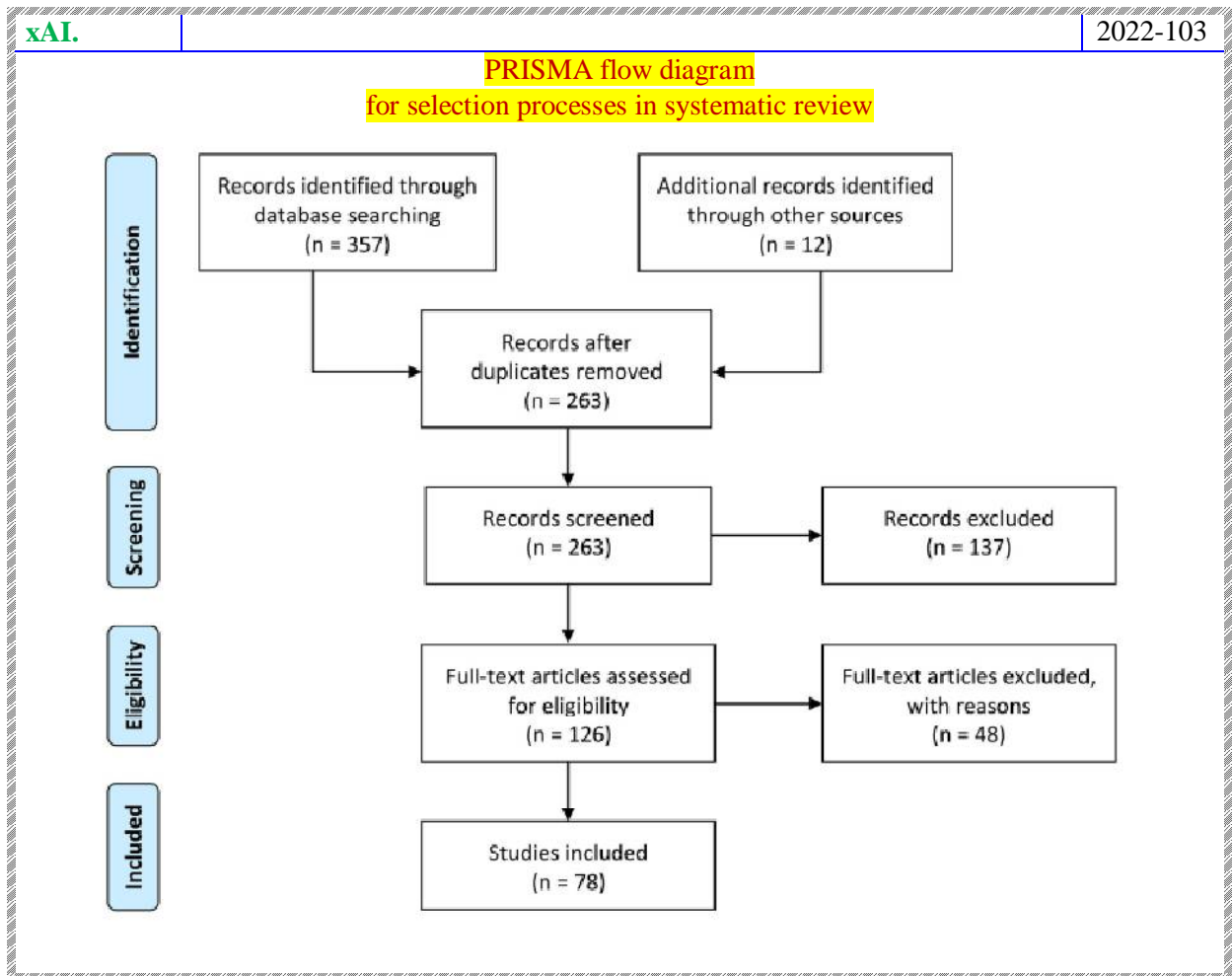
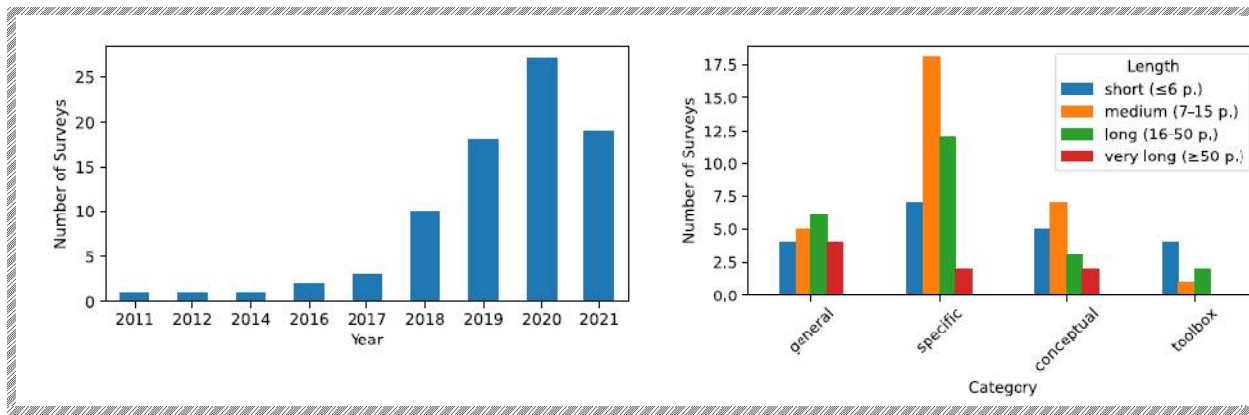
Bar chart representation of the number of ML and DL models used in conjunction with each XAI technique

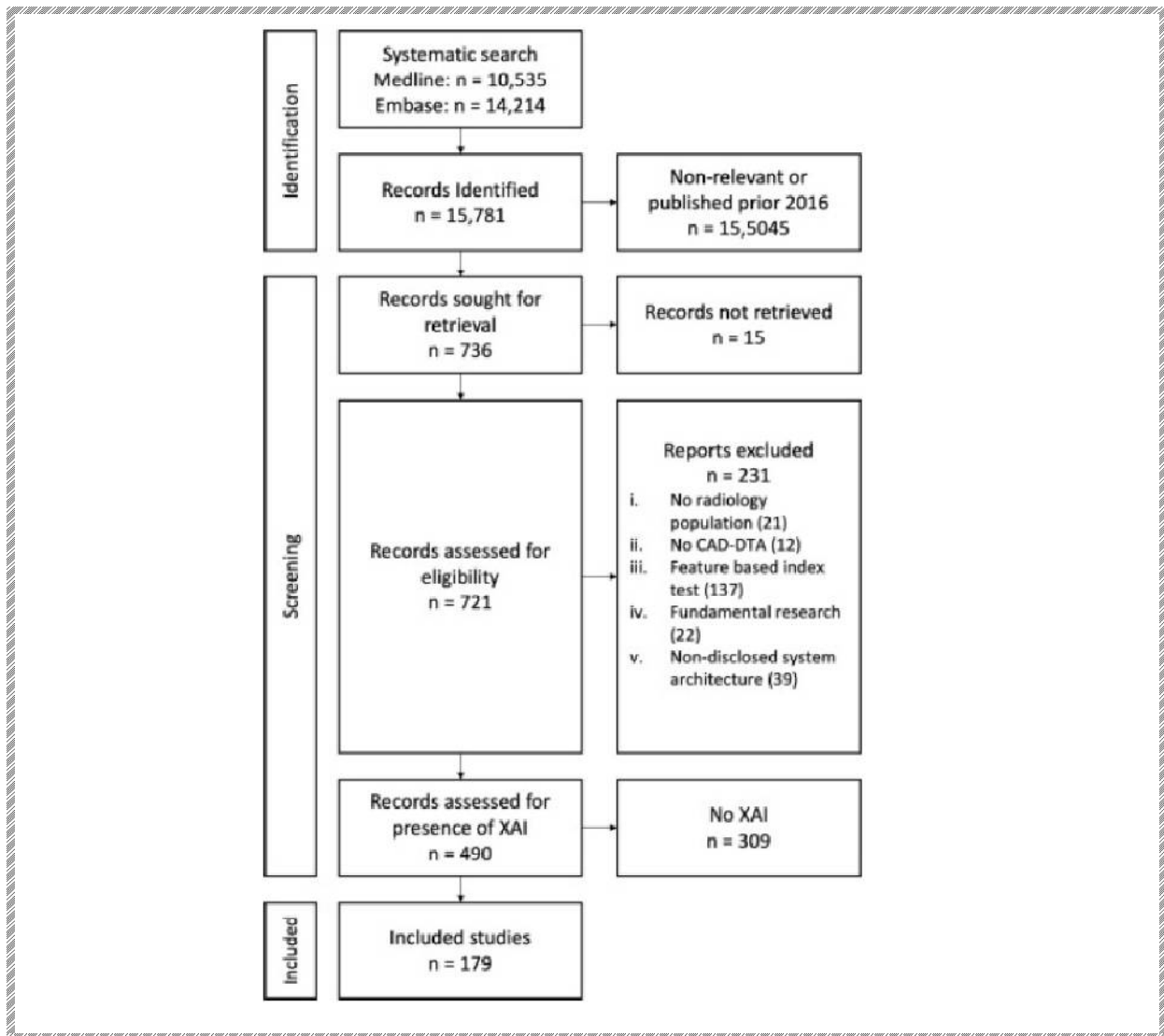


	Saliency Map	Rule	NLP	CBR	EBM	Fuzzy	LRP	GradCAM	LIME	SHAP
ML	0	2	0	0	2	2	1	2	2	35
DL	9	0	3	1	0	0	2	23	5	10

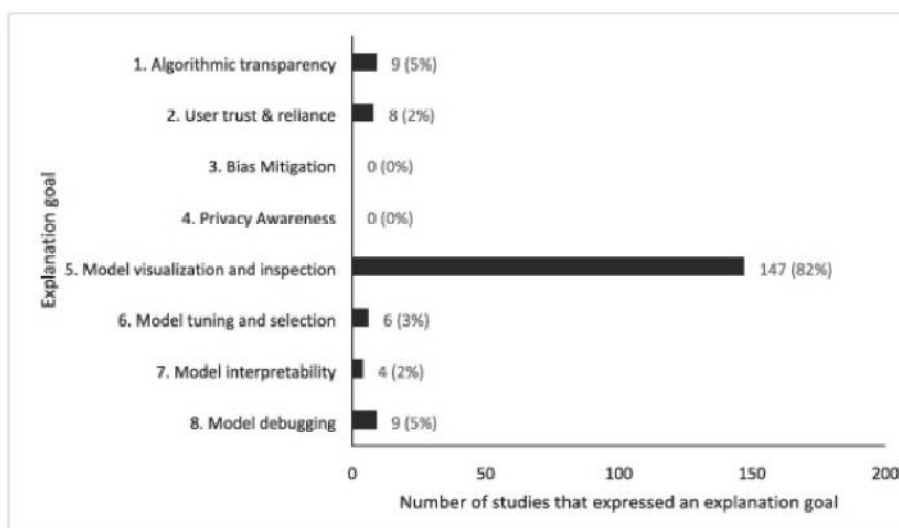
Distribution of the reviewed surveys over years 2010 to 2021

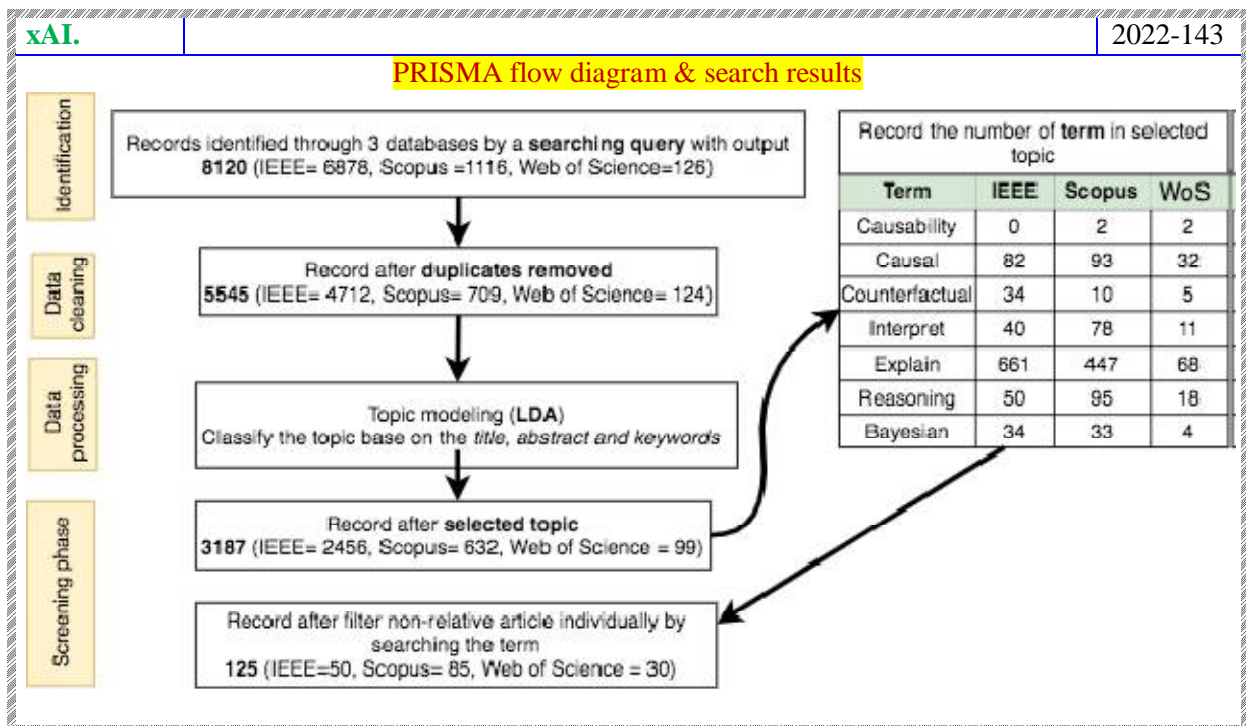
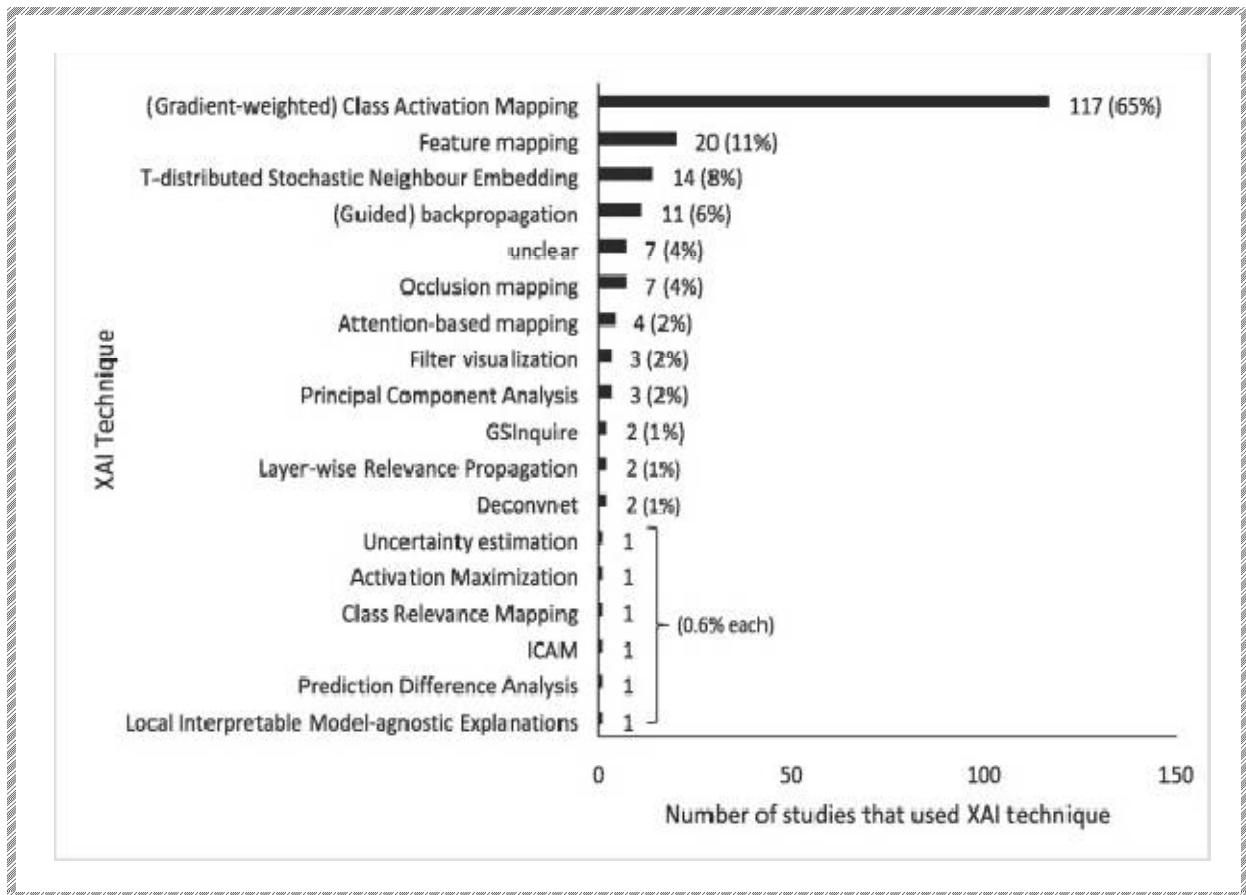
Distribution of the survey Lengths by general focus category

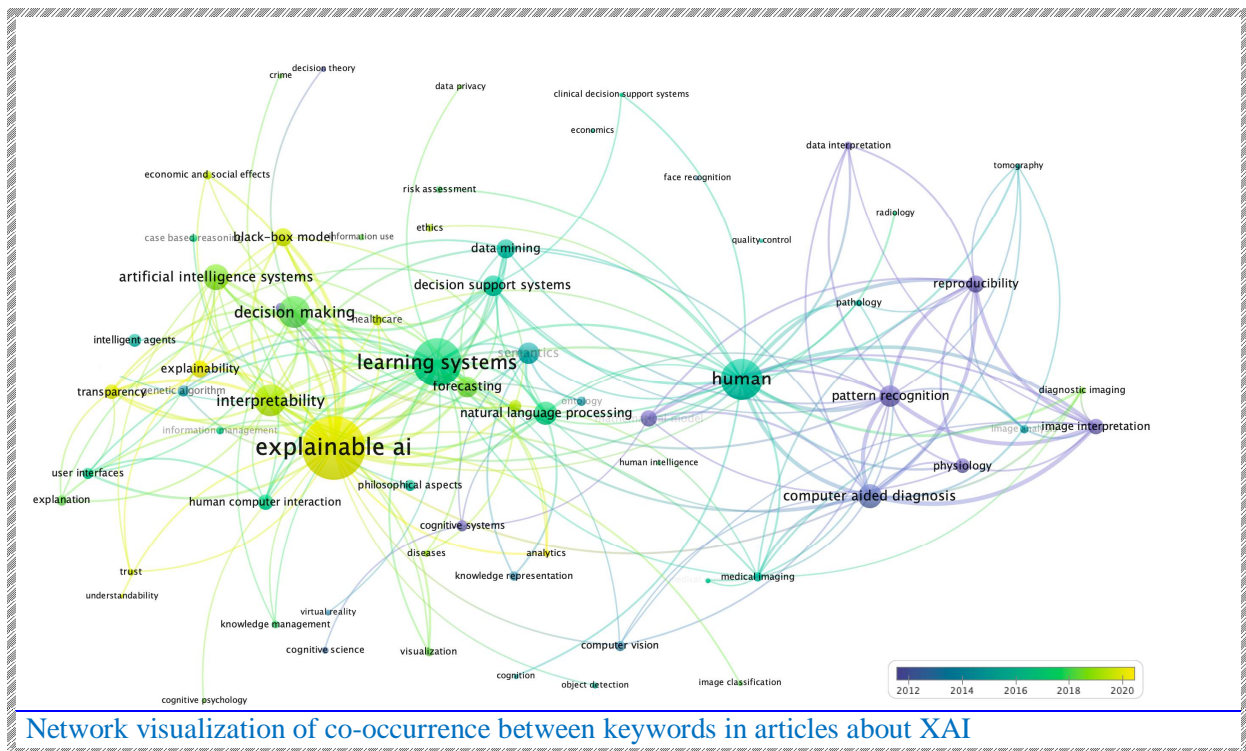
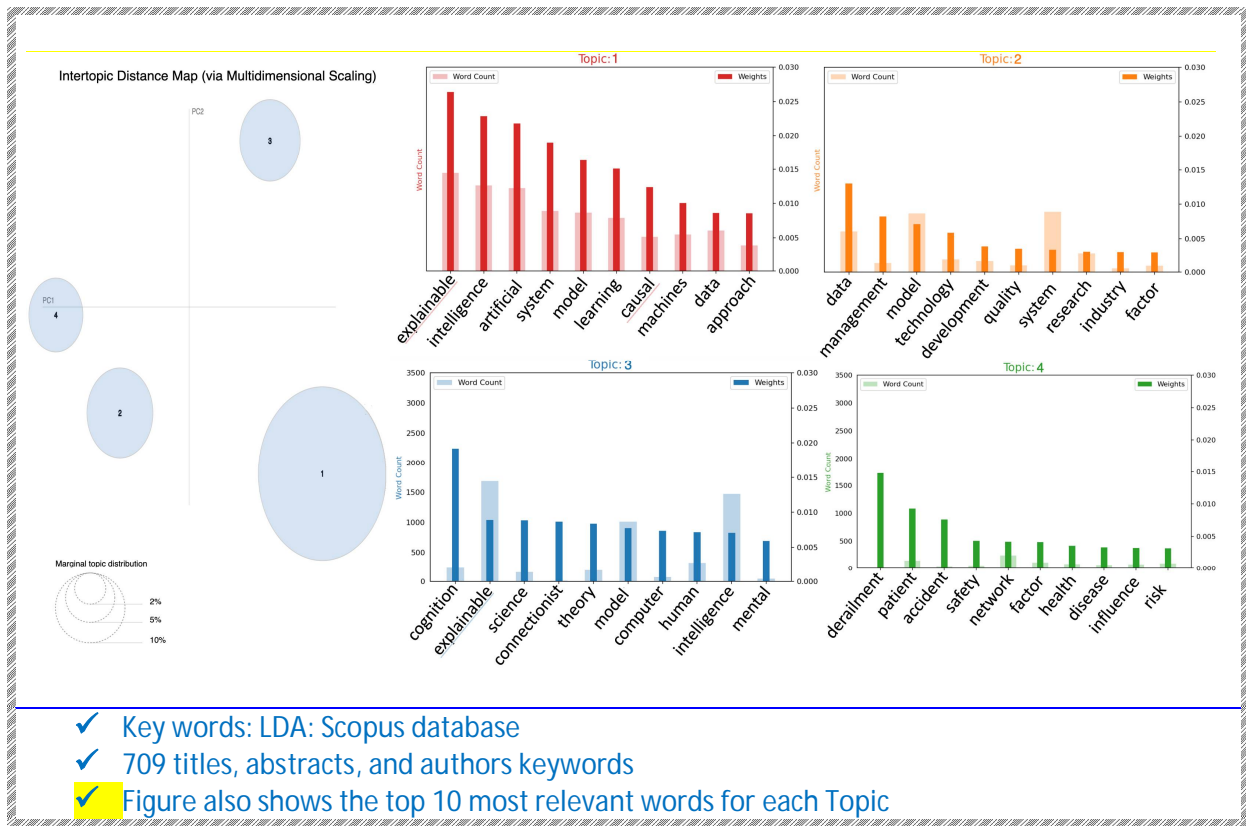


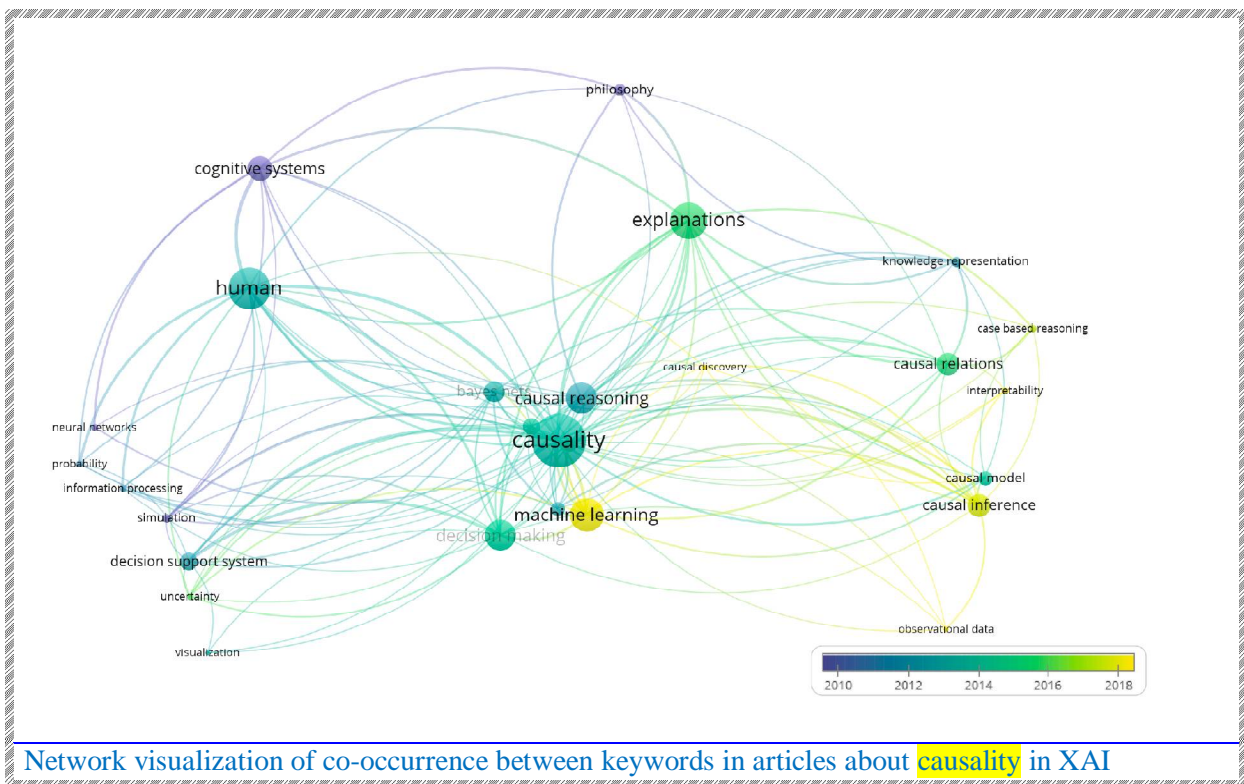
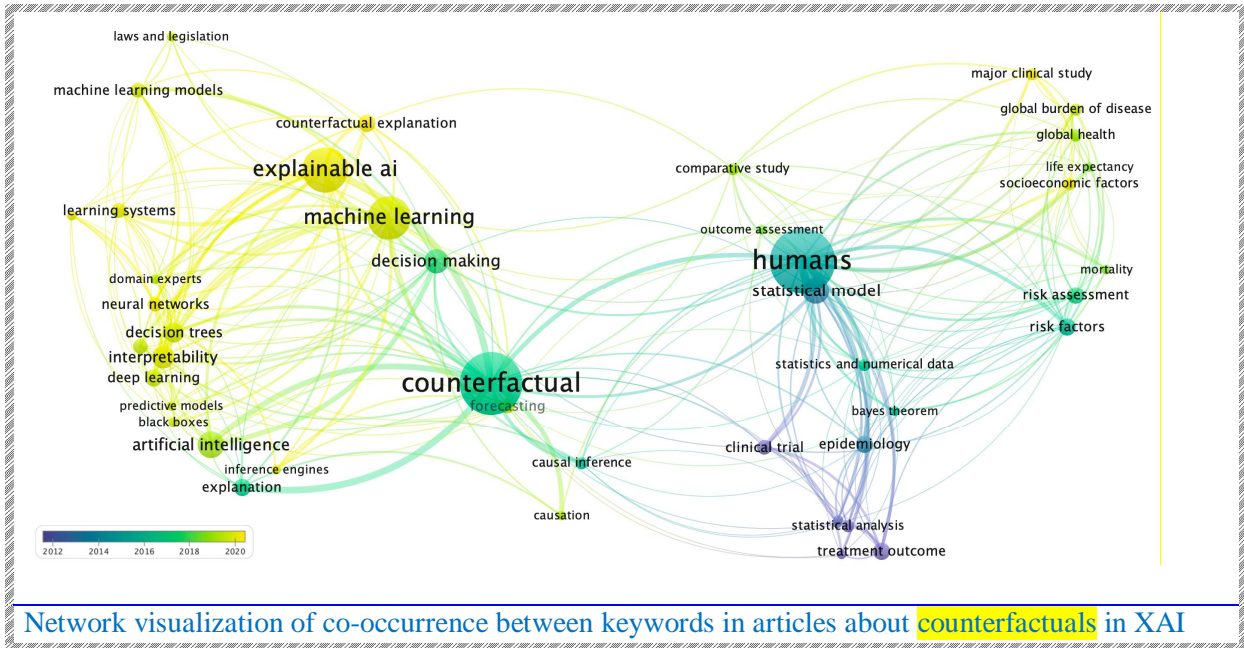


Explanation goals

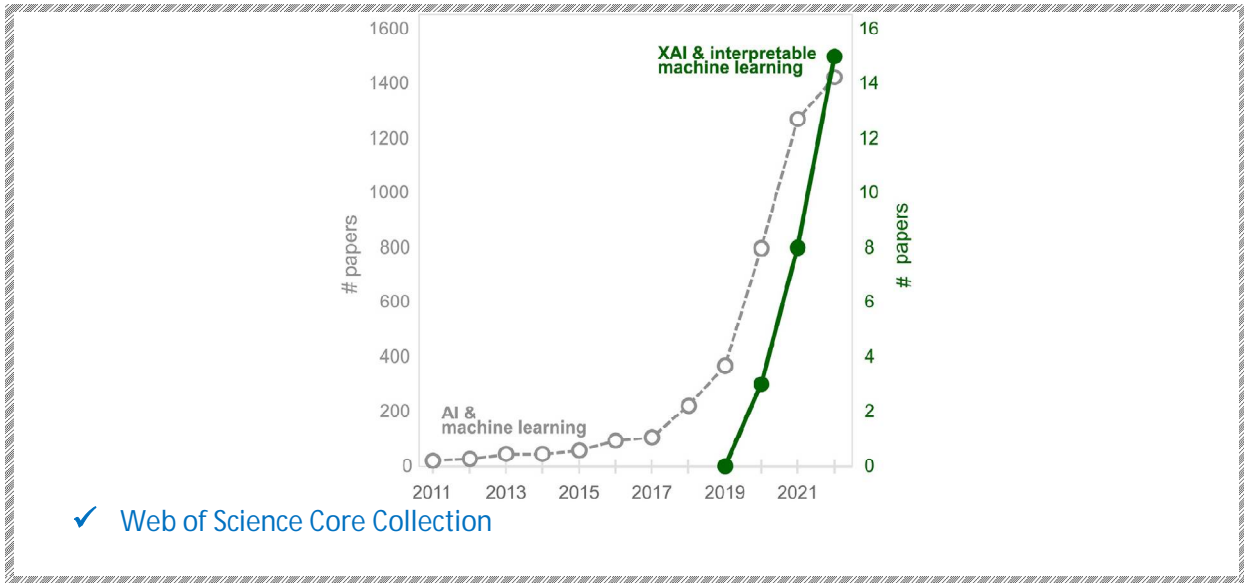




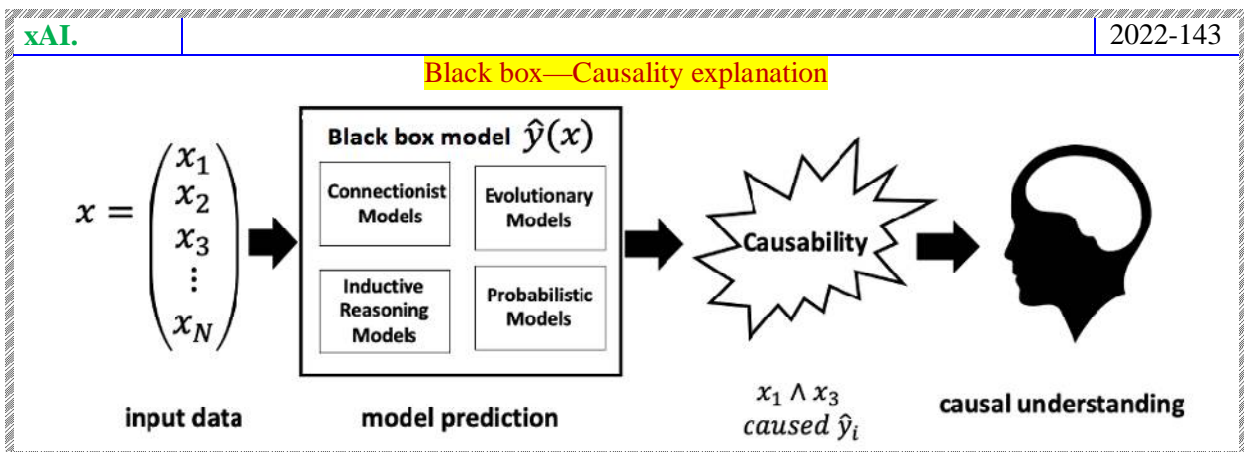
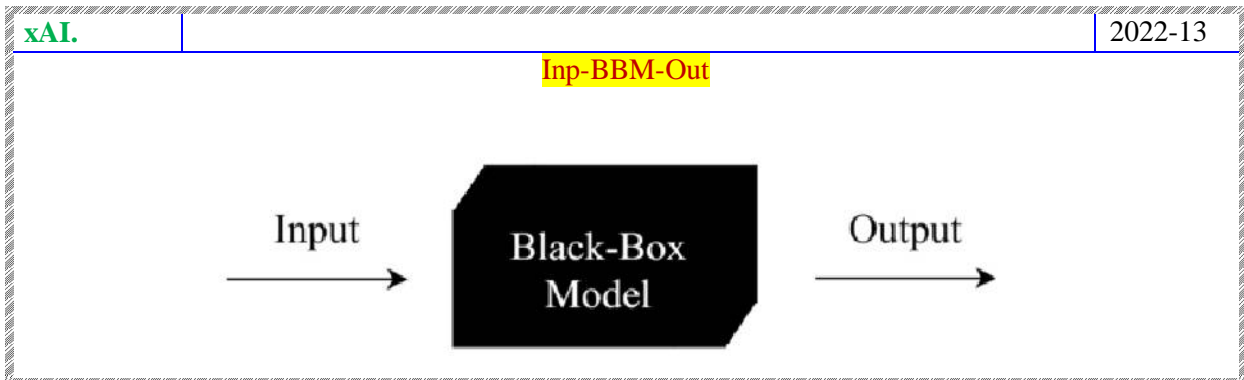




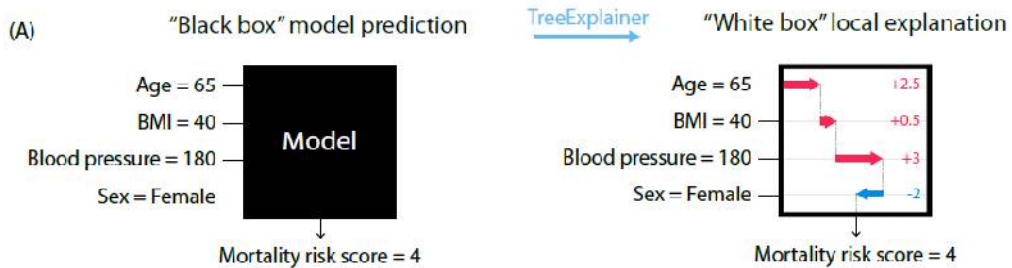
xAI.		2022-153
<p>“AI and machine learning” “XAI and interpretable machine learning” in agricultural science</p>		



Black--To--white box *through* grey box AI methods



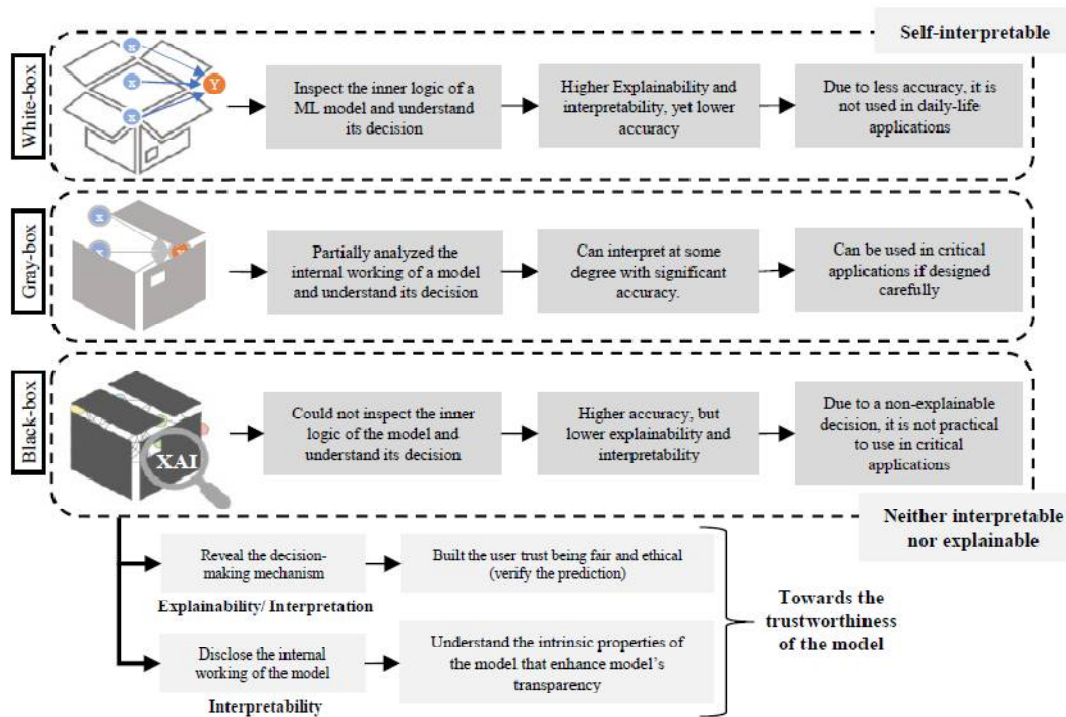
Black—White box Methods



xAI.

2023-030

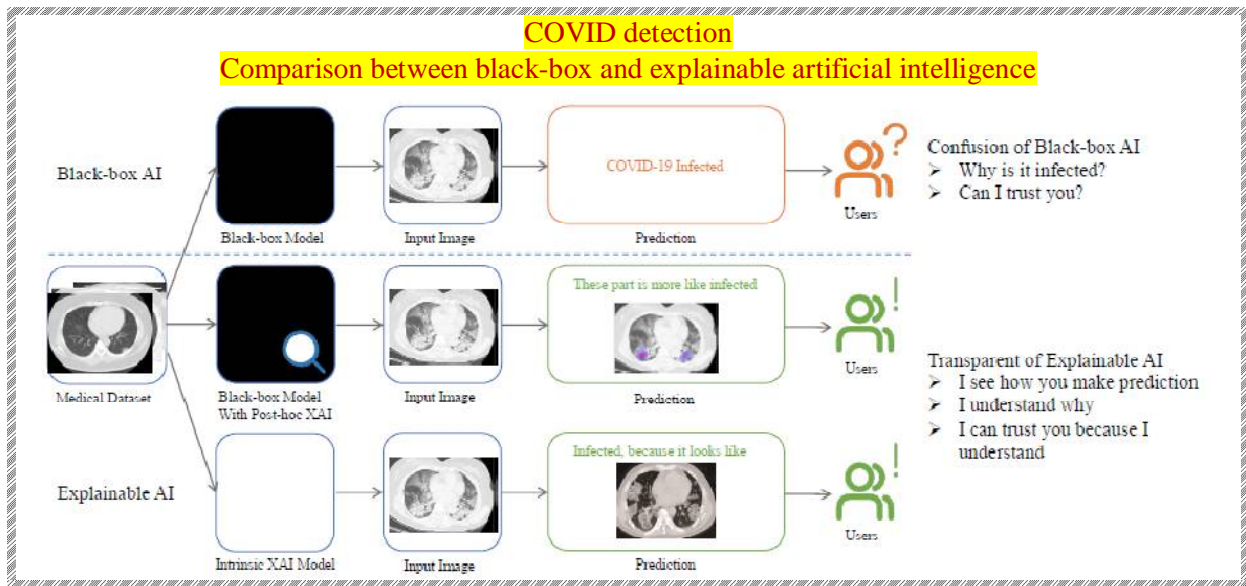
White-box, gray-box, black-box models



- 🔔 Black-box models are more accurate but less interpretable
- 🔔 White-box models are interpretable by design. → making their outputs easier to understand but less accurate.
- 🔔 Gray-box models yield a good interpretability/accuracy tradeoff
- 🔔 Future: More complex XAI techniques are required for creating trustworthy models.

xAI.

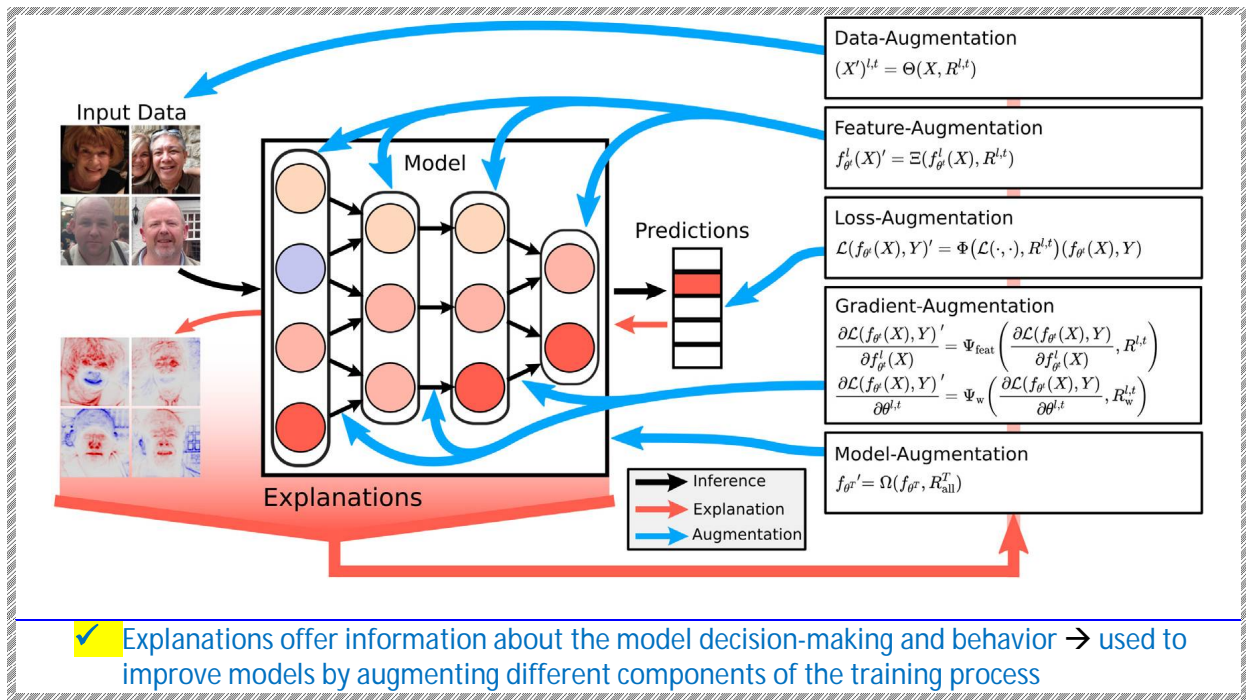
2023-



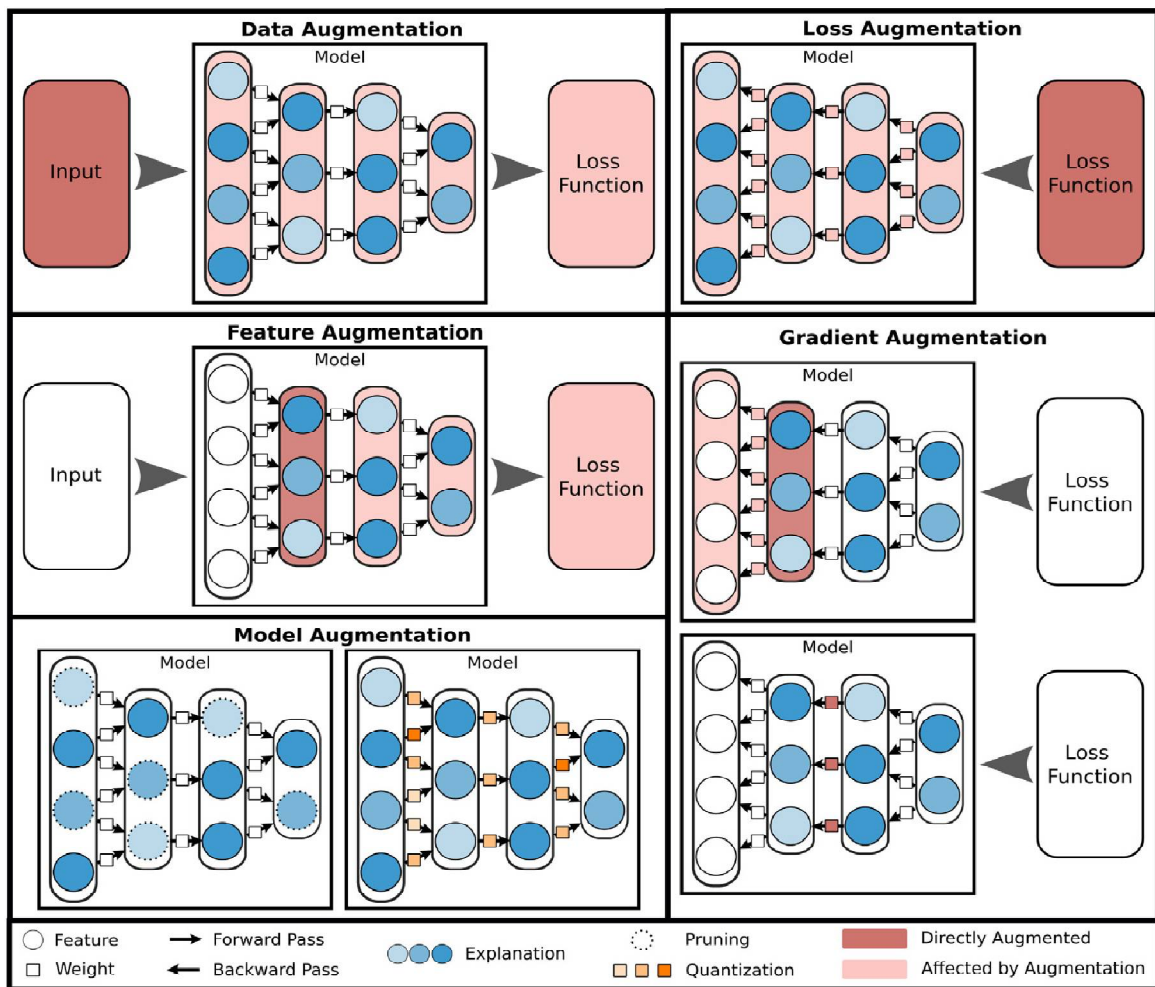
- Top branch shows the process of a black-box model. It provides only results such as classes (e.g., COVID or non-COVID).
- Middle and bottom branches: Two XAI methods;
- Middle branch: Example of saliency map
- Bottom branch: Prototype method

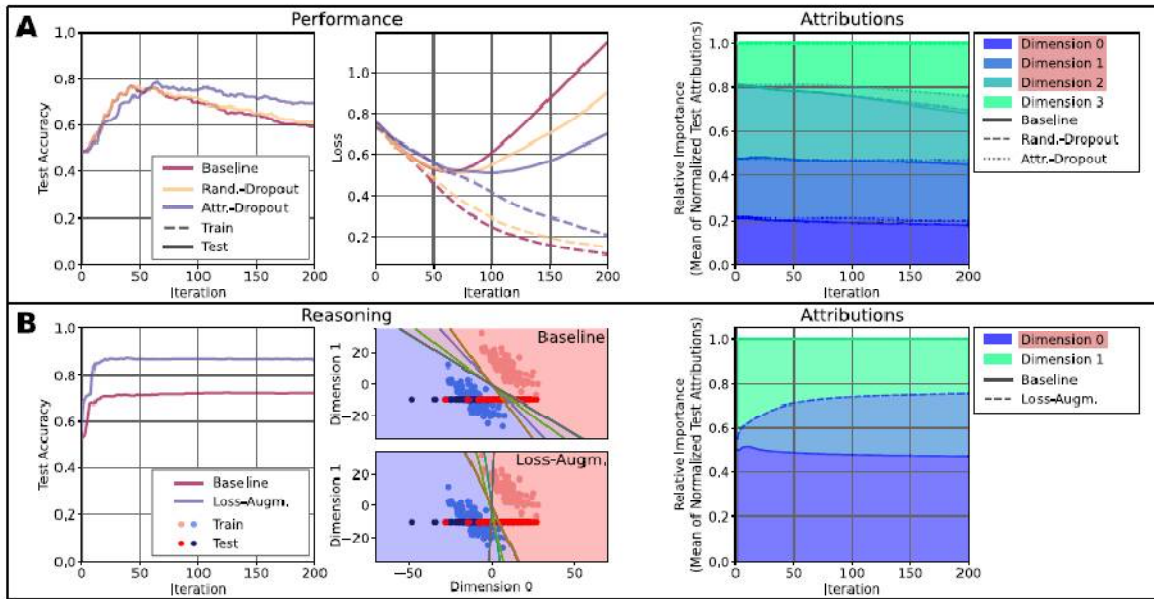
xAI.		Black-box to white box xAI			2023-
Explanation Type	Black-Box Model	Method	Scope	Functionality	
Feature Importance	Any	LIME	Local	Surrogate Model	
		LORE	Local	Surrogate Model	
		Anchors	Local	Surrogate Model	
		Occlusion	Local	Input Perturbation	
		Permutation Feature Importance	Global	Input Perturbation	
		Shapley Feature Importance	Global	Game Theory	
		SHAP	Both	Game Theory	
	Neural Network	Guided Backpropagation	Local	Backpropagation	
		Integrated Gradients	Local	Backpropagation	
		Layerwise Relevance Propagation	Local	Backpropagation	
CNN	DeepLift	Local	Backpropagation		
	Testing with Concept Activation Vectors	Global	Human Concepts		
	Activation Maximization	Global	Forward propagation		
Transformer	Deconvolution	Local	Backpropagation		
	Class Activation Map	Local	Backpropagation		
	Grad-CAM	Local	Backpropagation		
White-Box Model	Any	Attention Flow / Attention Rollout	Local	Network Graph	
		Transformer Relevance Propagation	Local	Backpropagation	
		Rule Extraction	Global	Simplification	
	CNN	Tree Extraction	Global	Simplification	
Example Based	Any	Model Distillation	Global	Simplification	
		Attention Network	Global	Model Adaption	
		Attention Network	Global	Model Adaption	
Visual Explanations	Any	Prototypes	Global	Example (Train Data)	
		Criticisms	Global	Example (Train Data)	
		Counterfactuals	Global	Fictional data point	
		Partial Dependence Plot	Global	Marginalization	
Visual Explanations	Any	Individual Conditional Expectation	Global	Marginalization	
		Accumulated Local Effects	Global	Accumulation	

xAI.		Model improvement with XAI			2023-079
------	--	----------------------------	--	--	----------



Types of XAI-based augmentation

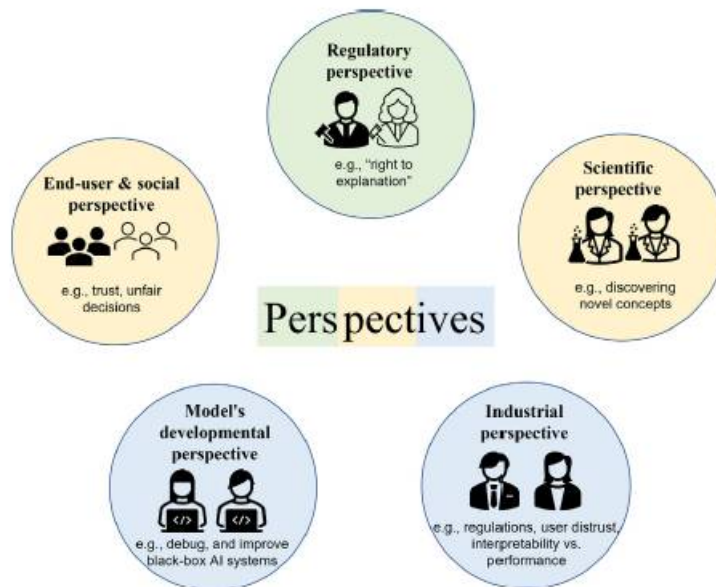




xAI.

2023-123

xAI- perspectives

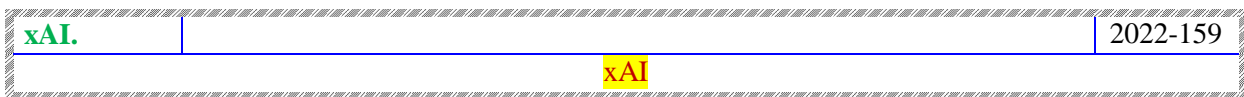
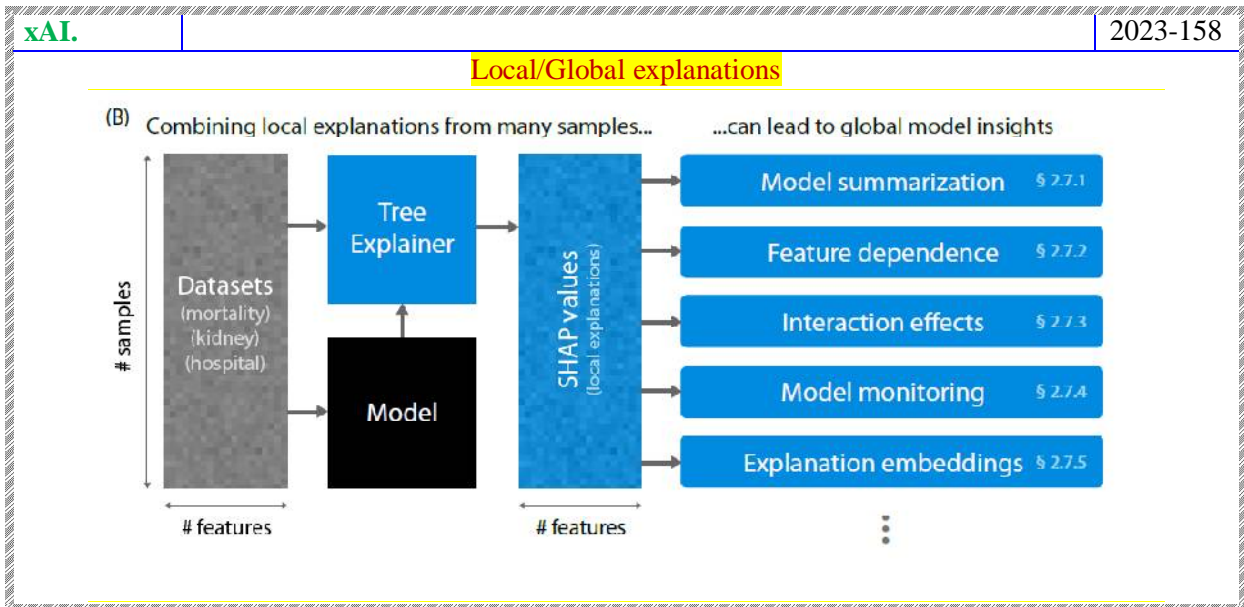
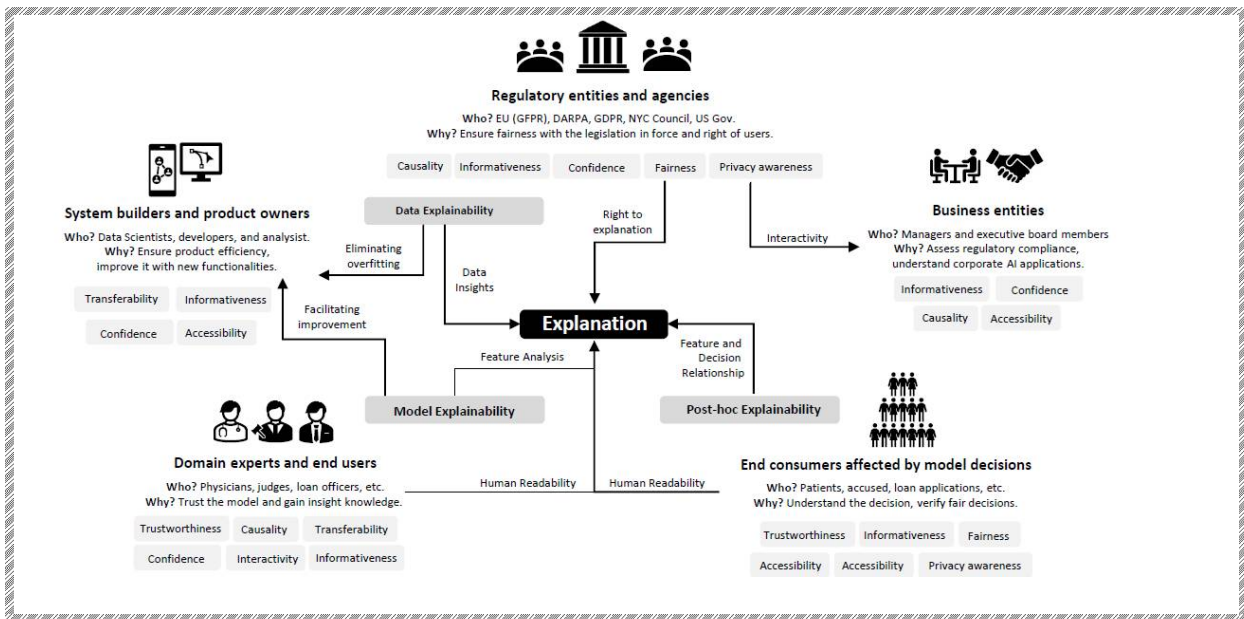


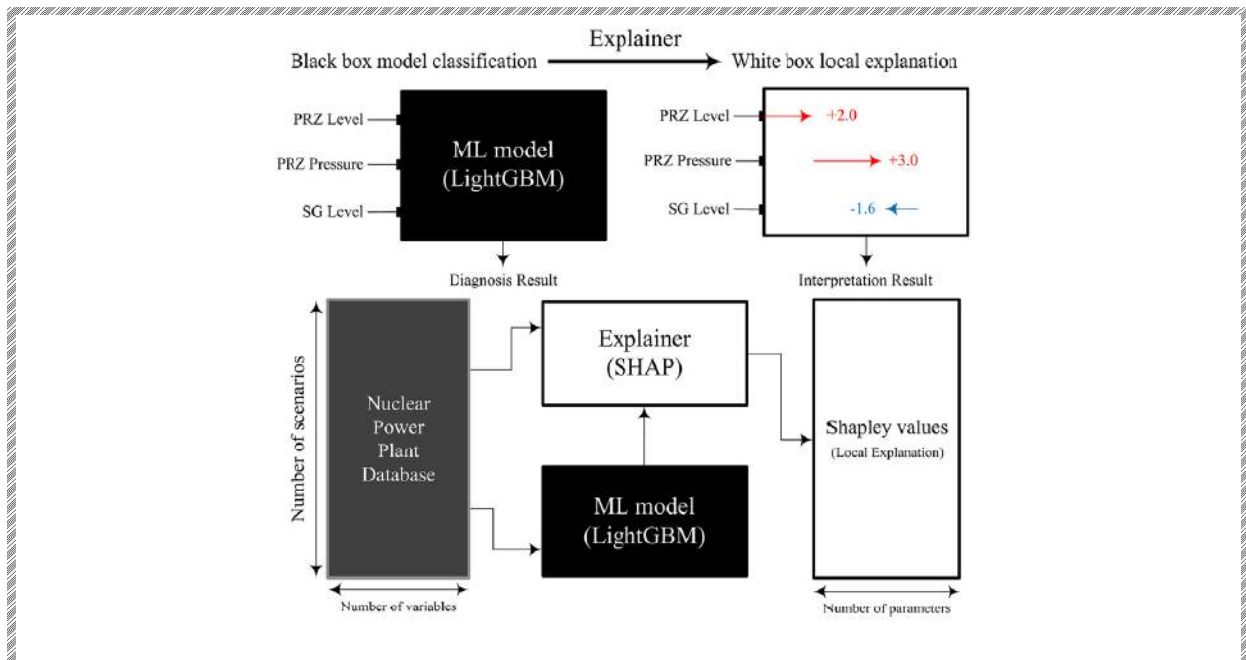
xAI.

2023-030

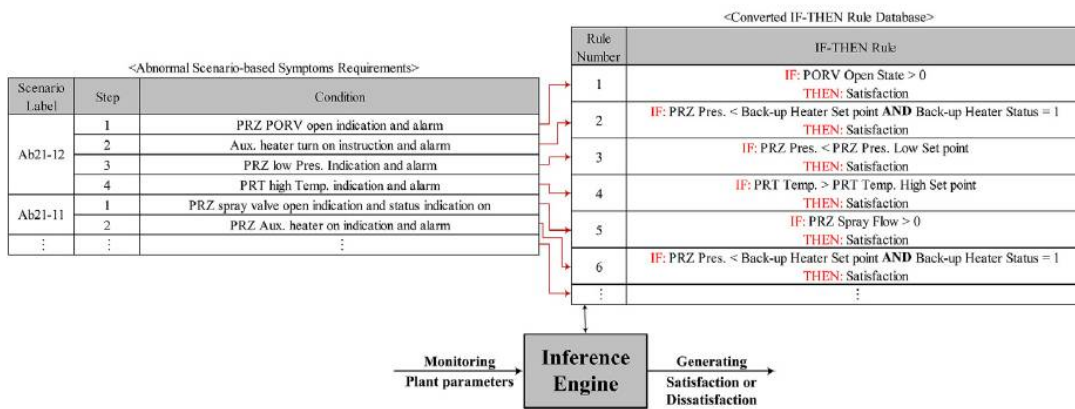
xAI stakeholders

(End-users, domain experts, developers, government bodies)

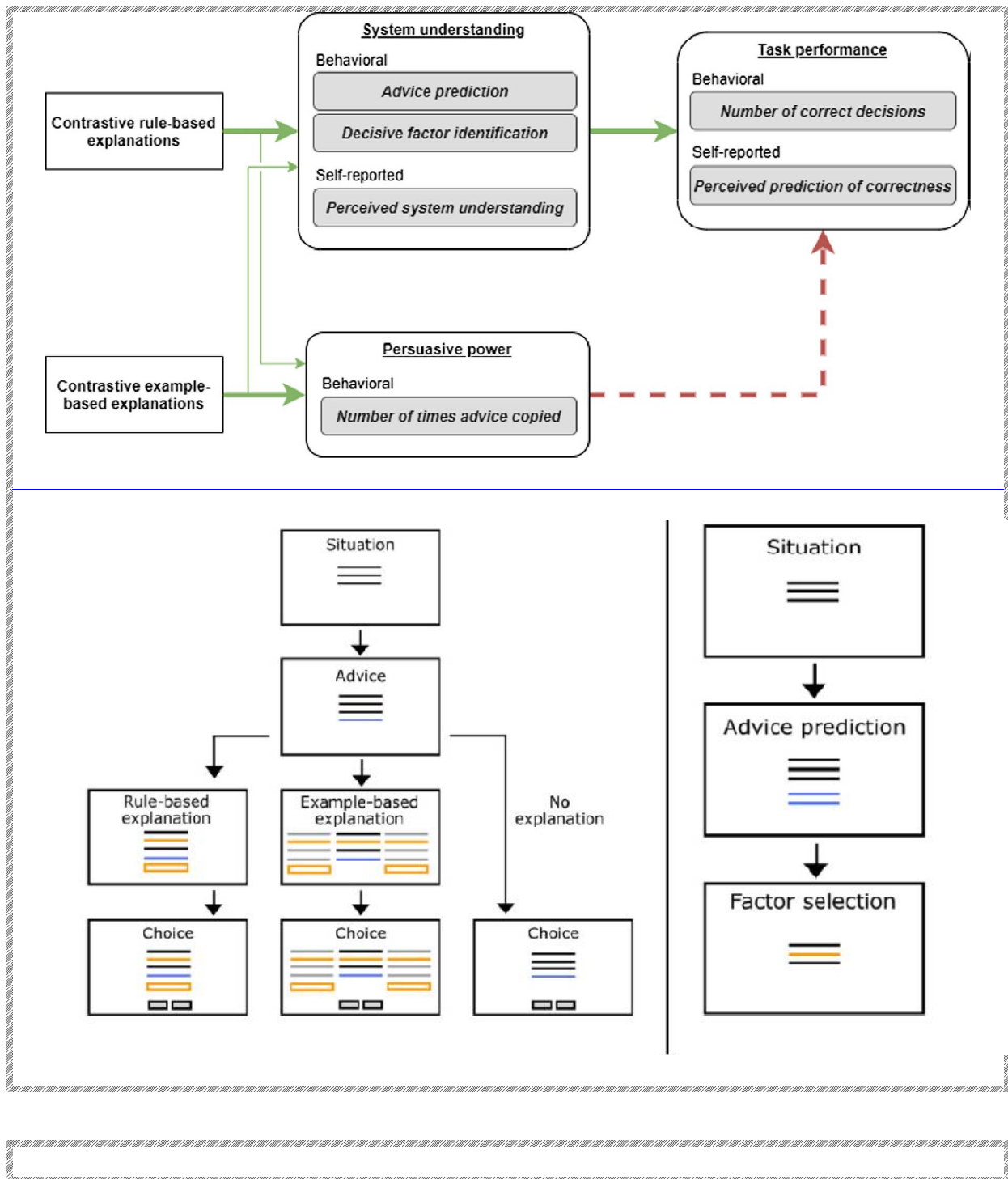


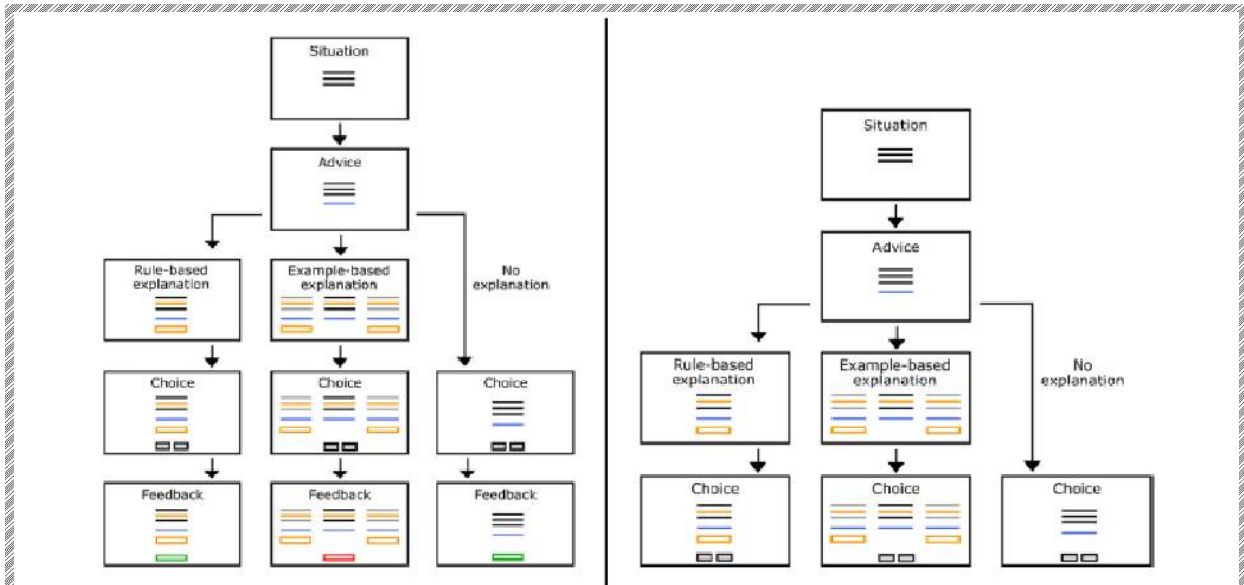


Structure of the rule-based system



Explanation modes





Rule-based explanation mode

Planned alcohol intake
3 units

Water intake so far
5 glasses













Hours slept
6 hours

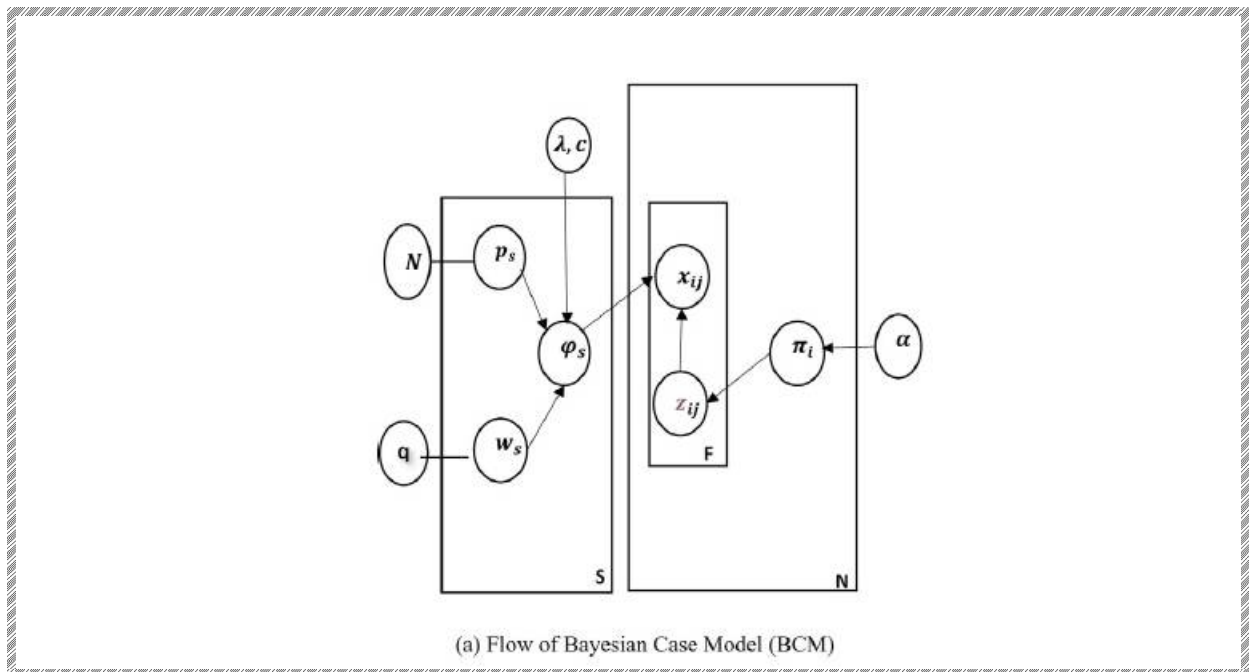
The system advises a
lower dose of insulin

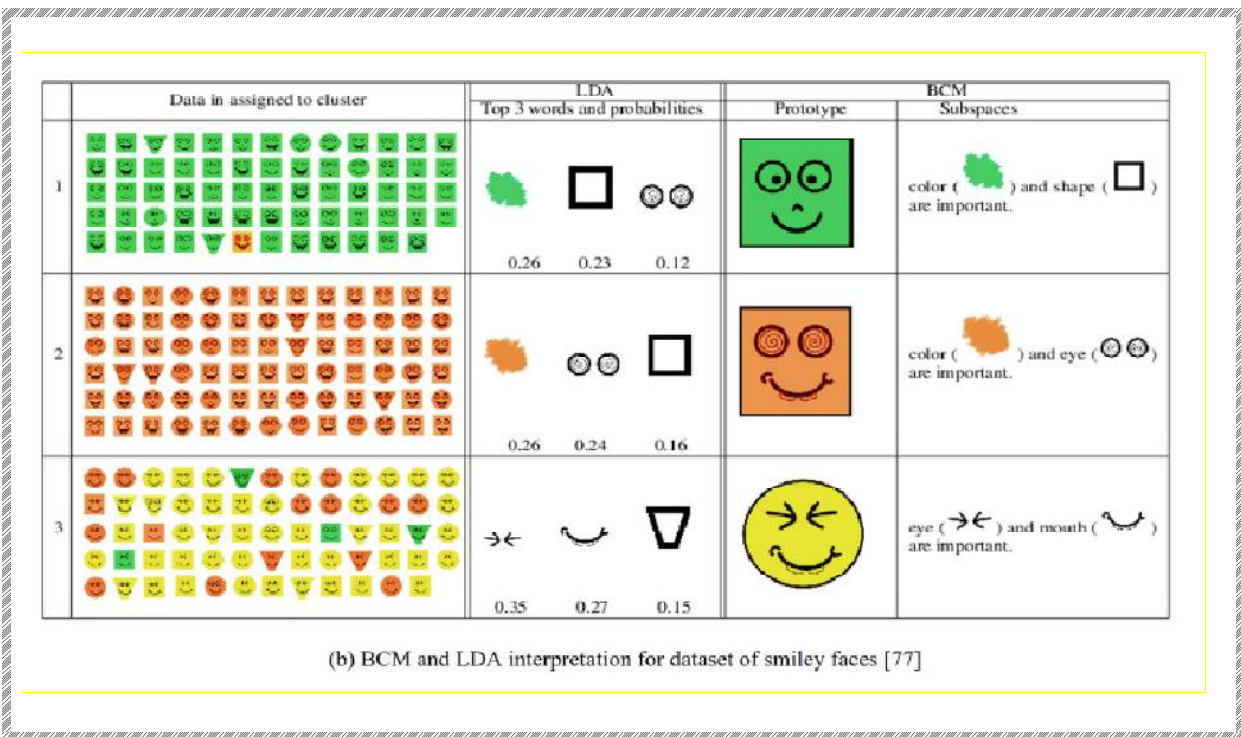
Your planned alcohol intake is more than 1 unit.

If this would have been 1 unit or less, the system would have advised a normal dose.

Example-based explanation styles

Comparable situation from your past	Current situation	Comparable situation from your past
 Planned alcohol intake 3 units	 Planned alcohol intake 3 units	 Planned alcohol intake 1 unit
 Water intake so far 5 glasses	 Water intake so far 5 glasses	 Water intake so far 4 glasses
 Hours slept 7 hours	 Hours slept 6 hours	 Hours slept 6.5 hours
 The system advises a lower dose of insulin	 The system advises a lower dose of insulin	 The system advises a normal dose of insulin
<p><i>Here, your planned alcohol intake was 3 units and the system also advised a lower dose of insulin.</i></p> <p><i>That advice had a positive effect on your blood sugar level.</i></p>		<p><i>Here, your planned alcohol intake was 1 unit and the system advised a normal dose of insulin instead.</i></p> <p><i>That advice had a positive effect on your blood sugar level.</i></p>

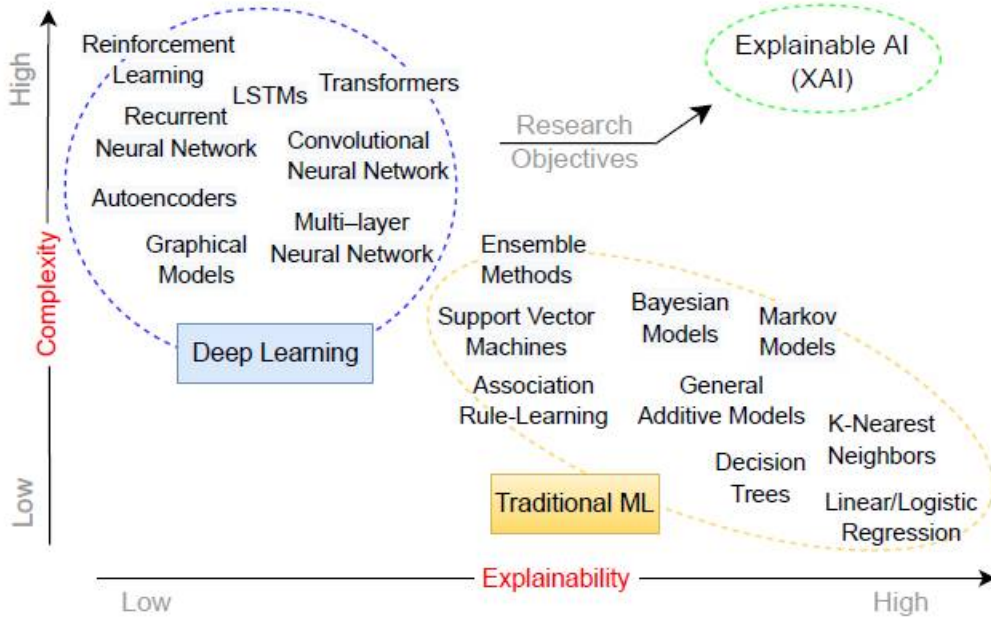




Workflows

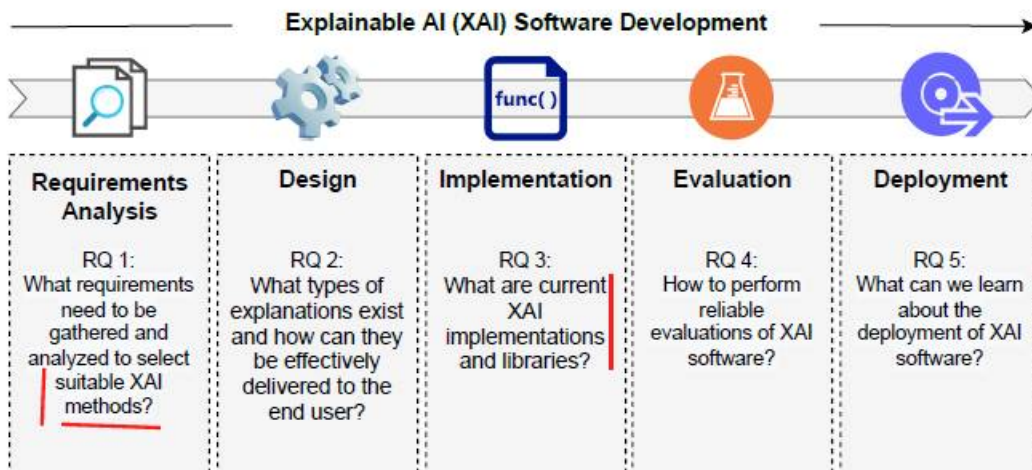
xAI.	Workflow of ECDM-SDAM methodology	2023-144
<div style="display: flex; flex-direction: column; align-items: center;"> <div style="display: flex; align-items: center; margin-bottom: 10px;"> <div style="border: 2px solid black; border-radius: 50%; width: 40px; height: 40px; display: flex; align-items: center; justify-content: center; margin-right: 10px;"> 1 </div> <div style="background-color: #e0f2f1; padding: 10px; border-radius: 10px; width: 70%;"> <p>Obtaining the expert opinions.</p> <ul style="list-style-type: none"> • Opinion extraction with the proposed ASAM model. • Opinion representation with the proposed BOC table. </div> </div> <div style="display: flex; align-items: center; margin-bottom: 10px;"> <div style="border: 2px solid black; border-radius: 50%; width: 40px; height: 40px; display: flex; align-items: center; justify-content: center; margin-right: 10px;"> 2 </div> <div style="background-color: #e0f2f1; padding: 10px; border-radius: 10px; width: 70%;"> <p>Crowd decision making.</p> <ul style="list-style-type: none"> • Collective aggregation. • Exploitation. </div> </div> <div style="display: flex; align-items: center;"> <div style="border: 2px solid black; border-radius: 50%; width: 40px; height: 40px; display: flex; align-items: center; justify-content: center; margin-right: 10px;"> 3 </div> <div style="background-color: #ffe0b2; padding: 10px; border-radius: 10px; width: 70%;"> <p>Explainable backward process.</p> <ul style="list-style-type: none"> • Identifying relevant criteria. • Identifying relevant aspect terms through subgroup discovery. • Identifying relevant sentences through attention mechanisms. </div> </div> </div> <div style="margin-top: 20px;"> <ul style="list-style-type: none"> ○ Model: Crowd Decision Making (CDM) ○ ECDM: Explainable CDM based on ○ Methodology : Subgroup Discovery and Attention Mechanisms (SDAM) ○ ASAM: Attention based Sentiment Analysis Method ○ BOC : Bag of Opinions by Criteria </div>		

Classification of AI models

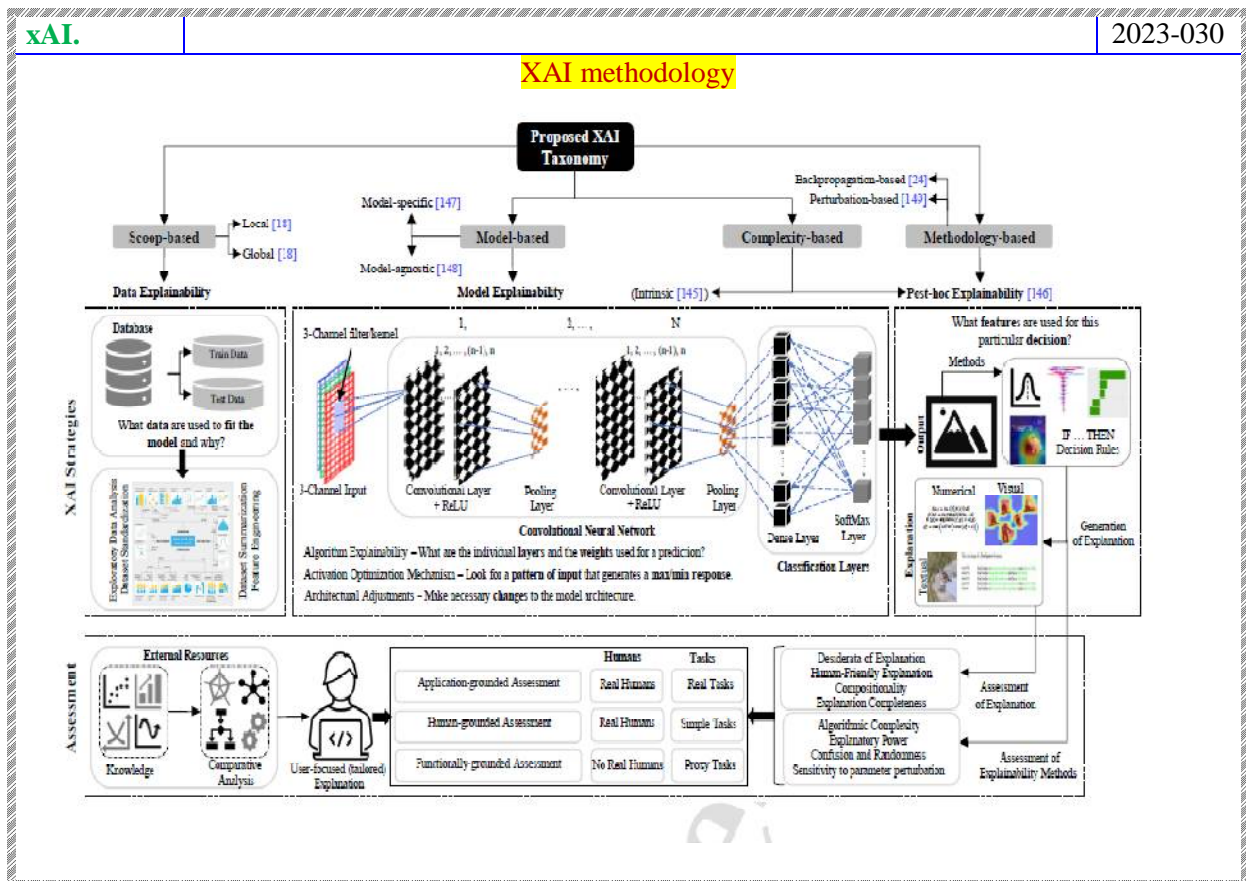
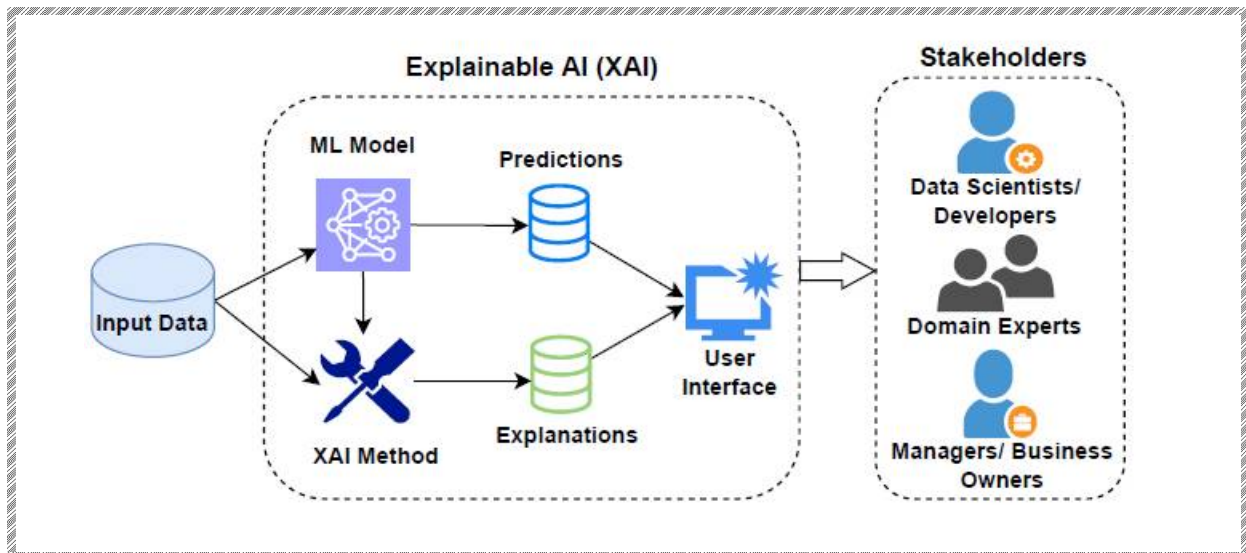


✓ Complexity, explainability, and potential applications

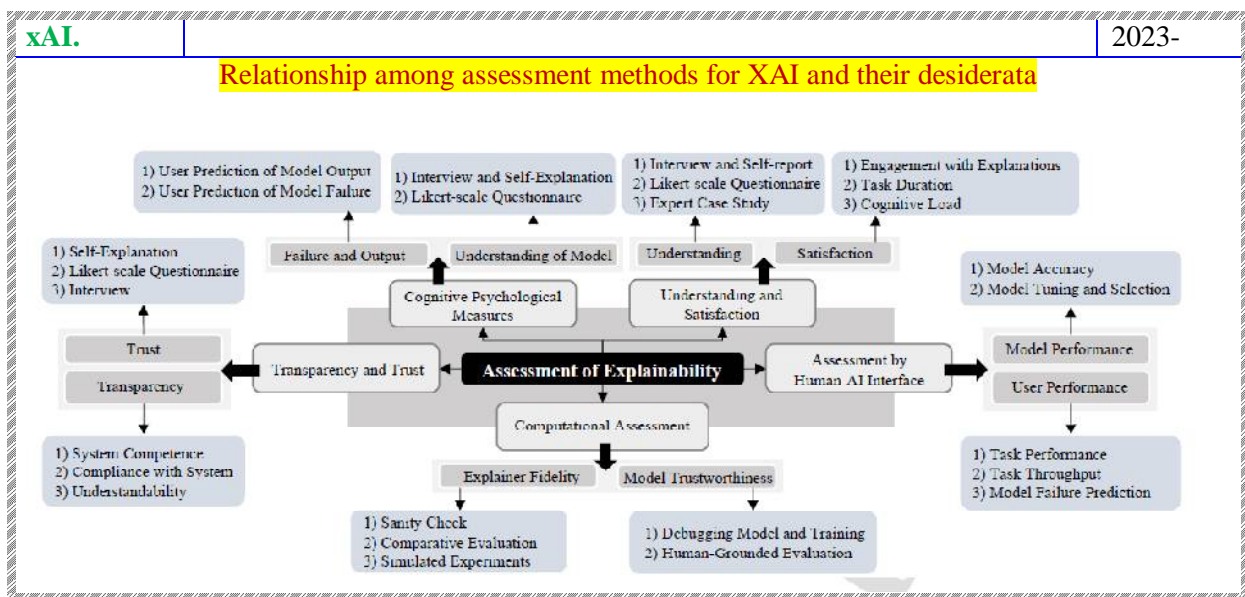
Addressed research questions in XAI software development process



Stakeholders for XAI

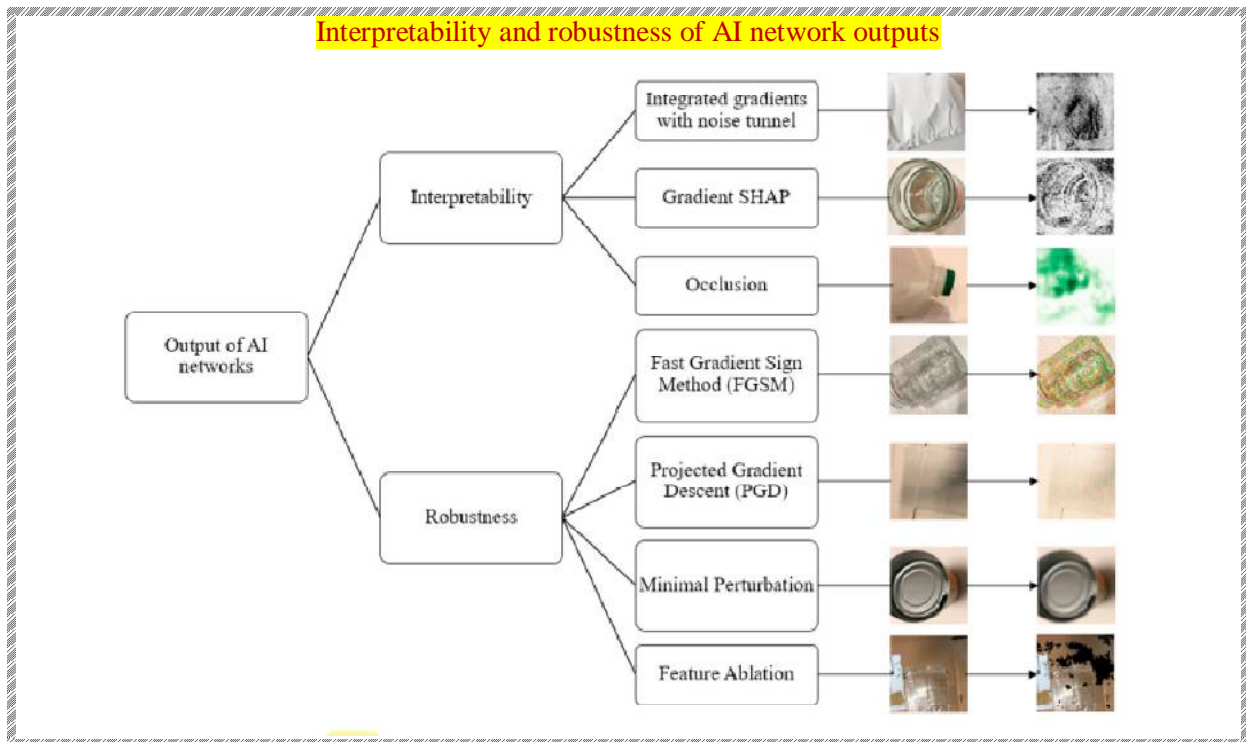
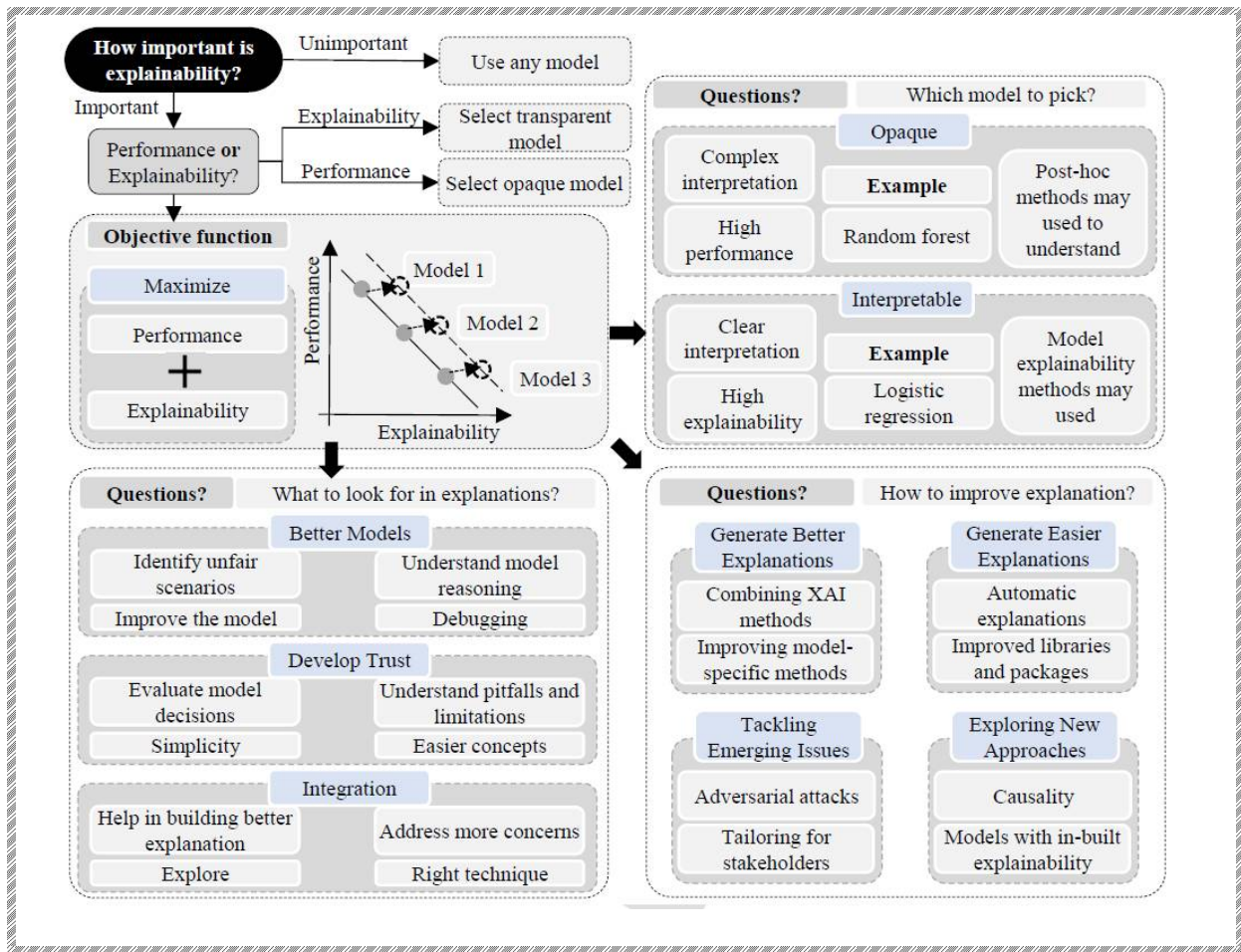


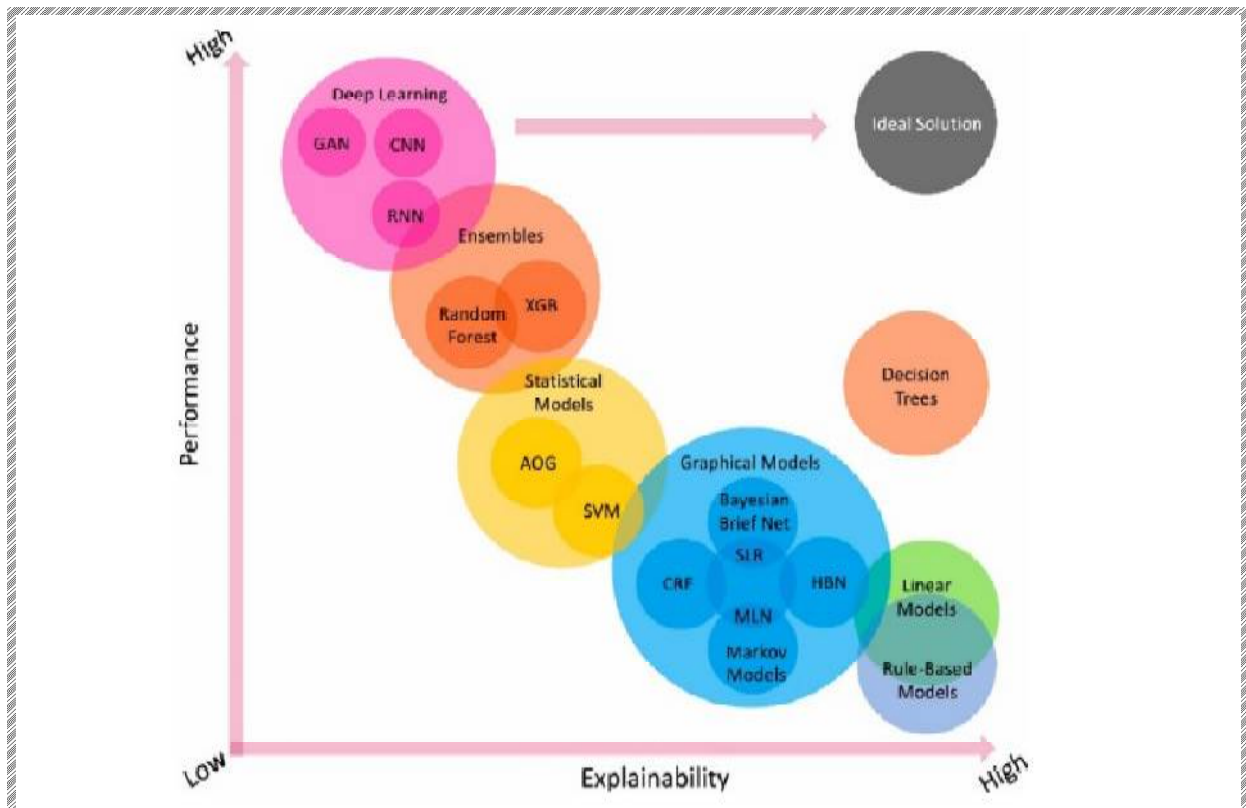
By Data Explainability	By Model Explainability	By Post-hoc Explainability
<p>D1: What sort of information do we have in the database?</p> <p>D2: What can be inferred from this data?</p> <p>D3: What are the most important portions of the data?</p> <p>D4: How is the information distributed?</p> <p>D5: Is it possible to increase the model's performance by lowering the number of dimensions?</p> <p>D6: Can a better explanation be offered by using data summarizing techniques?</p>	<p>M1: What makes a parameter, objective, or action important to the system?</p> <p>M2: When did the system examine a parameter, objective, or action, and when did the model reject it?</p> <p>M3: What are the consequences of making a different decision or adjusting a parameter?</p> <p>M4: How does the system carry out a certain action?</p> <p>M5: How do these model parameters, objectives, or actions relate to one another?</p> <p>M6: What factors does the system take into account (or disregard) when making a decision?</p> <p>M7: In order to achieve a goal/inference, which techniques does the system utilize or avoid?</p>	<p>P1: What is the reason behind the model's prediction?</p> <p>P2: What was the reason for occurrence X? What would happen if Y was the cause of occurrence X?</p> <p>P3: What variables have the most influence on the user's decision?</p> <p>P4: What if the information is altered?</p> <p>P5: To keep current results, what criteria must be met?</p> <p>P6: Is there anything that can be done to have a different outcome?</p> <p>P7: Why is it essential to make a certain conclusion or decision?</p>



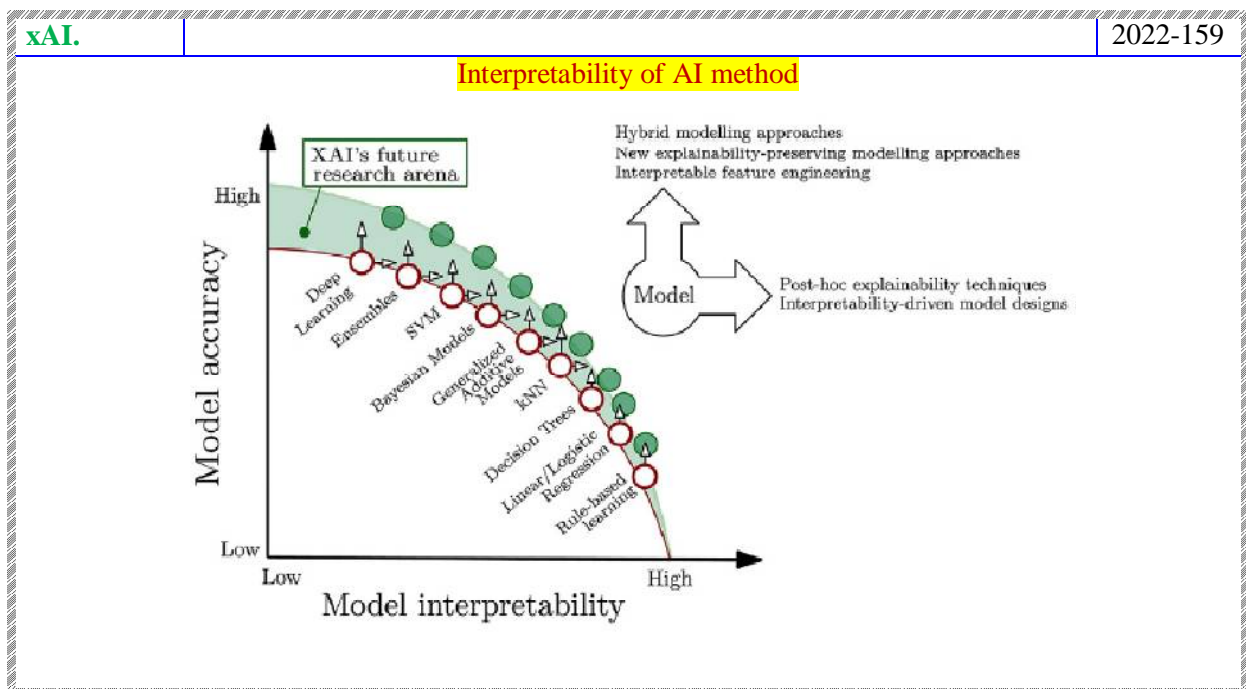
xAI. | 2023-

Application of XAI

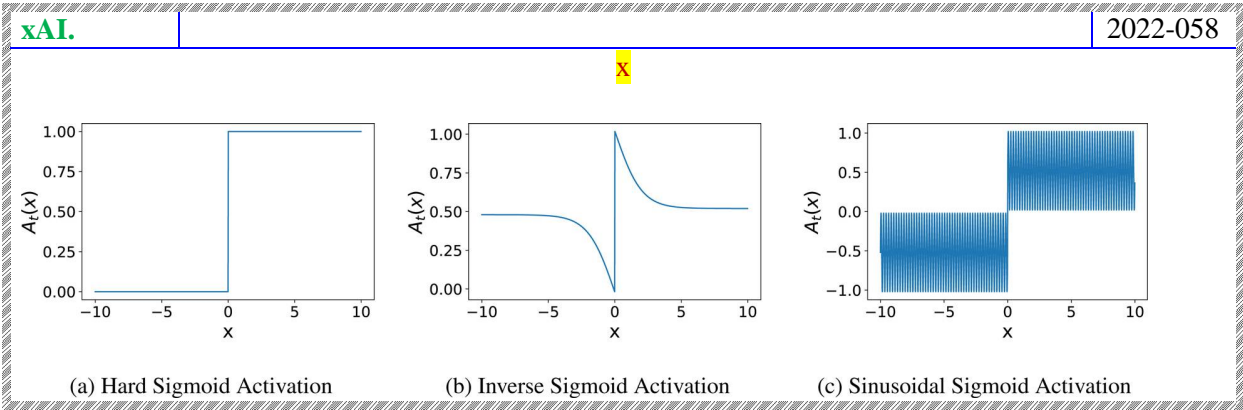




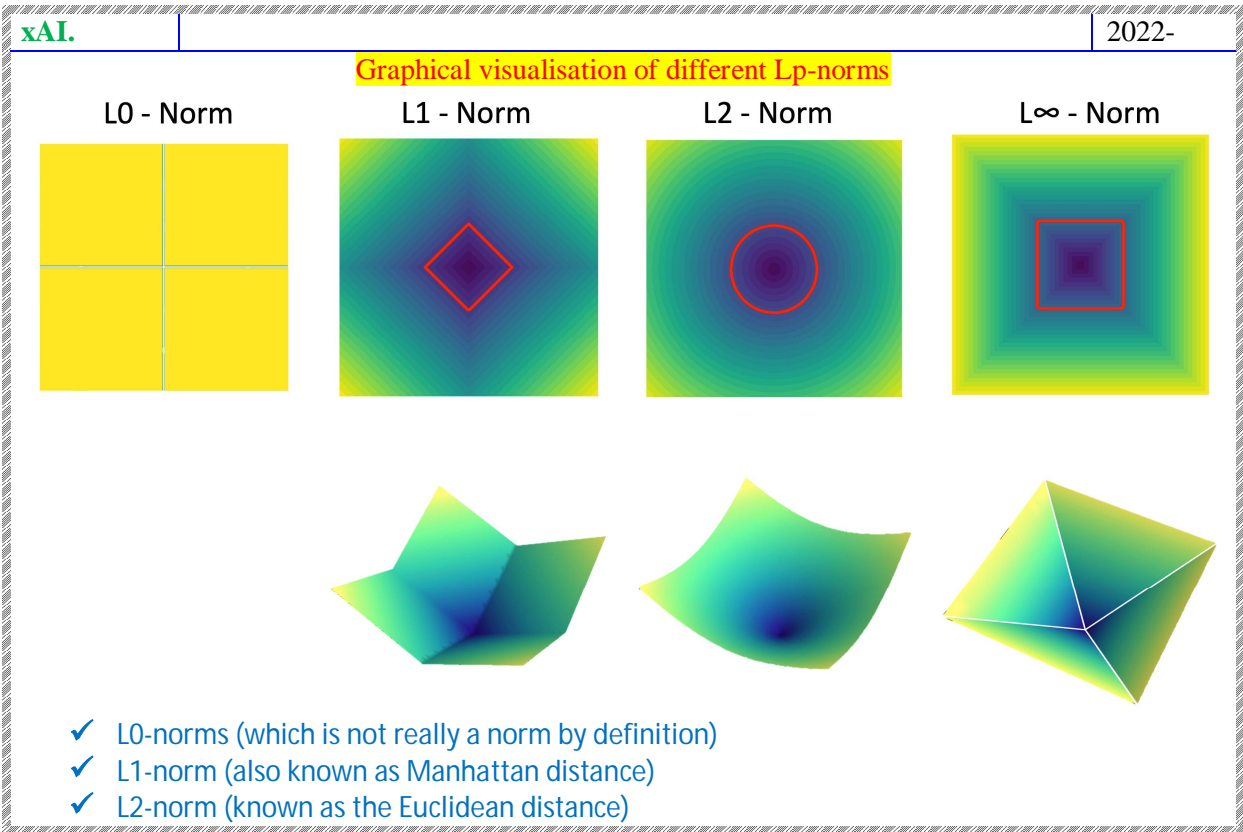
- ✓ HBN: hierarchical Bayesian networks;
- ✓ SLR: simple linear regression;
- ✓ CRF: conditional random fields;
- ✓ MLN: Markov logic network;
- ✓ AOG: stochastic and-or graphs;
- ✓ XGB: extreme gradient boosting;
- ✓ GAN: generative adversarial network;



Transfer functions (TFs)



L0, L1, L2, Linf Norms



Architectures

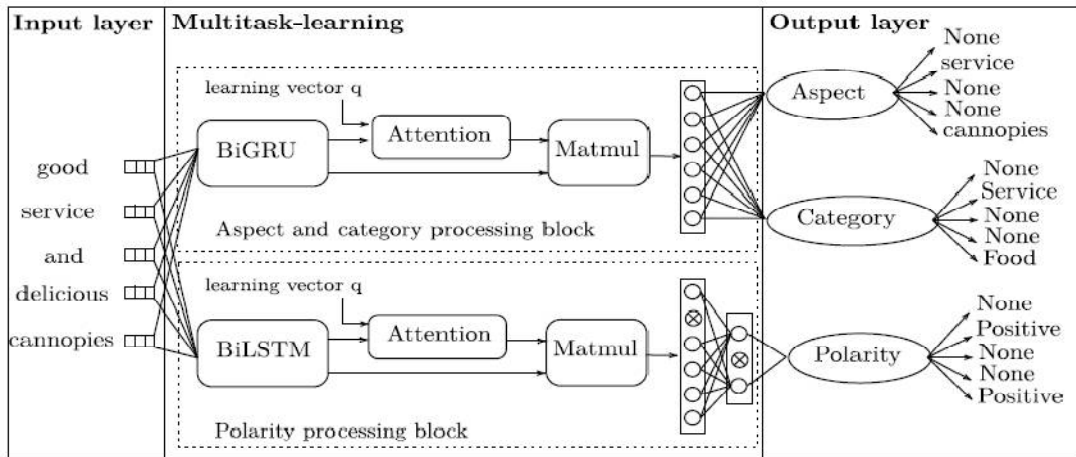
Computational Nets (CNs)

xAI.

2023-144

Architecture of ASAM

Attention based Sentiment Analysis Method (ASAM)



Models

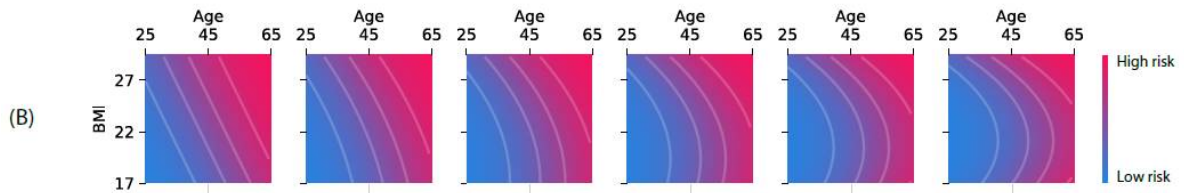
xAI.

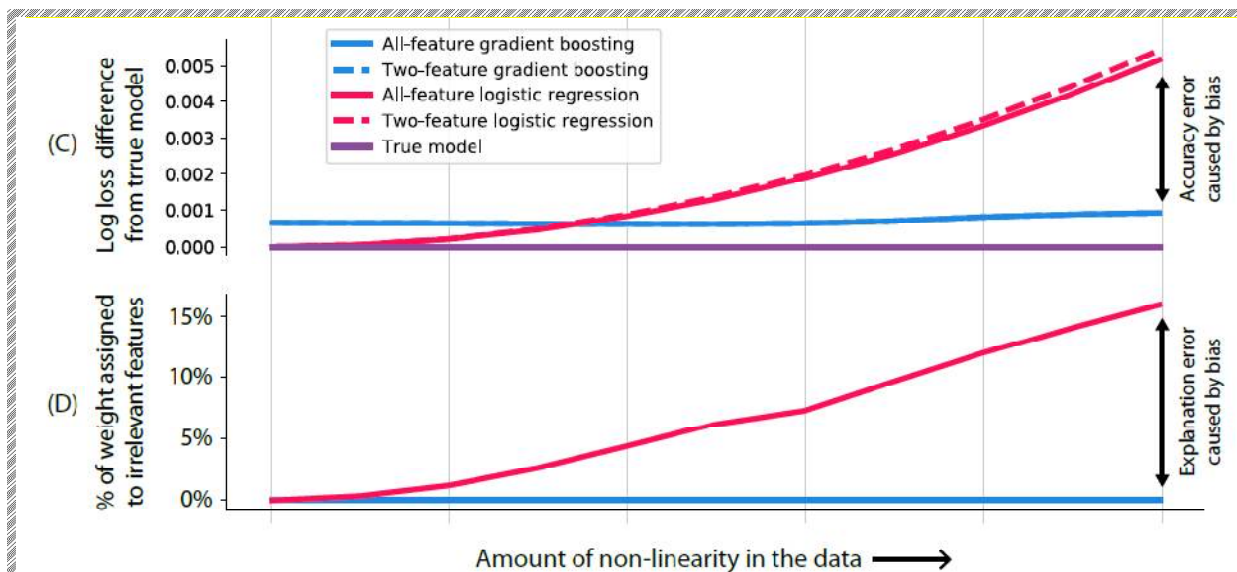
2023-158

Gradient boosted tree models

- + More accurate than neural networks
- + More interpretable than linear models

	Gradient Boosted Trees	Linear Model	Neural Network
(A) NHANES I Mortality (C-statistic)	0.821	0.813	0.816
CRIC Kidney Disease (area under the PR curve)	0.890	0.871	0.872
Hospital Procedure Duration (R ² value)	0.674	0.595	0.629





- ! (A) Gradient boosted tree models outperform both linear models and neural networks on medical datasets
- ! (B-D) Linear models exhibit explanation error as well as accuracy error in the presence of non-linearity
- ! (B) Data generating models used for the simulation, ranging from linear to quadratic along the body mass index (BMI) dimension
- ! (C) The test performance of linear logistic regression (red) is better than gradient boosting (blue) up until a specific amount of non-linearity. Not surprisingly, the bias of the linear model is higher than the gradient boosting model as shown by the steeper slope with increase the non-linearity.
- ! (D) As the true function becomes more non-linear the linear model assigns more credit (coefficient weight) to features that were not used by the data generating model.

Explanatory texts by ECDM-SDAM methodology to justify its final ranking

Positive explanatory text

- The user should book at the restaurant *The Ivy* since it obtains the highest overall rating.
- Its criterion of greatest interest, *food*, reaches a rating of 8.9 out of 10.
- The pie, crab, steak tartare, and liver stand out positively.
- Two of the expert sentences that most benefit this restaurant being selected as the best are: “good service and delicious food” and “constantly the best.”

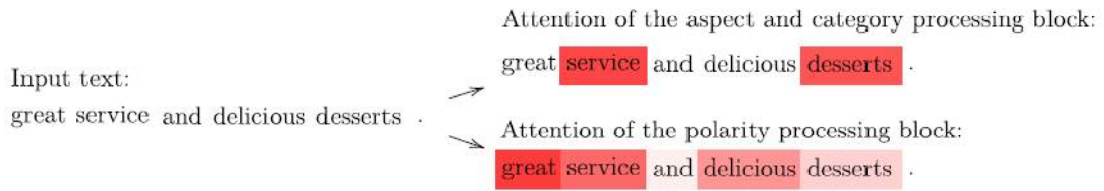
Negative explanatory text

- The restaurant *The Oxo Tower* is high quality although it is the last one of the ranking, so we identify its weakest points.
- Its most detrimental criterion is *drinks*.
- The acoustics, manager, and waiter stand out negatively.
- Two of the expert sentences that most harm to this restaurant being selected as the last one are: “poor service, meagre portions” and “not worth the trip”.

✓ SDAM: Subgroup Discovery and Attention Mechanisms

✓ ECDM: Explainable Crowd Decision Making

Visualization of weights of two attention layers of ASAM model



✓ Color intensity indicates relevancy or attention value of the word

Weighted Collective Evaluation matrix with identified criteria from the explanations

	restaurant	food	service	drinks	ambience	location	
Oxo Tower	0.27	0.28	0.181	0.0283	0.075	0.02837	<div style="display: flex; flex-direction: column; align-items: center;"> <div style="border: 1px solid green; border-radius: 50%; width: 15px; height: 15px; margin-bottom: 5px;"></div> BestCriterion <div style="border: 1px solid red; border-radius: 50%; width: 15px; height: 15px; margin-top: 5px;"></div> WorstCriterion </div>
J. Sheekey	0.28	0.31	0.187	0.028	0.0776	0.028	
The Wolseley	0.274	0.29	0.186	0.028	0.071	0.028	
The Ivy	0.283	0.3	0.19	0.029	0.077	0.028	

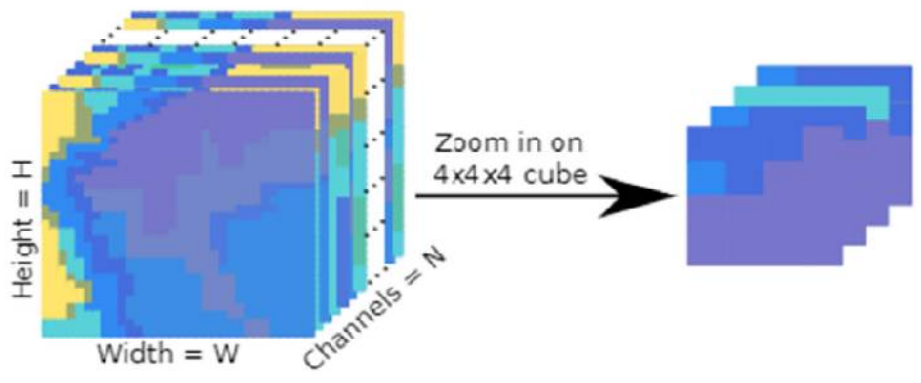
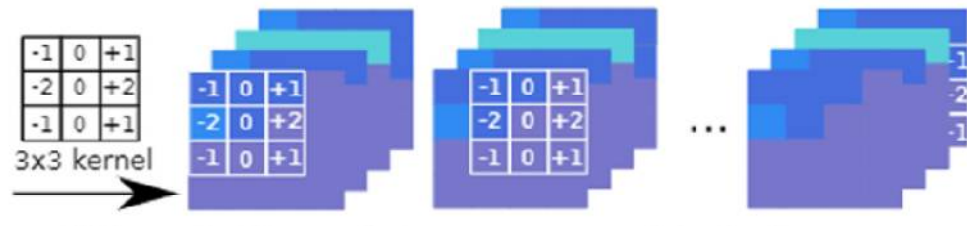
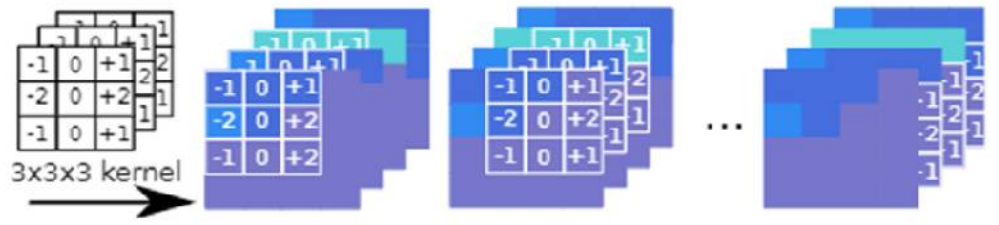
Significant rules extracted from the BOC tables using the Apriori-SD algorithm.

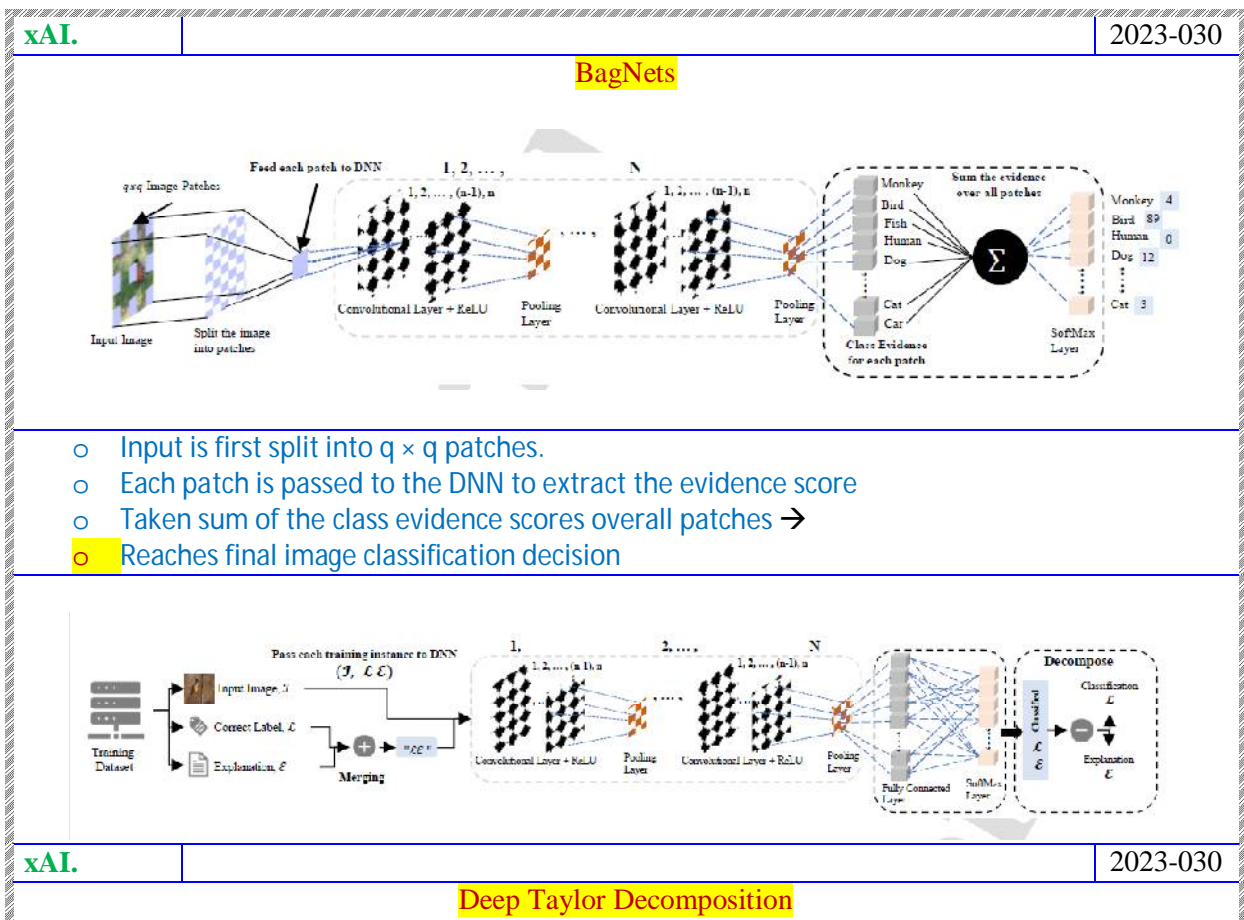
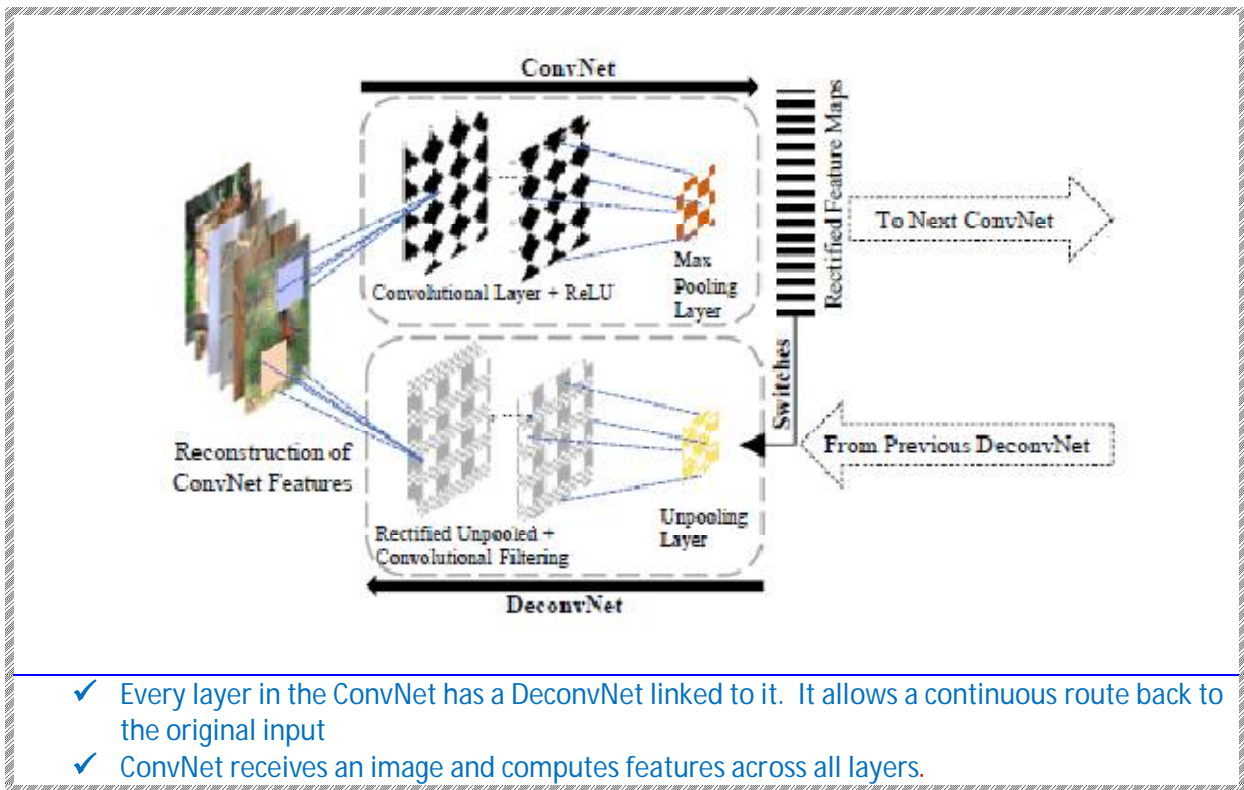
Restaurant	Rule	NWRAcc	Support	Confidence
The Ivy	{meal} → positive	0.52	0.027	1
The Ivy	{pie} → positive	0.51	0.022	1
The Ivy	{crab} → positive	0.51	0.022	1
The Ivy	{menu} → positive	0.51	0.022	1
The Ivy	{steak tartare} → positive	0.508	0.16	1
The Ivy	{liver} → positive	0.508	0.16	1
The Oxo Tower	{drinks} → negative	1	0.036	1
The Oxo Tower	{acoustics} → negative	0.75	0.023	1
The Oxo Tower	{manager} → negative	0.58	0.018	1
The Oxo Tower	{waiter} → negative	0.58	0.018	1

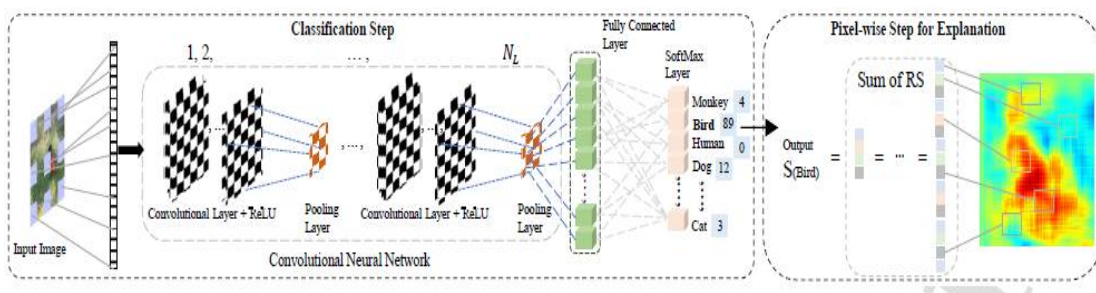
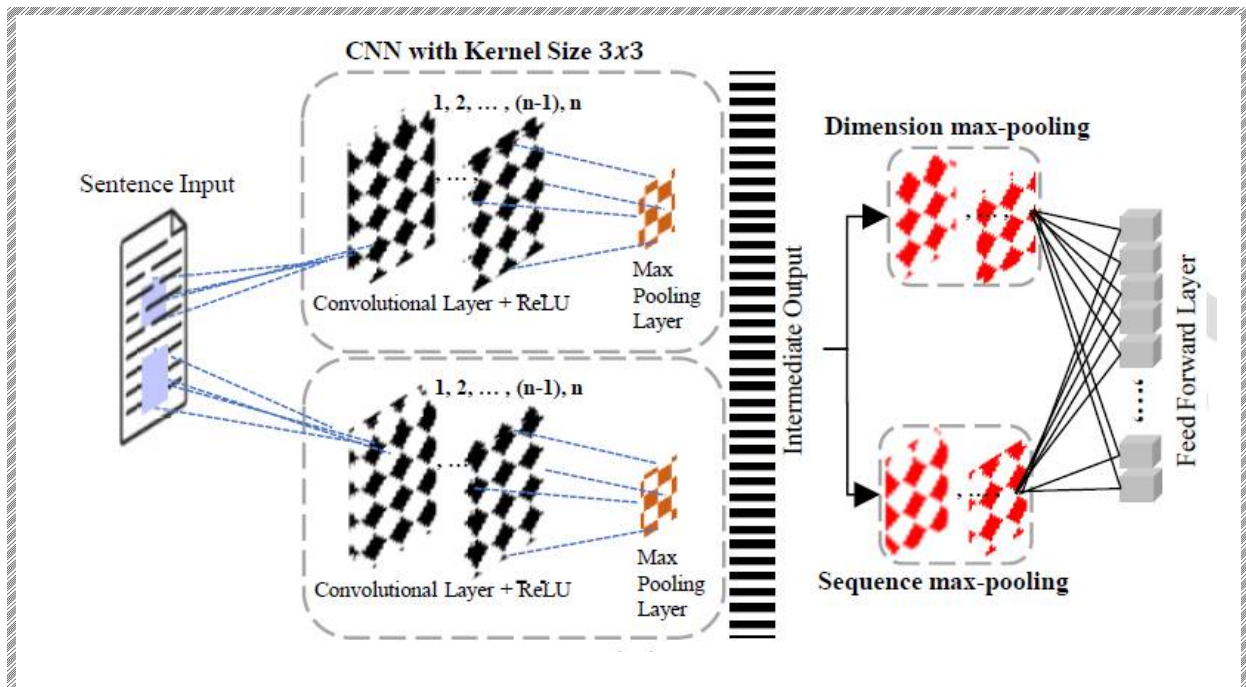
Significant top 3 positive and negative sentences.

<p>Ranking 1 Positive sentences for The Ivy (food criterion)</p>	<ol style="list-style-type: none"> 1. Good service and delicious food. 2. Food and service as always terrific. 3. The food, service and sense of occasion was truly perfect.
<p>Ranking 2 Positive sentences for The Ivy (restaurant criterion)</p>	<ol style="list-style-type: none"> 1. Constantly the best 2. Always very good 3. Will definitely revisit for a special occasion.
<p>Ranking 3 Negative sentences for The Oxo Tower (any criterion)</p>	<ol style="list-style-type: none"> 1. Not worth the trip 2. Poor service, meagre portions 3. Overall, therefore, it is poor value and plays to the tourist market.

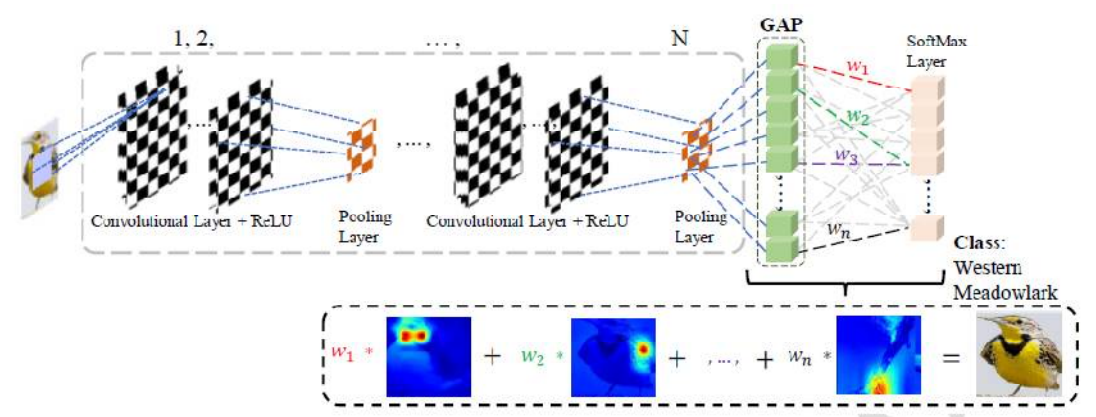
Convolution- 2D-3D

xAI.		2022-
2D and 3D convolution on 3D data cube to generate feature maps		
		
(a) Example of a 3d raster, or data cube.		
		
(b) Example 2D convolution to generate spatial-wise features.		
		
(c) Example 3D convolution to generate spatial-channel-wise features.		

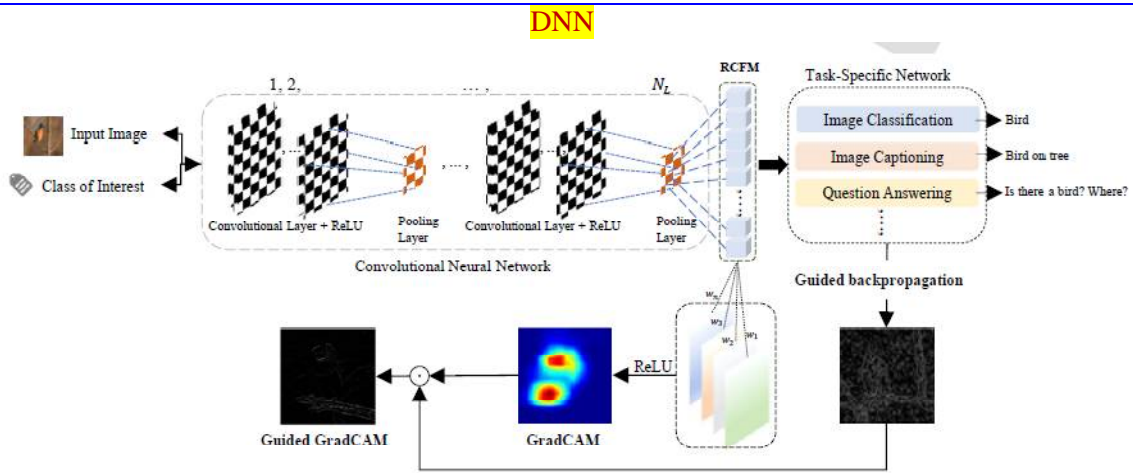




- 🔔 I step: input image has been identified as a bird
- 🔔 II step: model's reception of the features is shown as a heat map based on the relevance scores estimated from each hidden layer
- 🔔 Pixels surrounding the bird's location had a substantial impact on the outcome
- 🔔 Red regions: proved useful in the decision
- 🔔 Blue regions: not helpful in decision



✓ Discriminative areas, distinct to each class, are highlighted in the CAM

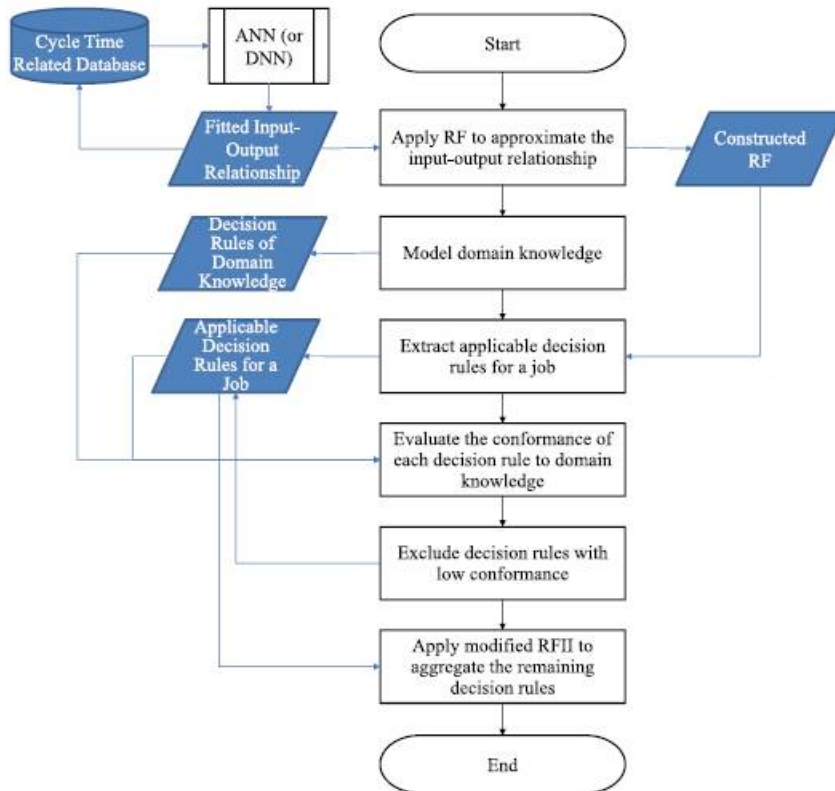


xAI.

Cycle time prediction

2023-005

Process of methodology



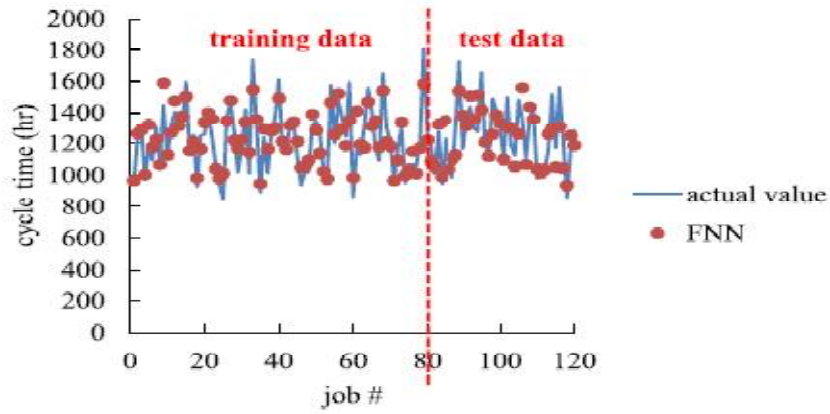


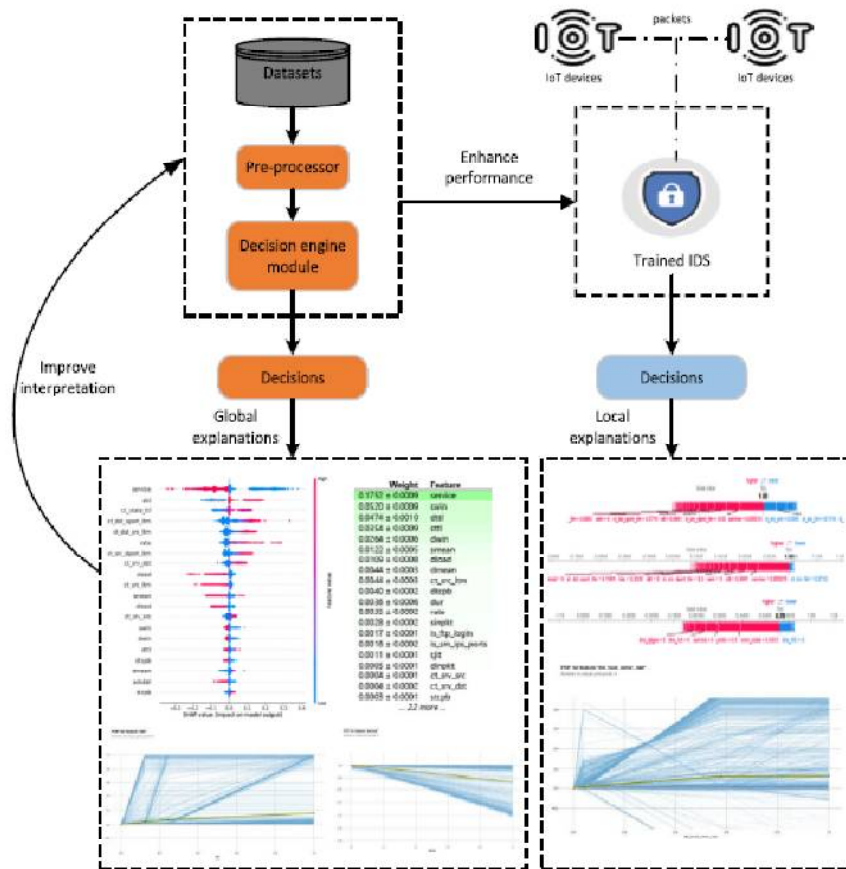
Fig. 2. Prediction results.

xAI.

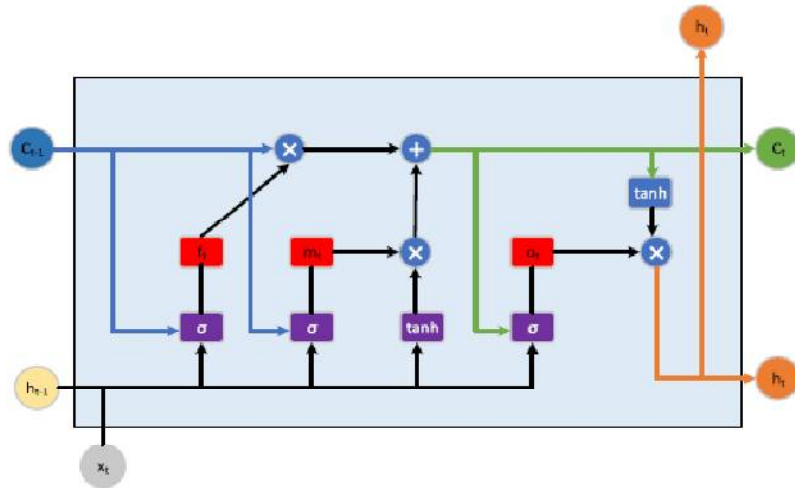
S: Shapley Additive exPlanations
 P: Permutation Feature Importance
 I: Individual Conditional Expectation
 P: Partial Dependence Plot

2023-057

Overview of the structure of SPIP framework



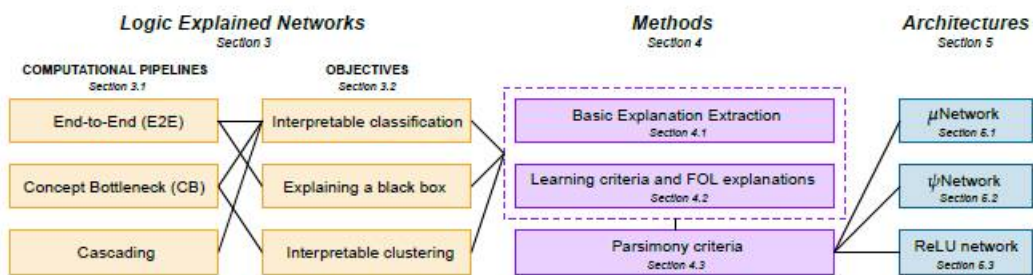
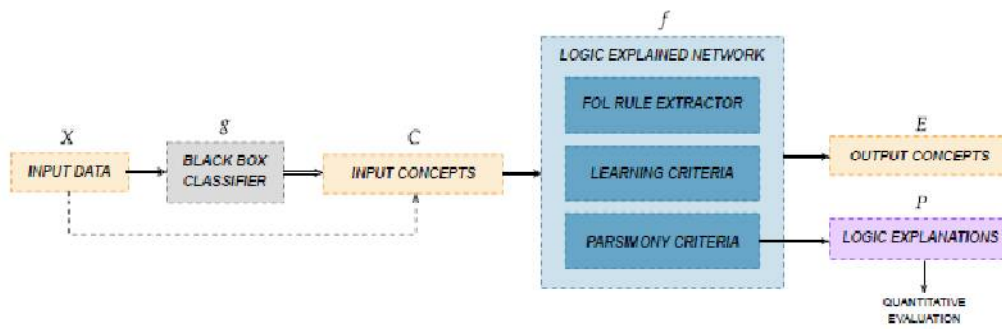
Long Short-Term Memory (RecNN)



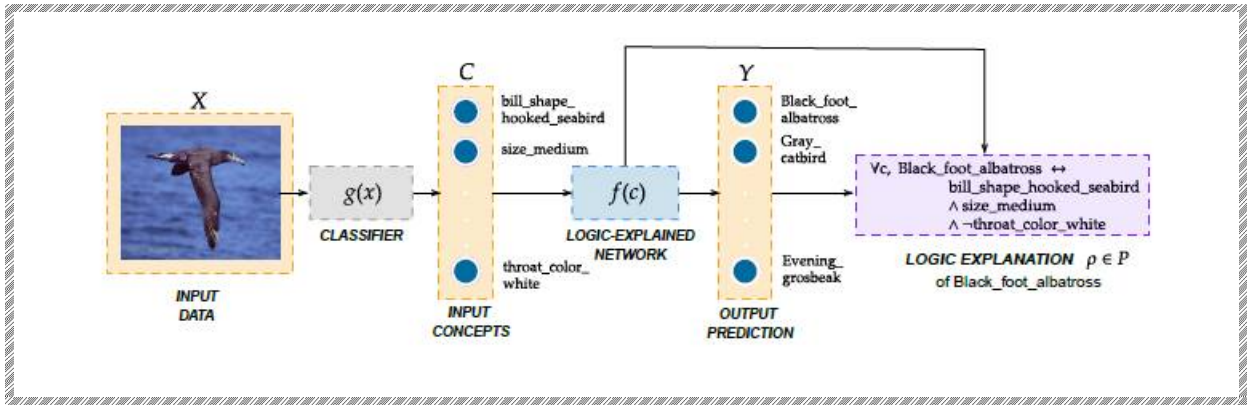
xAI.

2023-081

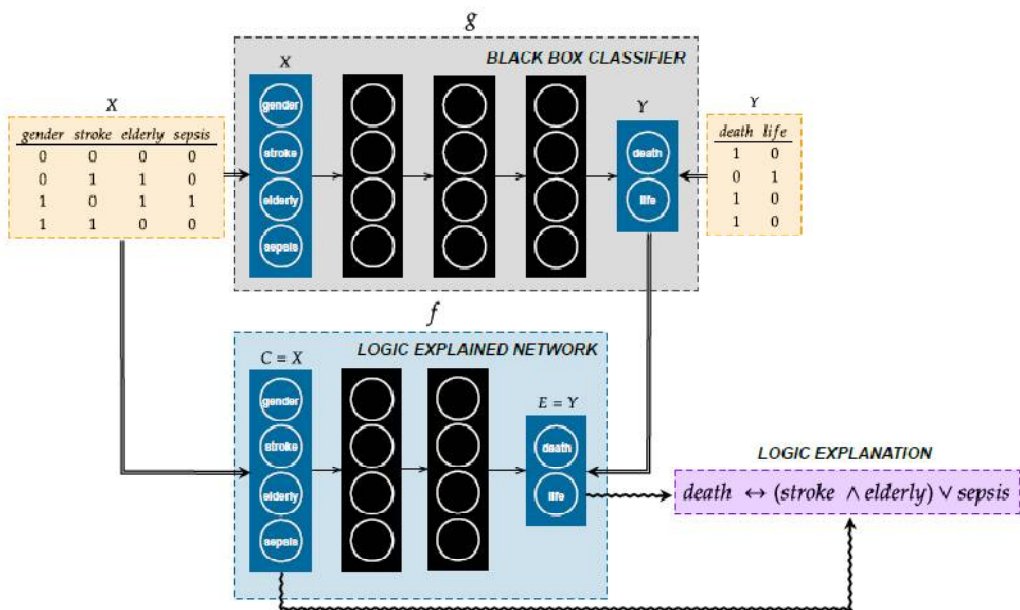
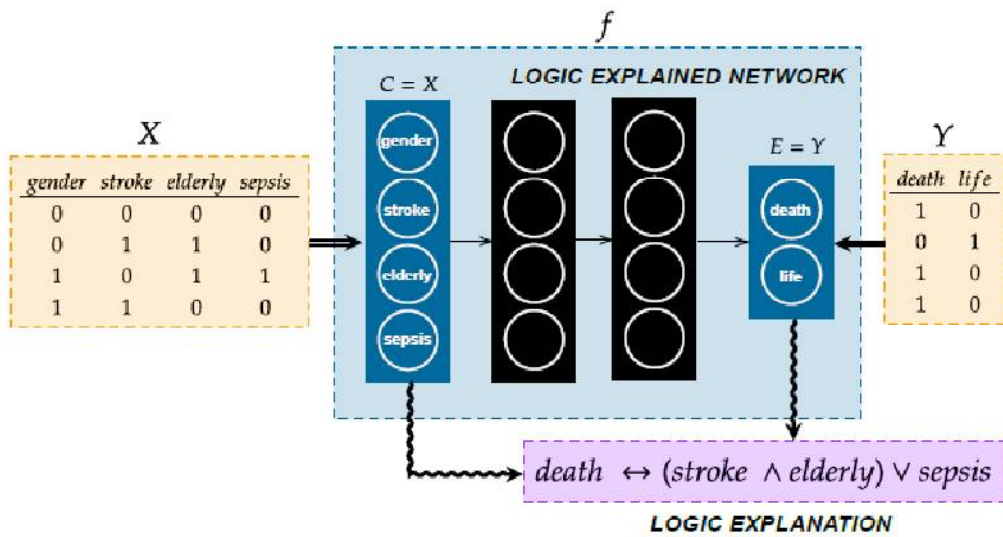
Logic Explained Networks (LENs, f)



- ✓ Starting from the nodes on the left, each path that ends to one of the nodes on the right → creating a specific instance of the LEN framework

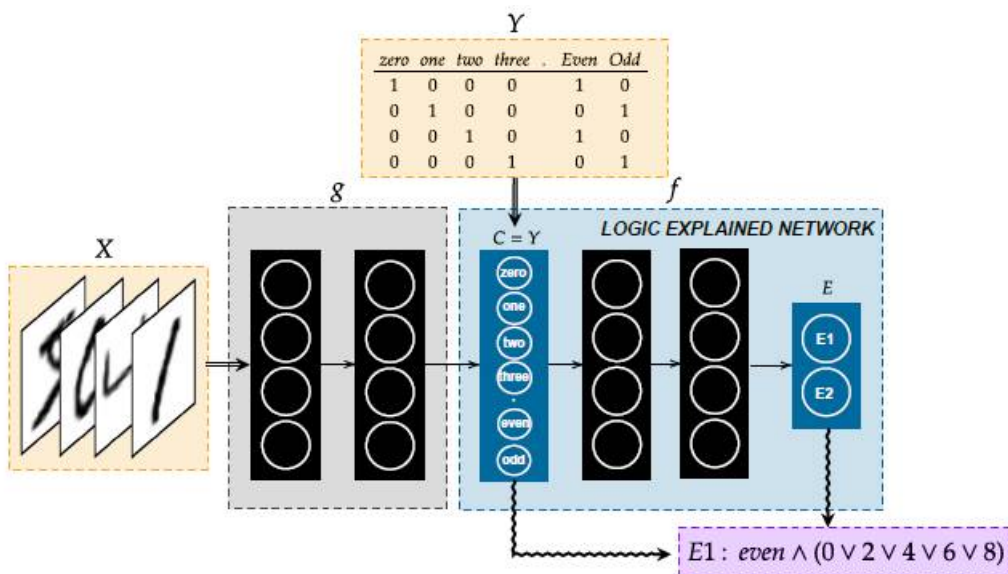
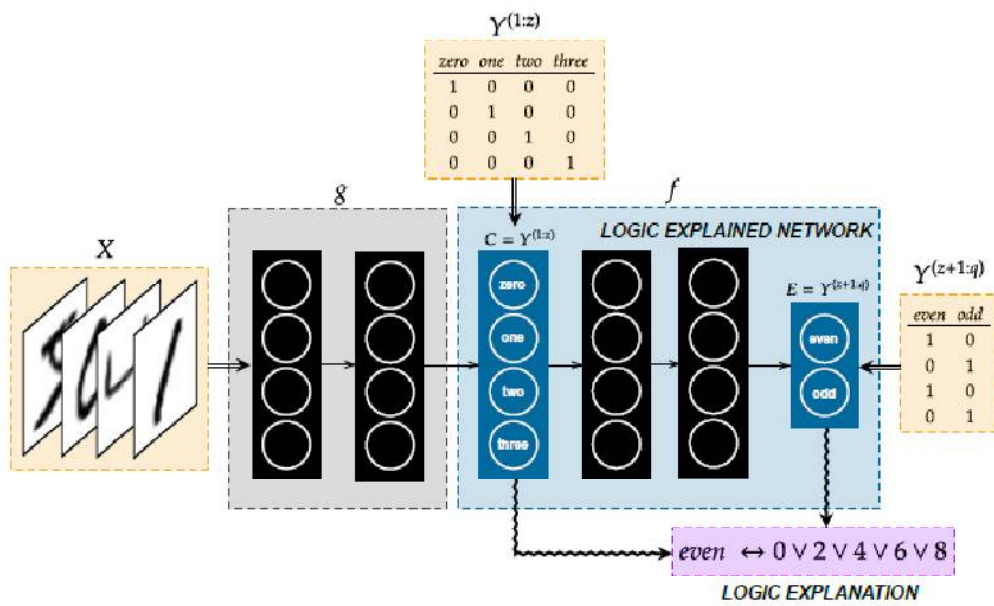


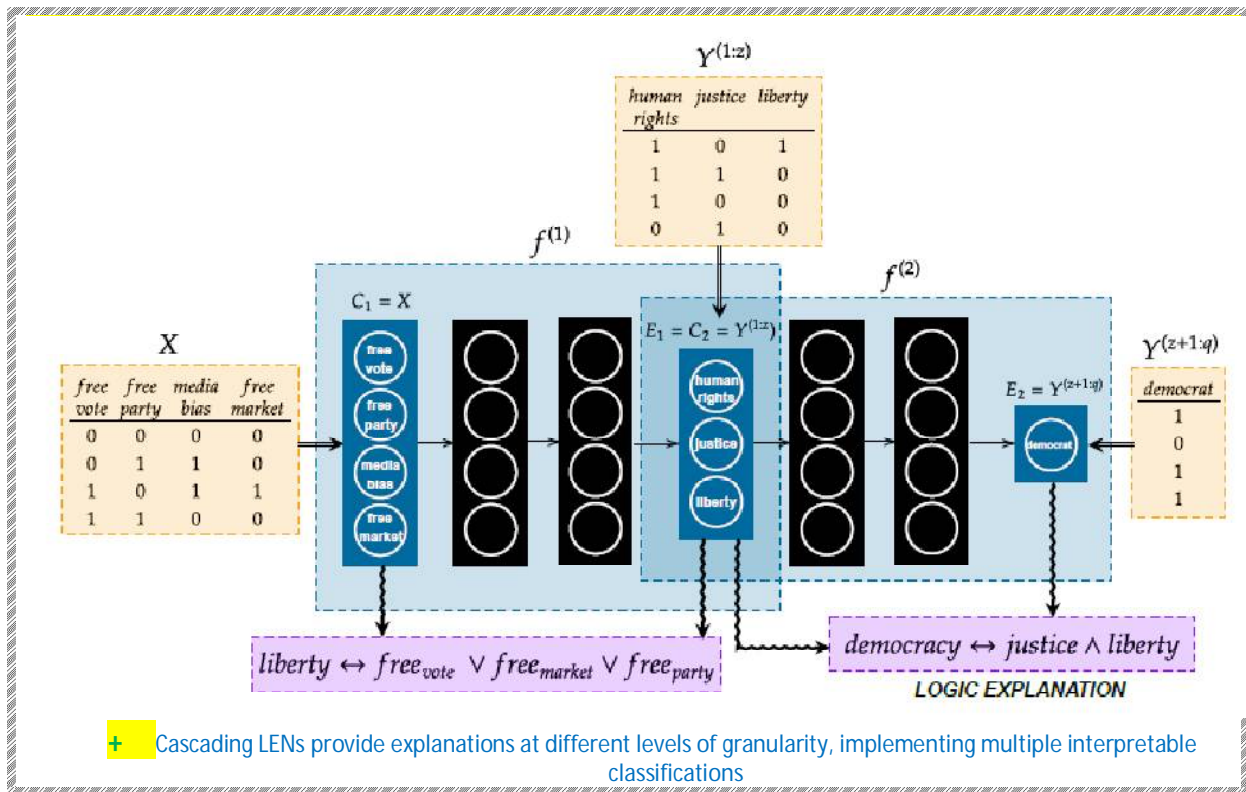
LEN (interpretable classification) solves classification problem



+ LEN provides explanations of a black-box classifier

Cascading LENS





xAI. 2023-081

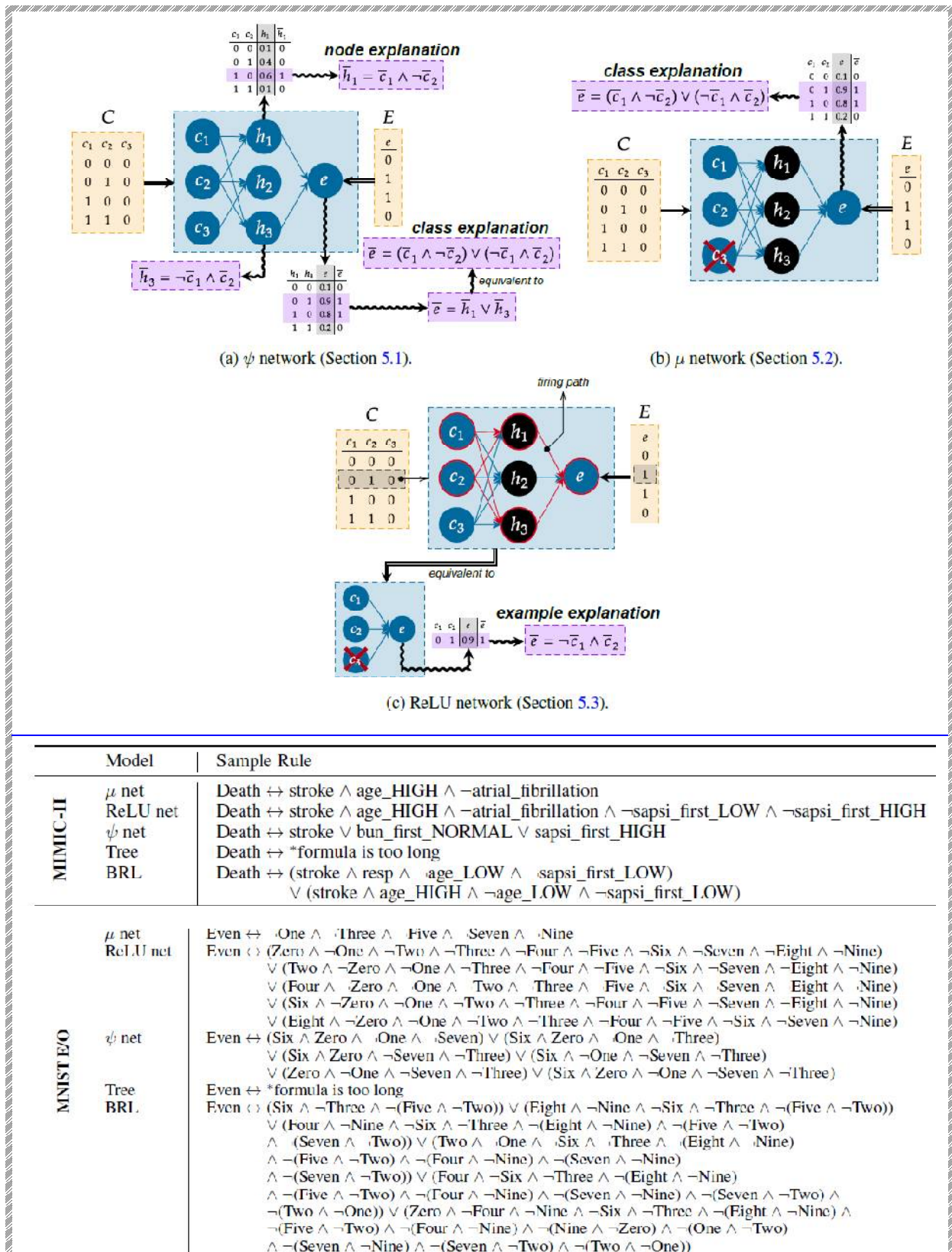
Empirical truth table T_i
of i-th LEN output f_i, with k = 4 and m_i = 2

Concept Data				Empirical Truth Table		
c ₁	c ₂	c ₃	c ₄	\bar{c}_1	\bar{c}_2	\bar{f}_i
0.2	0.1	0.3	0.7	0	0	0
0.2	0.9	0.1	0.9	0	1	1
0.3	0.8	0.2	0.9	0	1	1
0.7	0.3	0.1	0.7	1	0	1
0.9	0.1	0.1	0.8	1	0	1
0.9	0.9	0.3	0.9	1	1	1

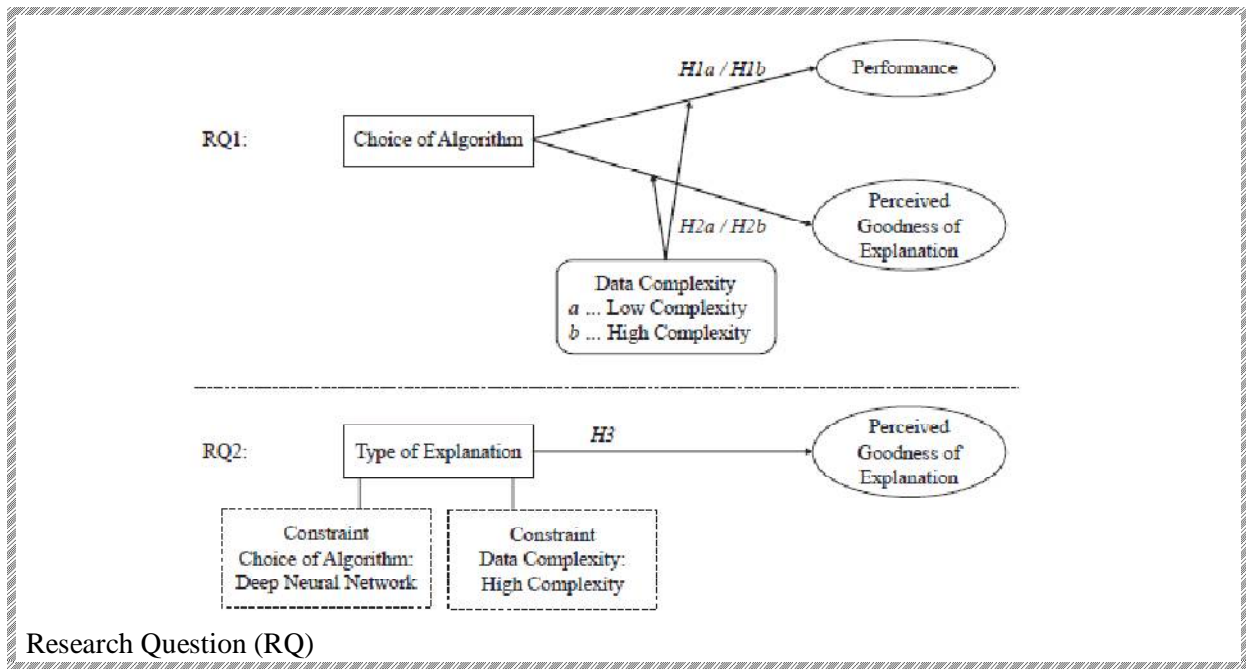
Booleanization + pruning

- class-level explanation → $(\neg \bar{c}_1 \wedge \bar{c}_2) \vee (\bar{c}_1 \wedge \neg \bar{c}_2) \vee (\bar{c}_1 \wedge \bar{c}_2)$
- set-level explanation → $(\bar{c}_1 \wedge \neg \bar{c}_2) \vee (\bar{c}_1 \wedge \bar{c}_2)$
- example-level explanation → $\bar{c}_1 \wedge \bar{c}_2$

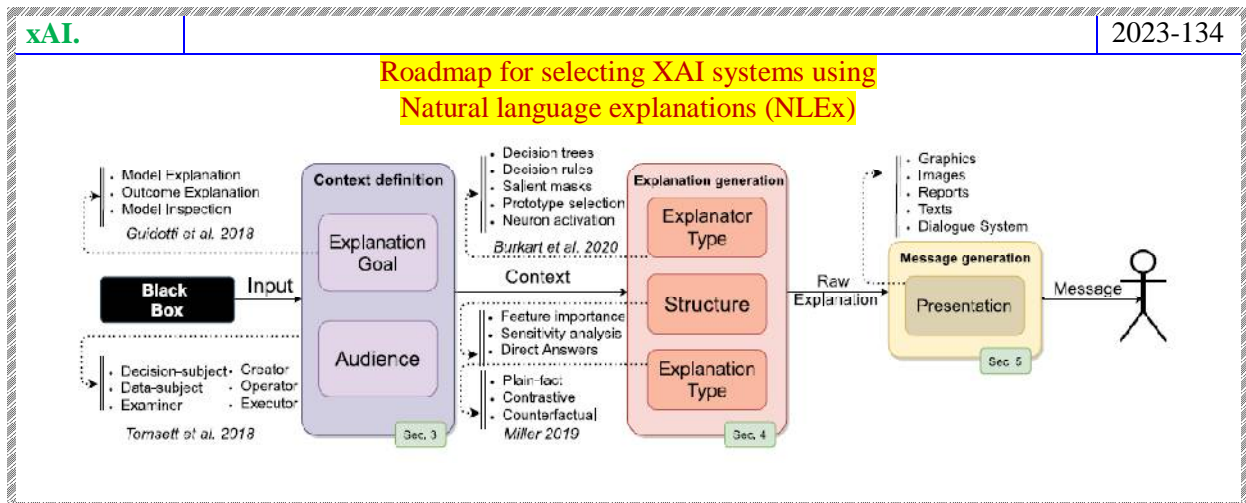
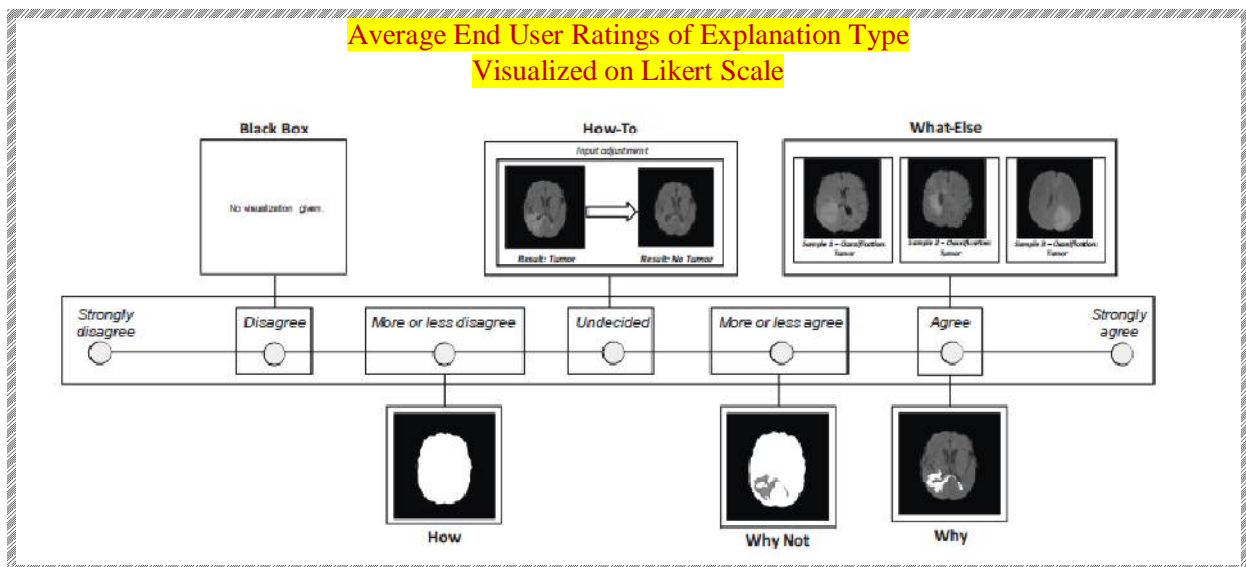
Out-of-the-box LENs



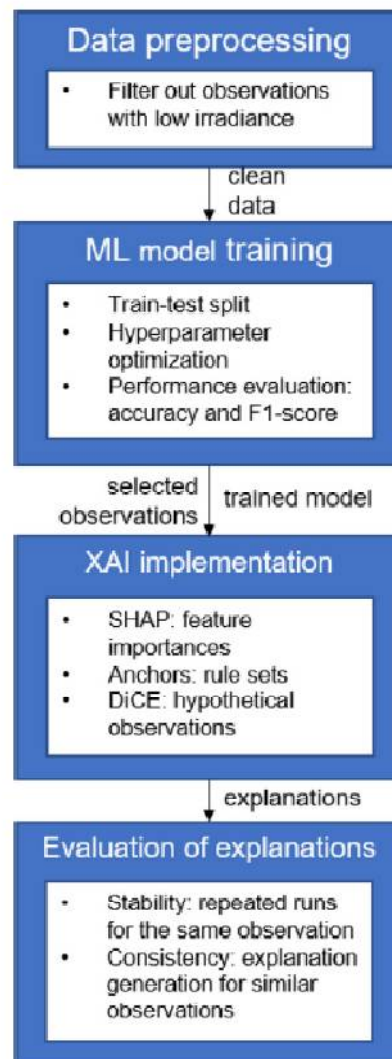
Model	Sample Rule	
MIMIC-II	μ net	Death \leftrightarrow stroke \wedge age_HIGH \wedge \neg atrial_fibrillation
	ReLU net	Death \leftrightarrow stroke \wedge age_HIGH \wedge \neg atrial_fibrillation \wedge \neg sapsi_first_LOW \wedge \neg sapsi_first_HIGH
	ψ net	Death \leftrightarrow stroke \vee bun_first_NORMAL \vee sapsi_first_HIGH
	Tree	Death \leftrightarrow *formula is too long
	BRI	Death \leftrightarrow (stroke \wedge resp \wedge age_LOW \wedge sapsi_first_LOW) \vee (stroke \wedge age_HIGH \wedge \neg age_LOW \wedge \neg sapsi_first_LOW)
MINIST/O	μ net	Even \leftrightarrow \neg One \wedge \neg Three \wedge \neg Five \wedge \neg Seven \wedge \neg Nine
	ReLU net	Even \leftrightarrow (Zero \wedge \neg One \wedge \neg Two \wedge \neg Three \wedge \neg Four \wedge \neg Five \wedge \neg Six \wedge \neg Seven \wedge \neg Eight \wedge \neg Nine) \vee (Two \wedge \neg Zero \wedge \neg One \wedge \neg Three \wedge \neg Four \wedge \neg Five \wedge \neg Six \wedge \neg Seven \wedge \neg Eight \wedge \neg Nine) \vee (Four \wedge Zero \wedge One \wedge Two \wedge Three \wedge Five \wedge Six \wedge Seven \wedge Eight \wedge \neg Nine) \vee (Six \wedge \neg Zero \wedge \neg One \wedge \neg Two \wedge \neg Three \wedge \neg Four \wedge \neg Five \wedge \neg Seven \wedge \neg Eight \wedge \neg Nine) \vee (Eight \wedge \neg Zero \wedge \neg One \wedge \neg Two \wedge \neg Three \wedge \neg Four \wedge \neg Five \wedge \neg Six \wedge \neg Seven \wedge \neg Nine)
	ψ net	Even \leftrightarrow (Six \wedge Zero \wedge \neg One \wedge \neg Seven) \vee (Six \wedge Zero \wedge One \wedge \neg Three) \vee (Six \wedge Zero \wedge \neg Seven \wedge \neg Three) \vee (Six \wedge \neg One \wedge \neg Seven \wedge \neg Three) \vee (Zero \wedge \neg One \wedge \neg Seven \wedge \neg Three) \vee (Six \wedge Zero \wedge \neg One \wedge \neg Seven \wedge \neg Three)
	Tree	Even \leftrightarrow *formula is too long
	BRI	Even \leftrightarrow (Six \wedge \neg Three \wedge \neg (Five \wedge \neg Two)) \vee (Eight \wedge \neg Nine \wedge \neg Six \wedge \neg Three \wedge \neg (Five \wedge \neg Two)) \vee (Four \wedge \neg Nine \wedge \neg Six \wedge \neg Three \wedge \neg (Eight \wedge \neg Nine) \wedge \neg (Five \wedge \neg Two) \wedge (Seven \wedge \neg Two)) \vee (Two \wedge \neg One \wedge \neg Six \wedge \neg Three \wedge (Eight \wedge \neg Nine) \wedge \neg (Five \wedge \neg Two) \wedge \neg (Four \wedge \neg Nine) \wedge \neg (Seven \wedge \neg Nine) \wedge \neg (Seven \wedge \neg Nine) \wedge \neg (Five \wedge \neg Two) \wedge \neg (Four \wedge \neg Nine) \wedge \neg (Seven \wedge \neg Nine) \wedge \neg (Seven \wedge \neg Two) \wedge \neg (Two \wedge \neg One)) \vee (Zero \wedge \neg Four \wedge \neg Nine \wedge \neg Six \wedge \neg Three \wedge \neg (Eight \wedge \neg Nine) \wedge \neg (Five \wedge \neg Two) \wedge \neg (Four \wedge \neg Nine) \wedge \neg (Nine \wedge \neg Zero) \wedge \neg (One \wedge \neg Two) \wedge \neg (Seven \wedge \neg Nine) \wedge \neg (Seven \wedge \neg Two) \wedge \neg (Two \wedge \neg One))



Research Question (RQ)

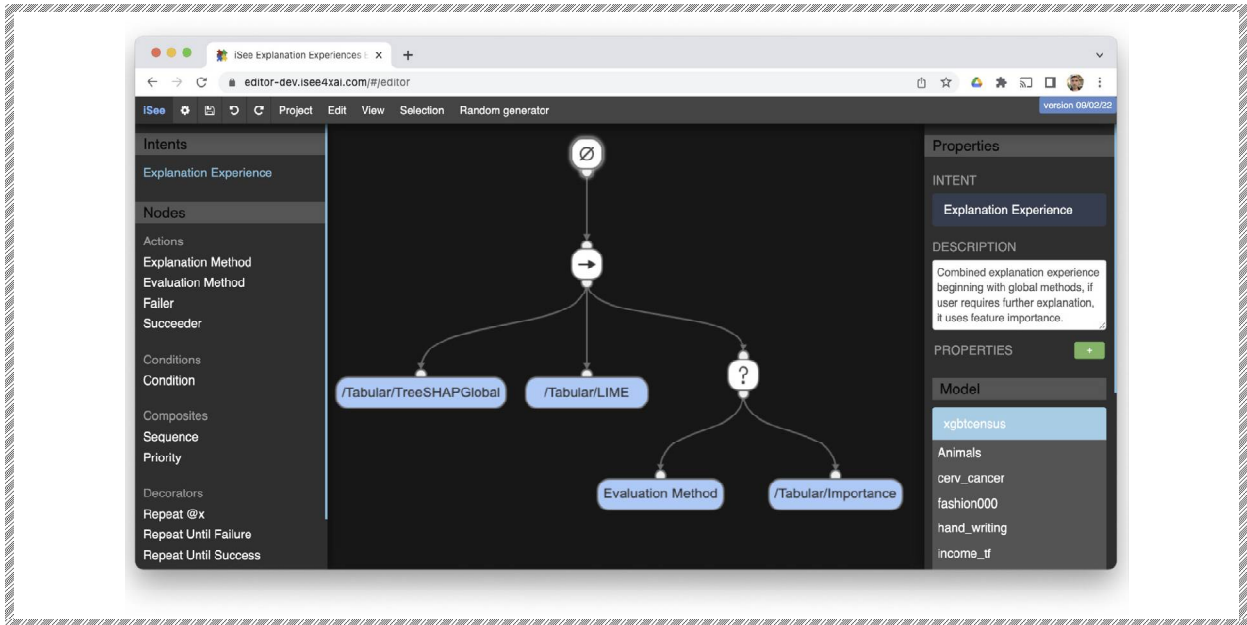


Method Flow xAI

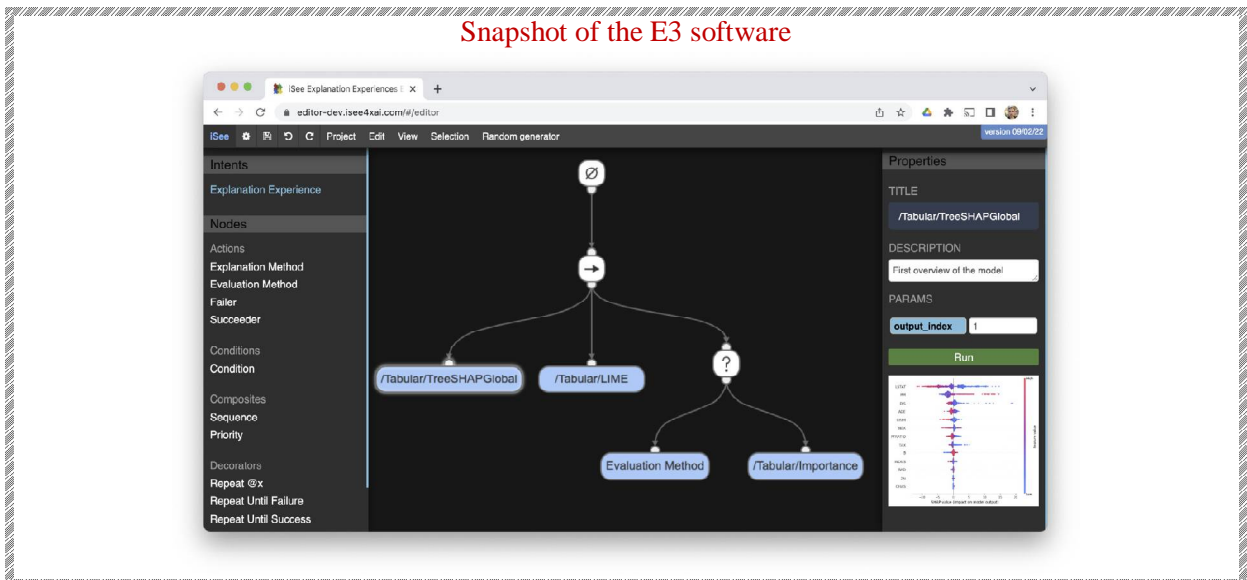


Software architecture Flow-diagram Evolving (Safe)





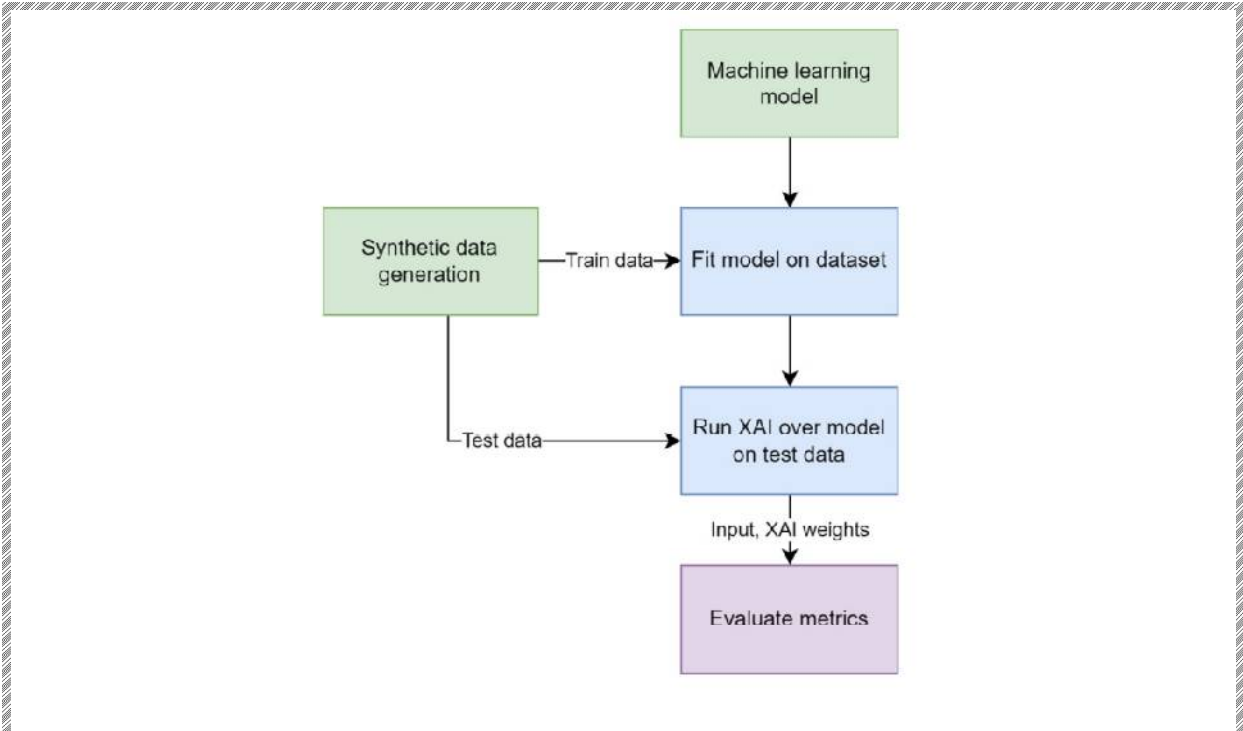
Snapshot of the E3 software



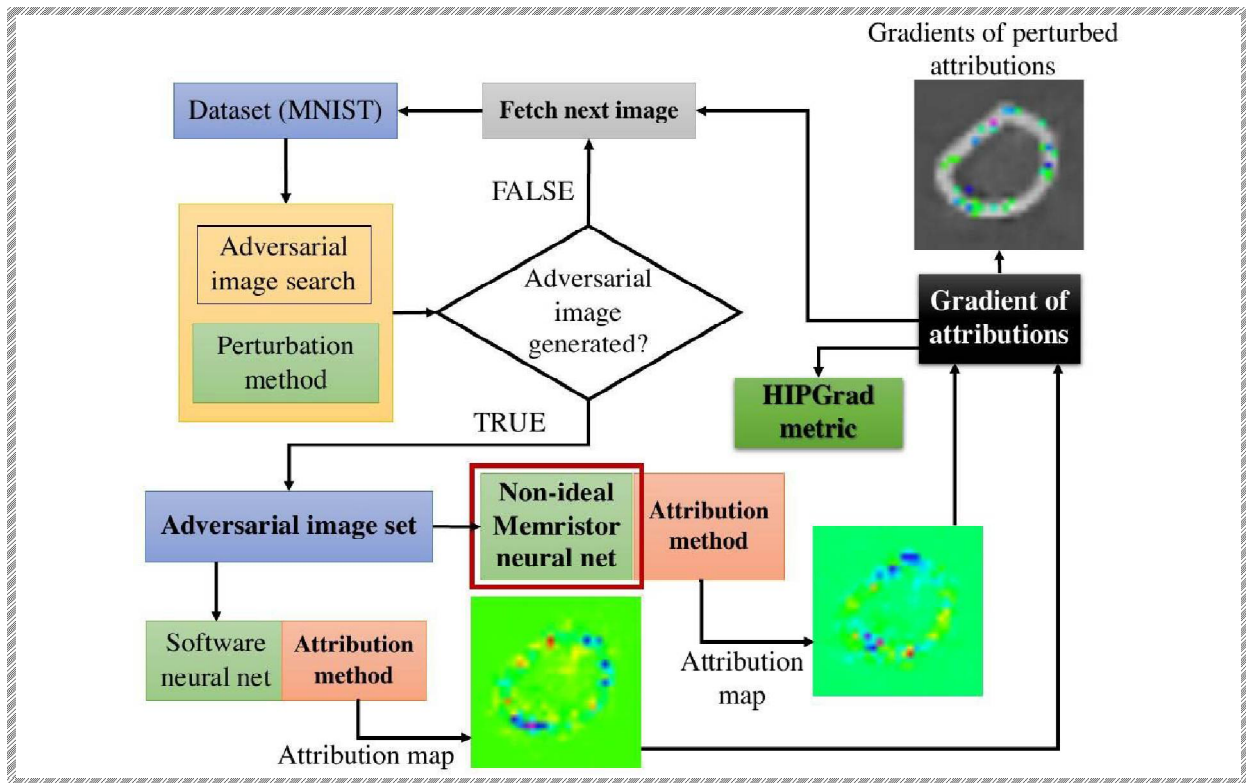
xAI.

2023-142

xAI-method



Library Name	Description	Library Version
NumPy	A library that implements linear algebra operations, mathematical functions, elements of statistical analysis	1.21.0
Matplotlib	Library for plotting various types of graphs	3.5.1
Scipy	Library designed to perform scientific and engineering calculations	1.8.0
Pandas	Library for working with tabular data structures	1.4.1
Shap	Library with implementation of the XAI SHAP method	0.40.0
Lime	Library with the implementation of the XAI LIME method	0.2.0.1
Scikit-learn	Library with tools for designing and training models	1.0.2



xAI. 2023-160

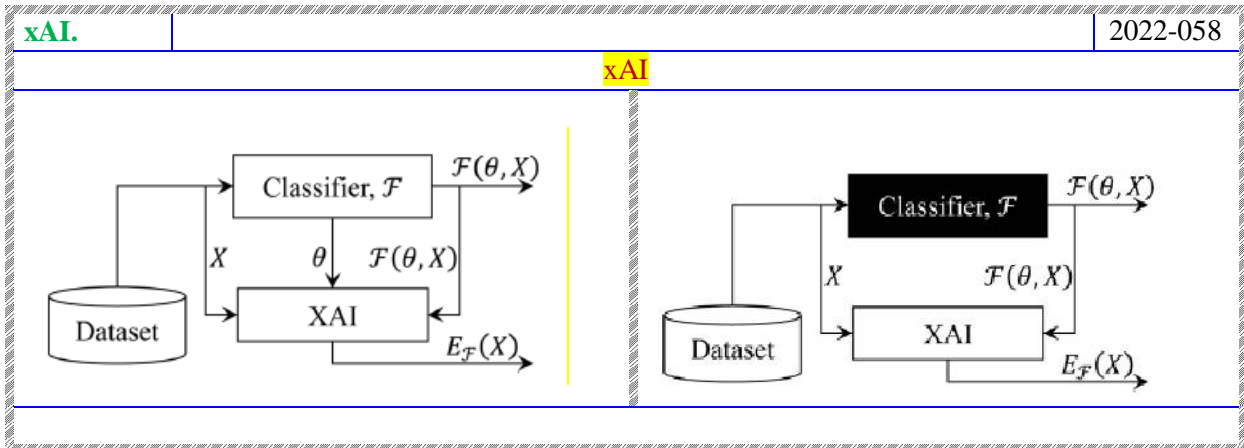
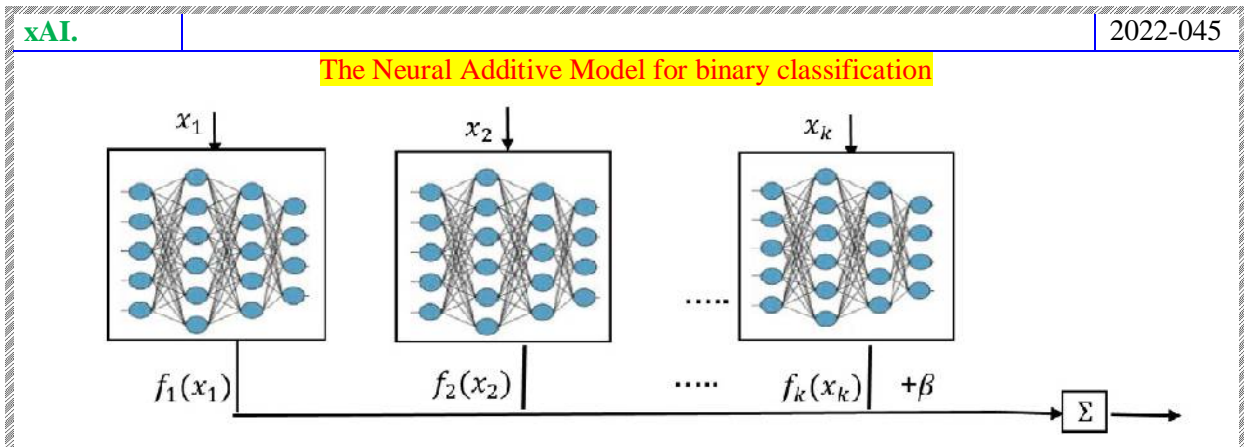
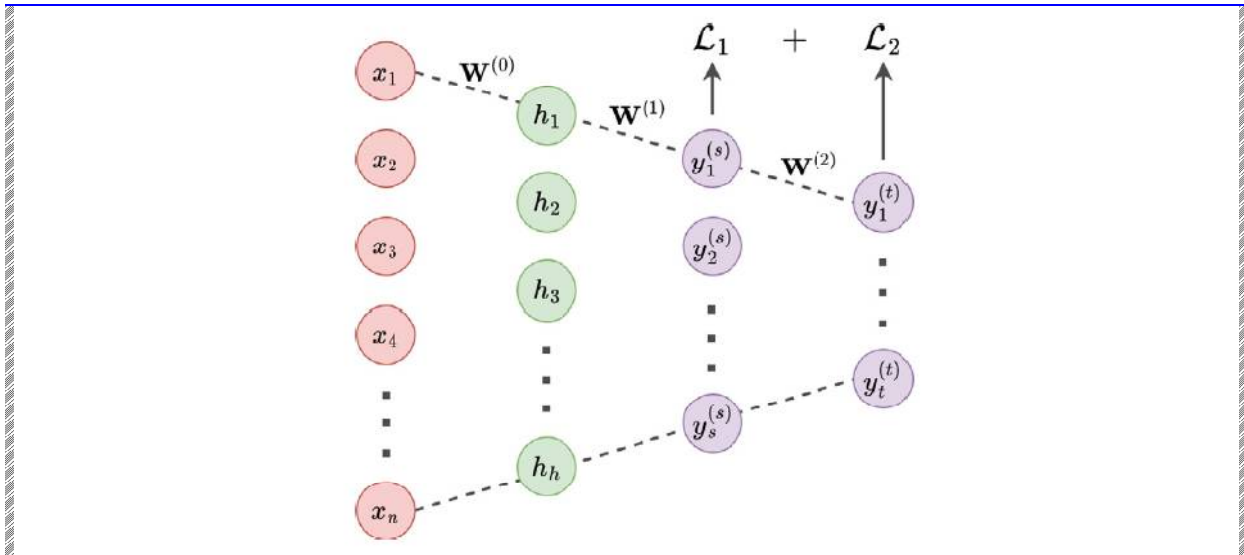
Libraries

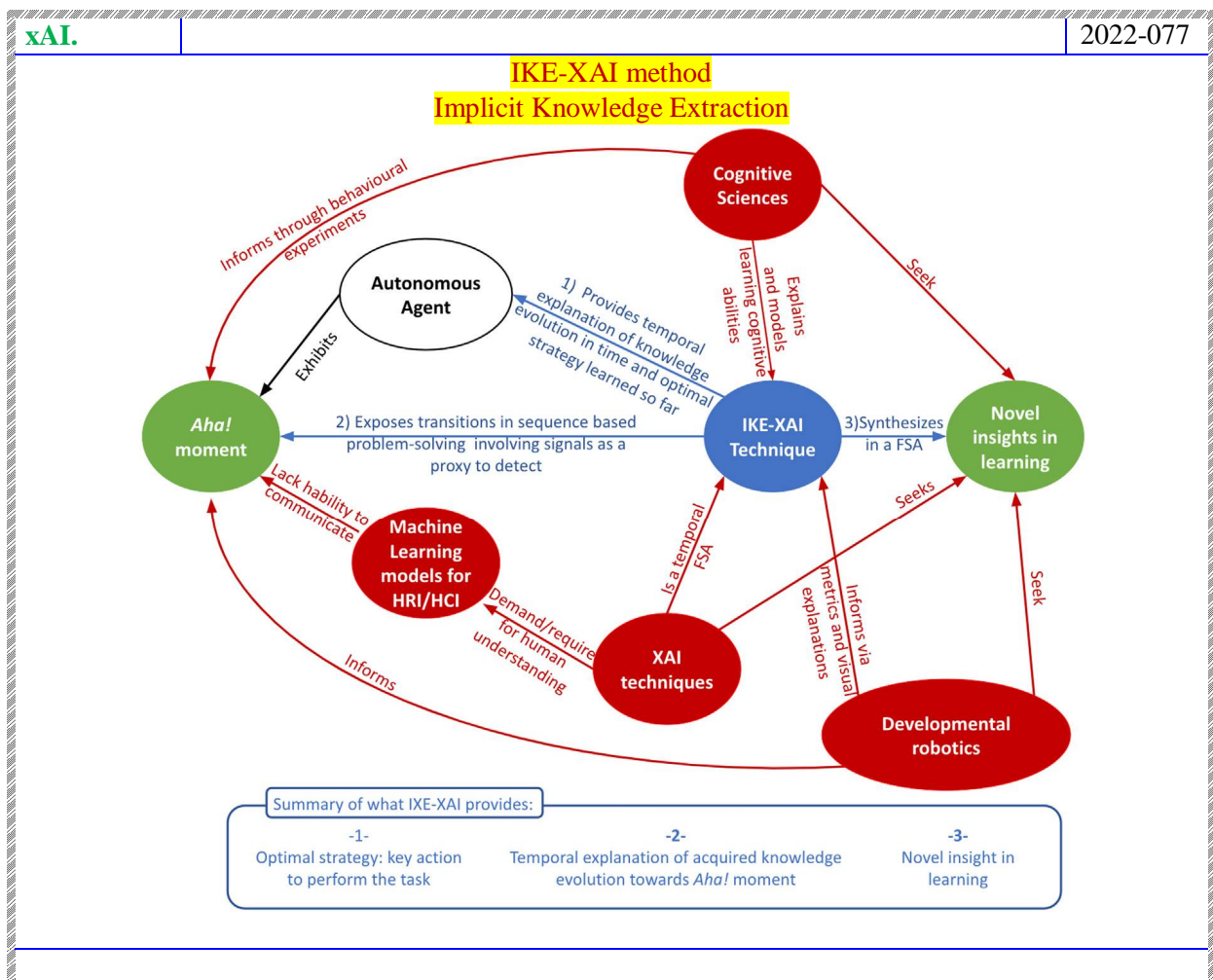
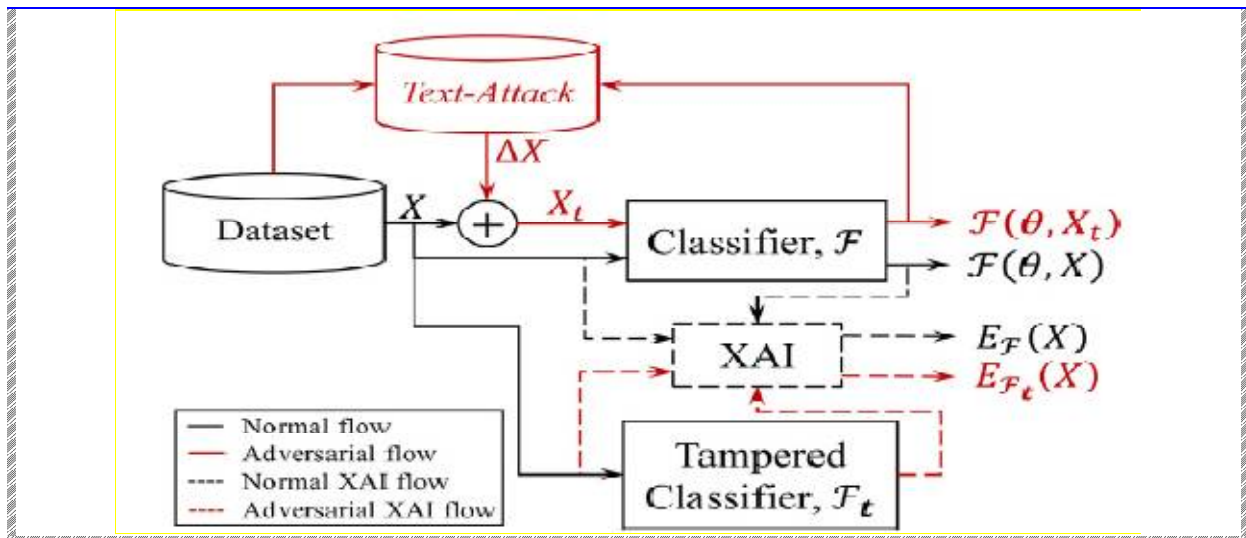
Library	Faithfulness	Robustness	Localisation	Complexity	Axiomatic	Randomisation
Captum (2)	1	1	0	0	0	0
AIX360 (2)	2	0	0	0	0	0
TorchRay (1)	0	0	1	0	0	0
Quantus (27)	9	4	6	3	3	2

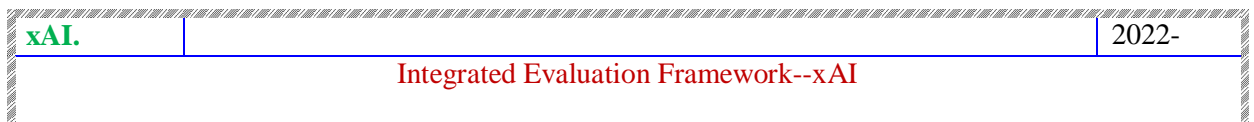
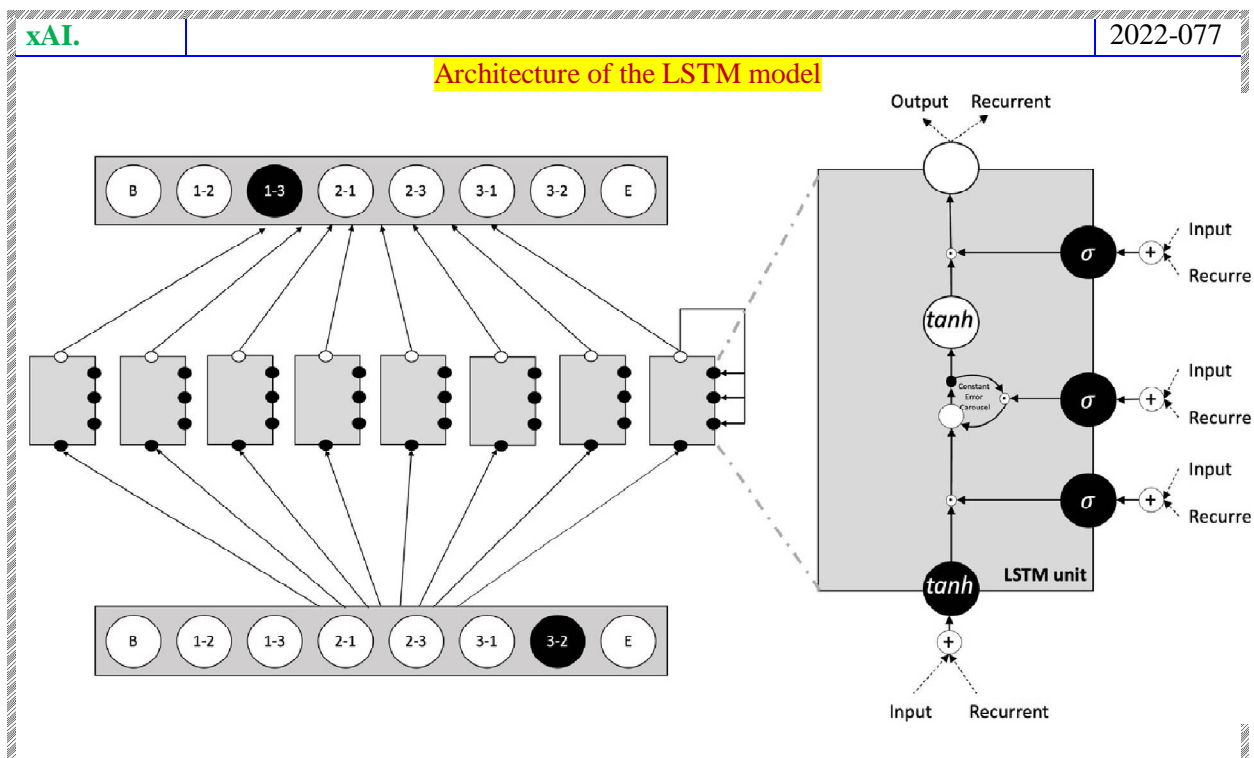
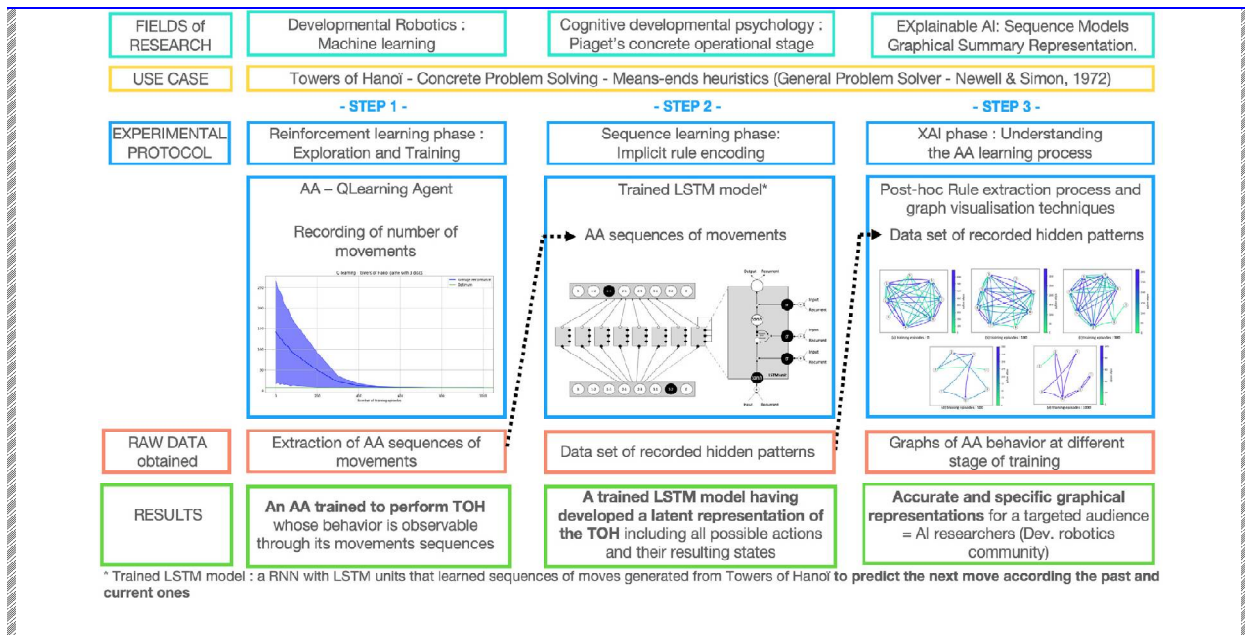
AIX360 (Arya et al., 2019), captum (Kokhlikyan et al., 2020), TorchRay (Fong et al., 2019) and Quantus)

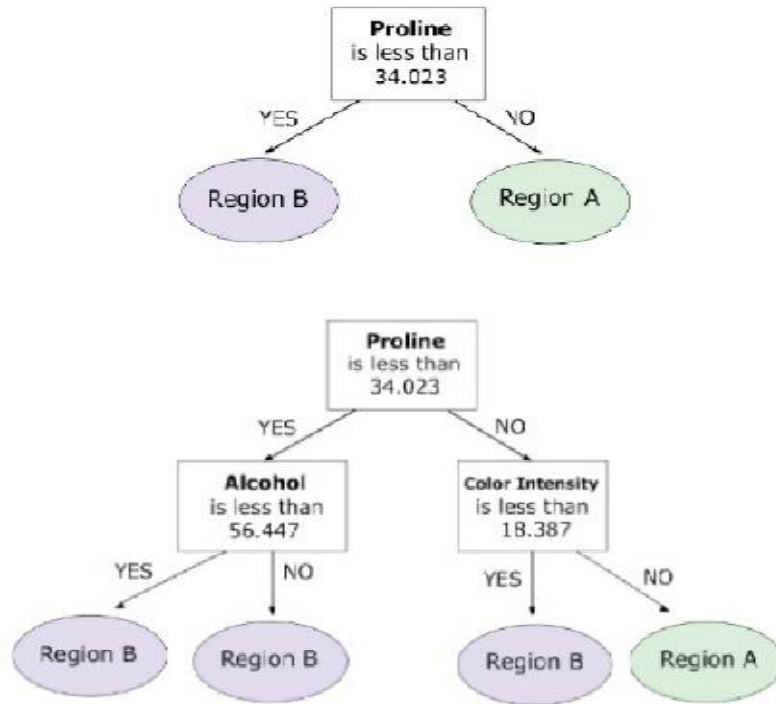
xAI. 2022-016

PLENARY
(exPLaining bLack-box modElS in Natural lAnguage thRough fuzzY linguistic summaries),
Architecture of multi-task neural network









Decision trees in categorization tasks

Please predict the region of origin of the following wine based on the values of its attributes and the decision tree.

Please indicate your answer by pressing A or B on your keyboard.

Total Phenols:	55.7
Alcohol:	57.1
Proline:	55.0
Flavonoids:	48.5
Hue:	35.2
Magnesium:	32.6
Color Intensity:	40.6

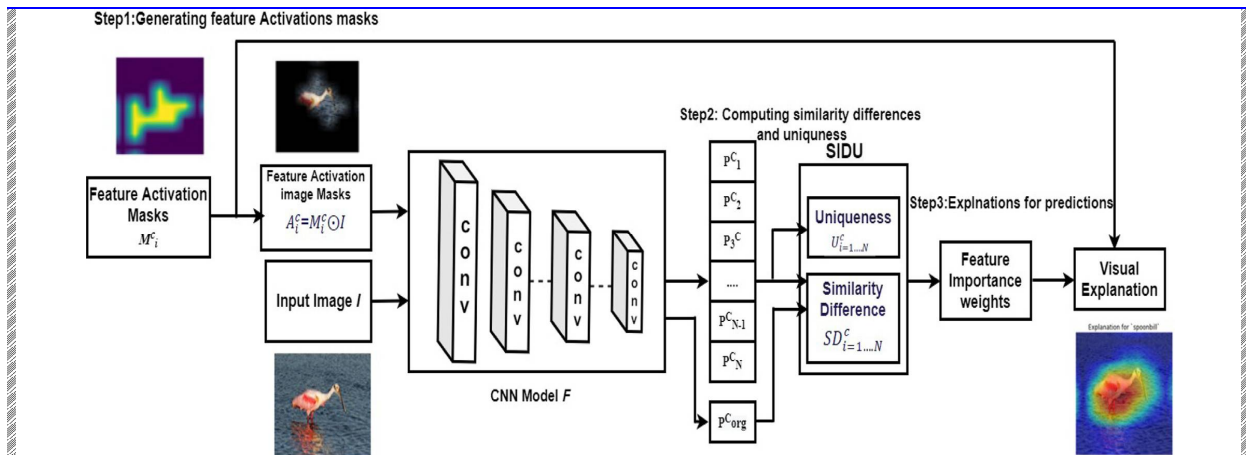
Your Answer: **Incorrect! AI's answer was A**
Press the space bar to continue...

xAI.

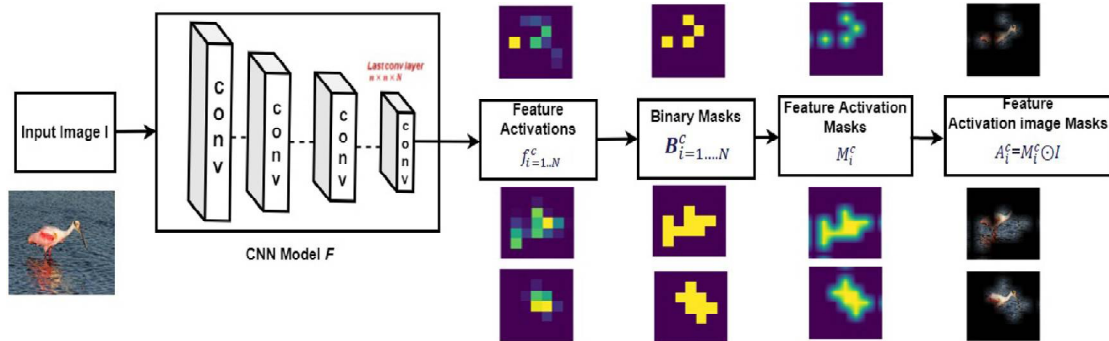
Similarity Difference and Uniqueness (SIDU)

2022-124

Block diagram of SIDU.



Generating feature image masks from last layer activation's of CNN model F

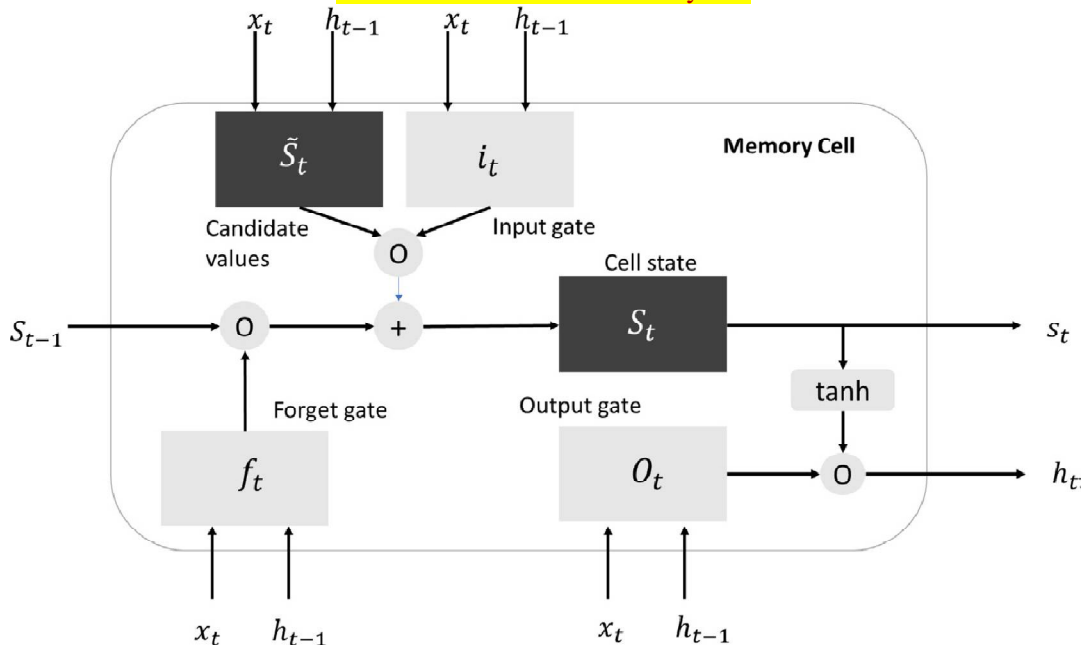


CNN model F is same of all the steps

xAI.

2022-126

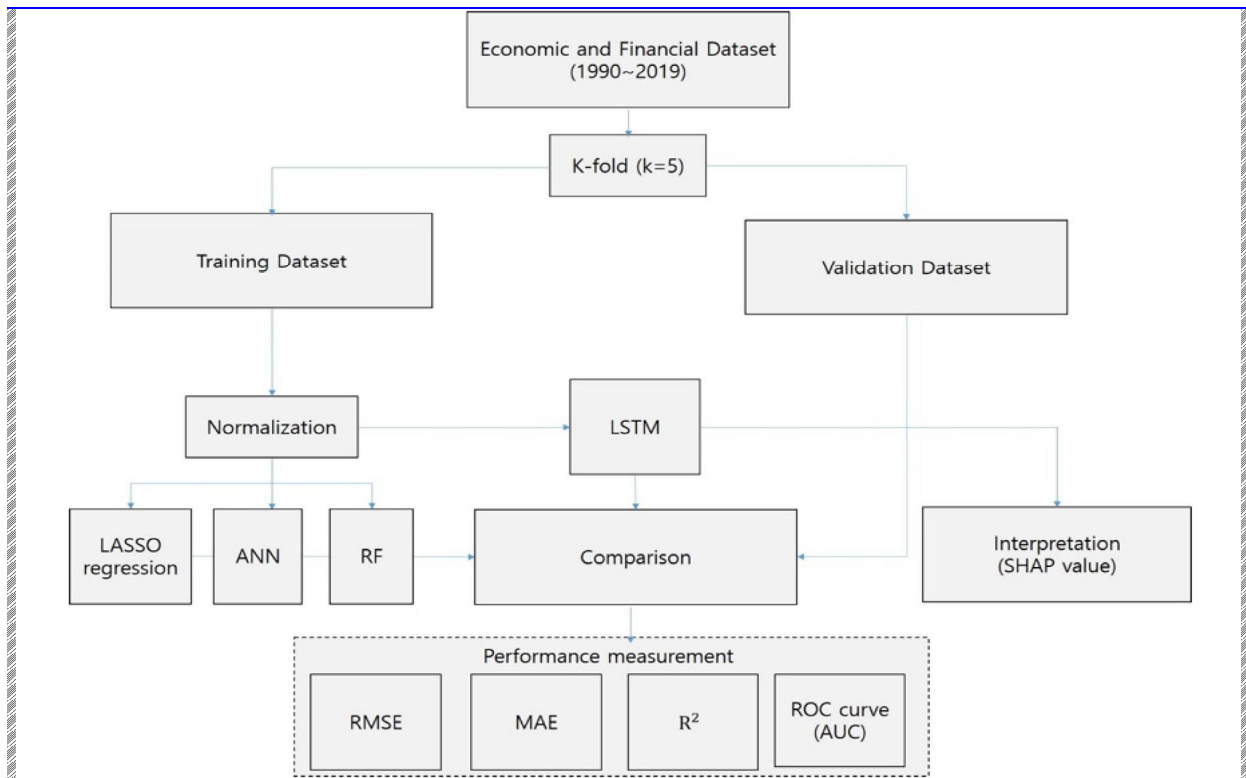
Structure of an LSTM memory cell



xAI.

2022-126

Flow diagram of prediction process in models



xAI.

2022-127

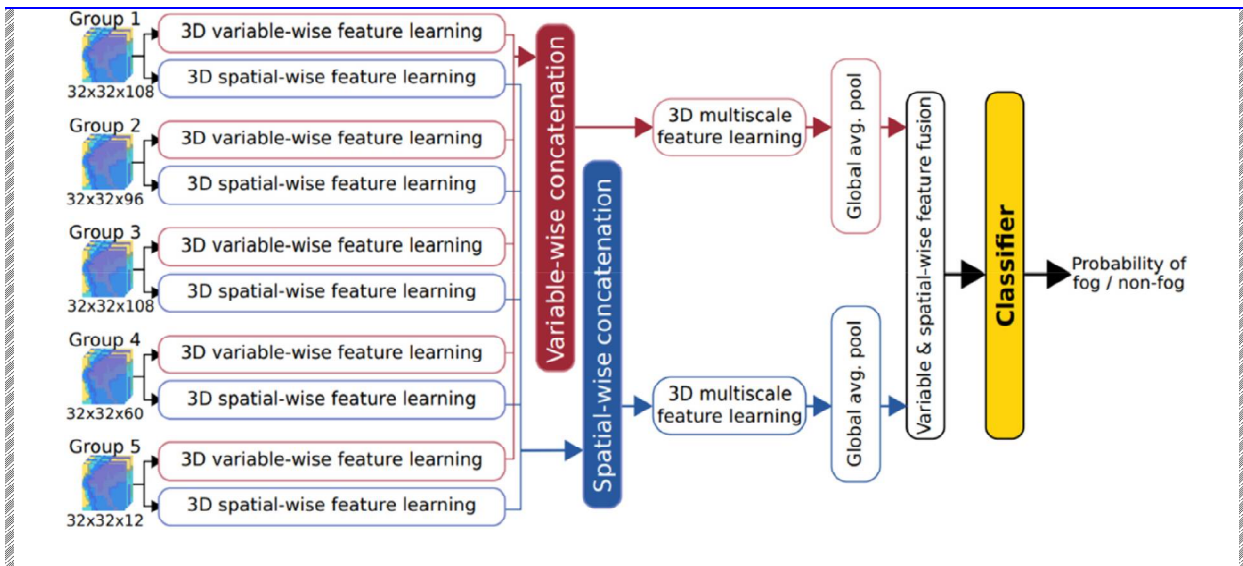
Principal hyperparameters of Denoising Generative Adversarial Network.

Parameter	Value
2D representation width	20
2D representation height	20
2D number of channels	1
Base number of filters	64
Kernel size	(3, 3)
Strides	(2, 2)
Dropout rate	0.5
Learning rate	2e-4
Max epochs	2000
Regularization L1 lambda	100

xAI.

2022-138

Overview of the FogNet3D parallel processing of features –
spatial-wise (blue) and variable-wise (red)




Animation showing the occurrences of each channel in the top-K channels
 [K : 1, 2, ..., 384]

FogNet XAI: Top Channels Animation

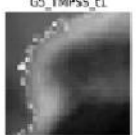
<https://gridftp.tamucc.edu/fognet/>

Channel-wise PartitionSHAP
(top 3 out of 384 channels)

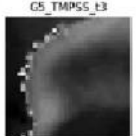
G4_WEL_050m_t1 SHAP:



G5_TMP55_t1 SHAP:



G5_TMP55_t3 SHAP:



FogNet: fog prediction model using 32x32x384 raster of atmospheric variables

XAI: what channels are influencing the model's decisions?




Channel-wise PartitionSHAP: assigns SHAP values to superpixels within channels

We then order the channels by maximum SHAP value in its cells

To get a global view of influential channels, plotting the number of times that each channel appears in the top K sorted channels across dataset samples

This animation shows the effect of increasing K from 1 channel to all (384)

Categories: hits, misses, false alarms, correct rejections

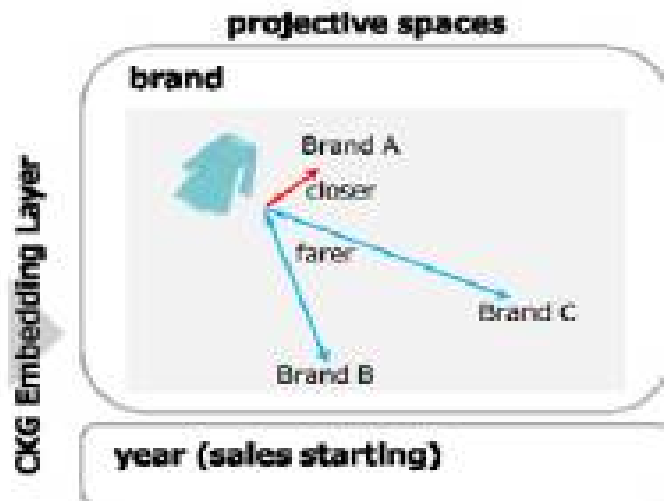
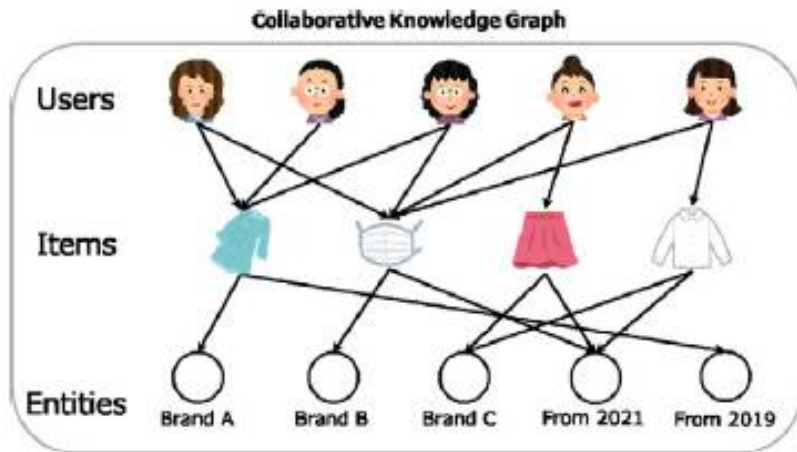


Illustration of attentive embedding propagation layer

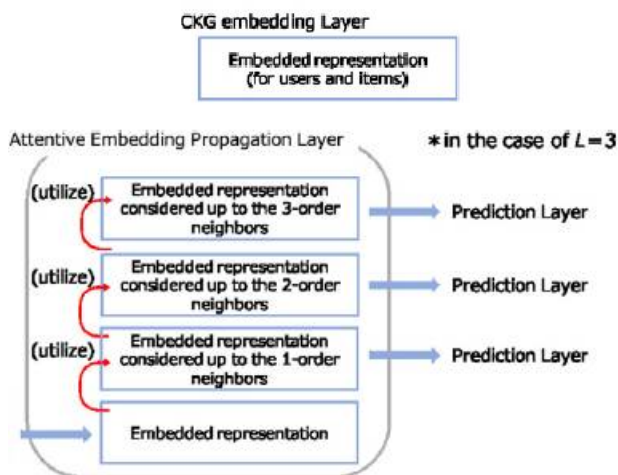


Illustration of prediction layer

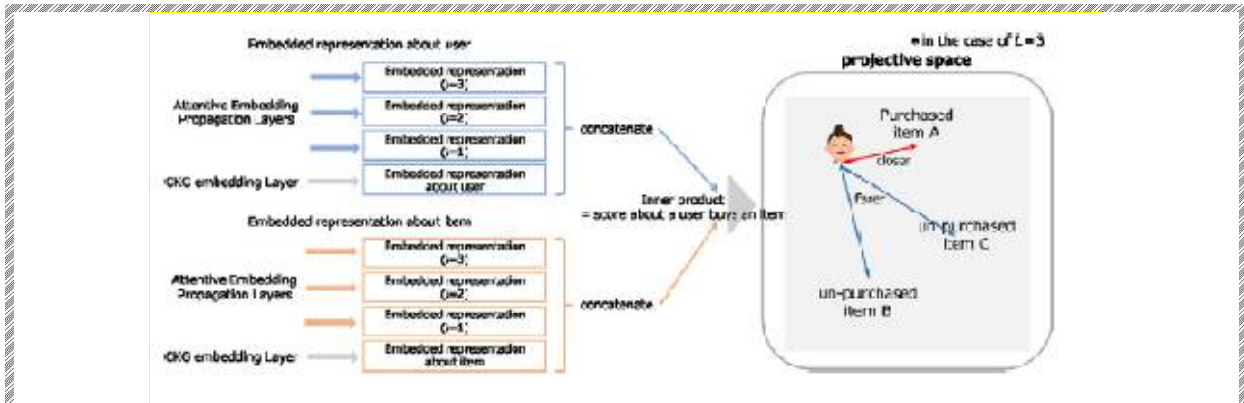
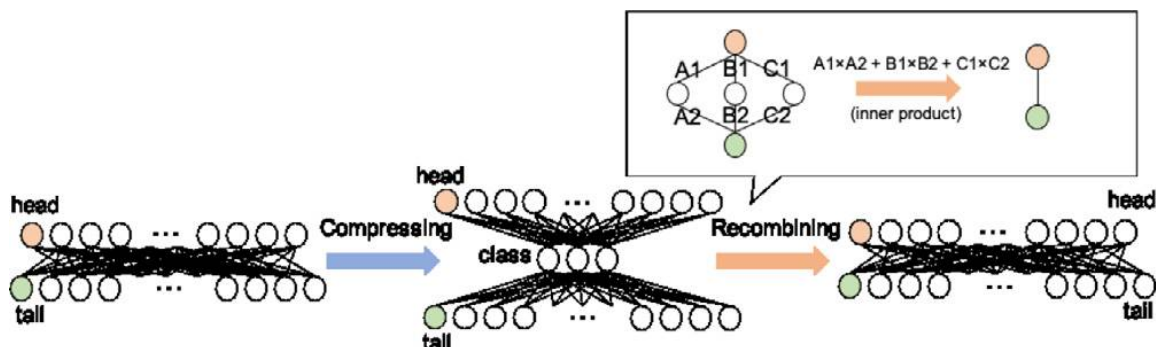
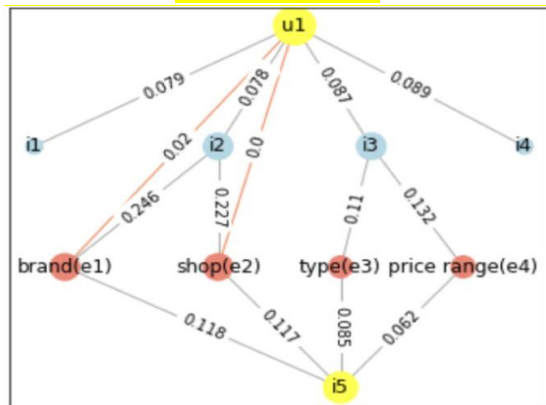


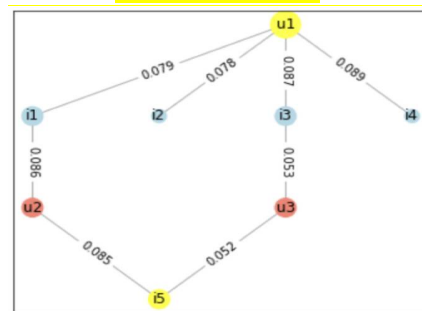
Illustration of compressing and recombining

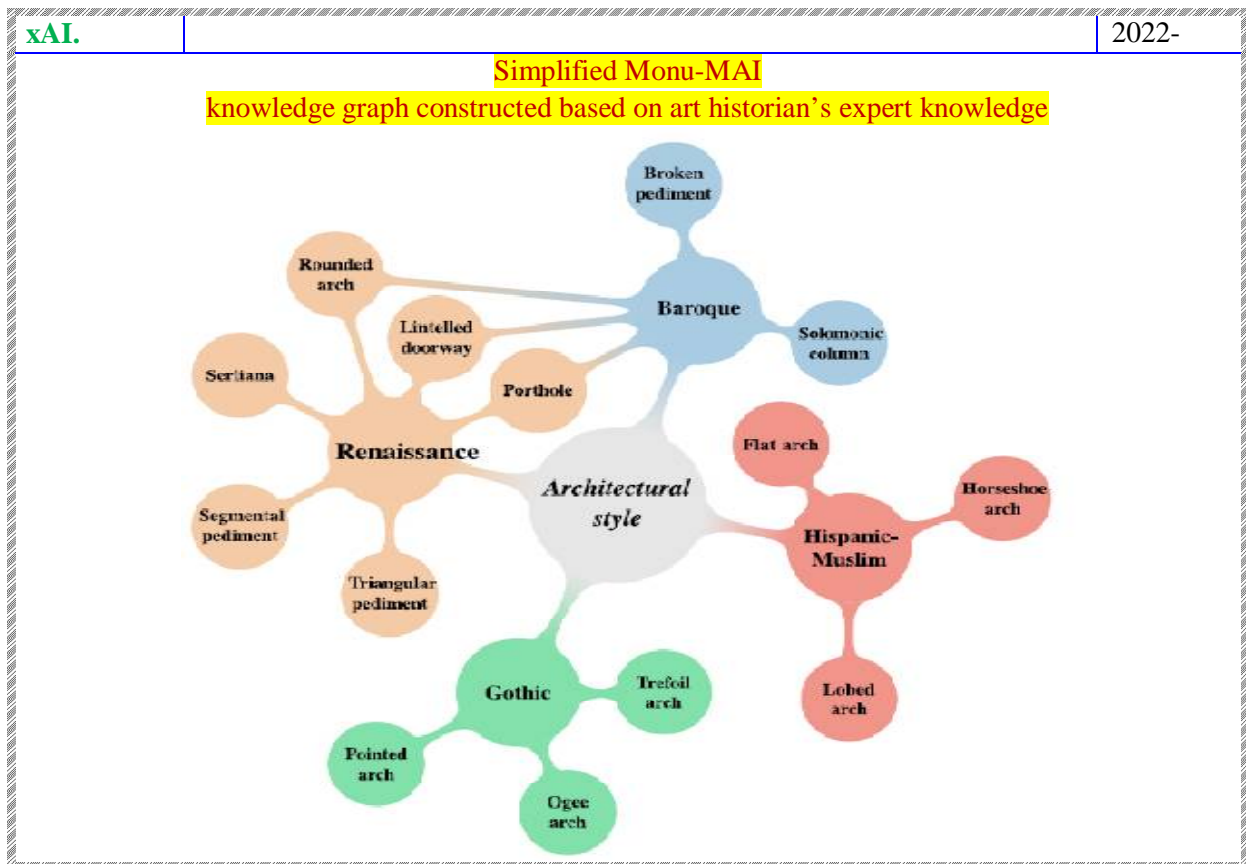
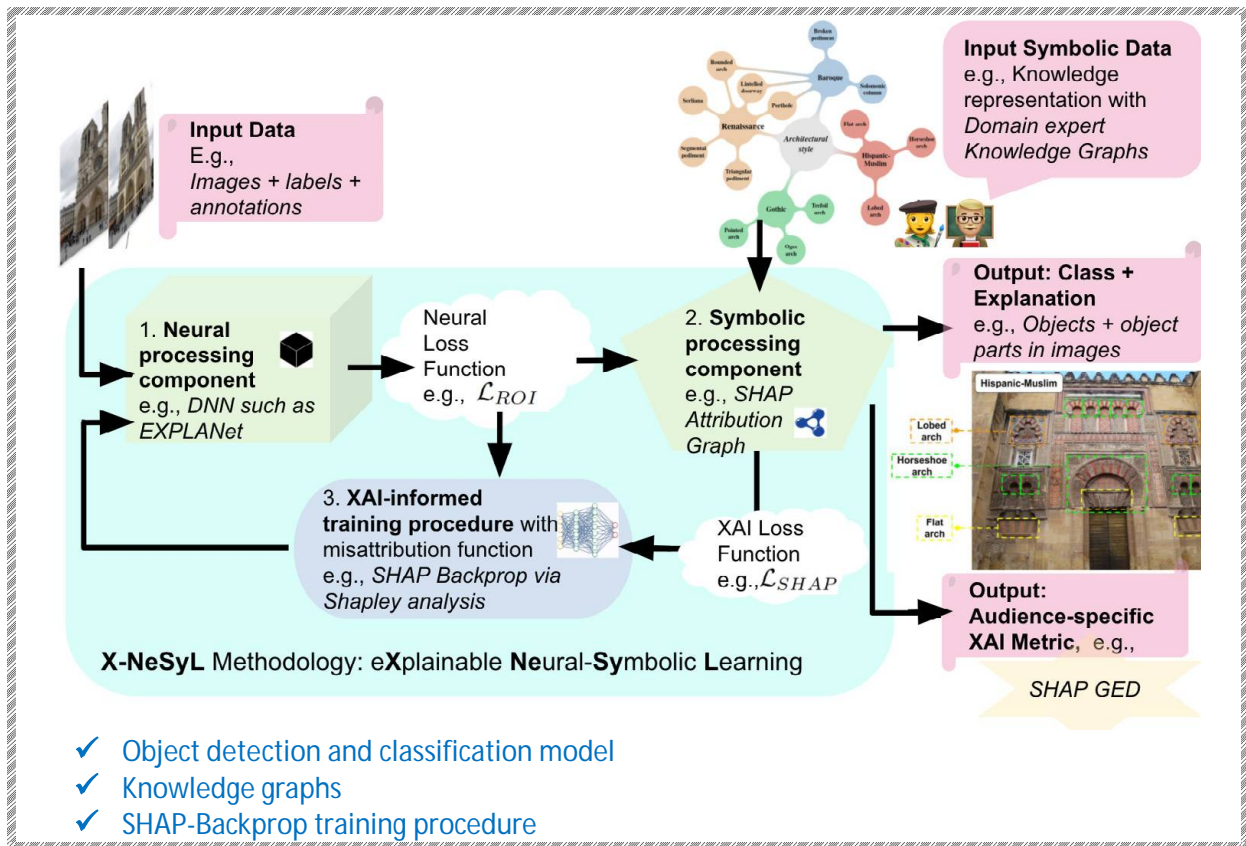


Example of a map visualizing attribute-based reason for recommendation.

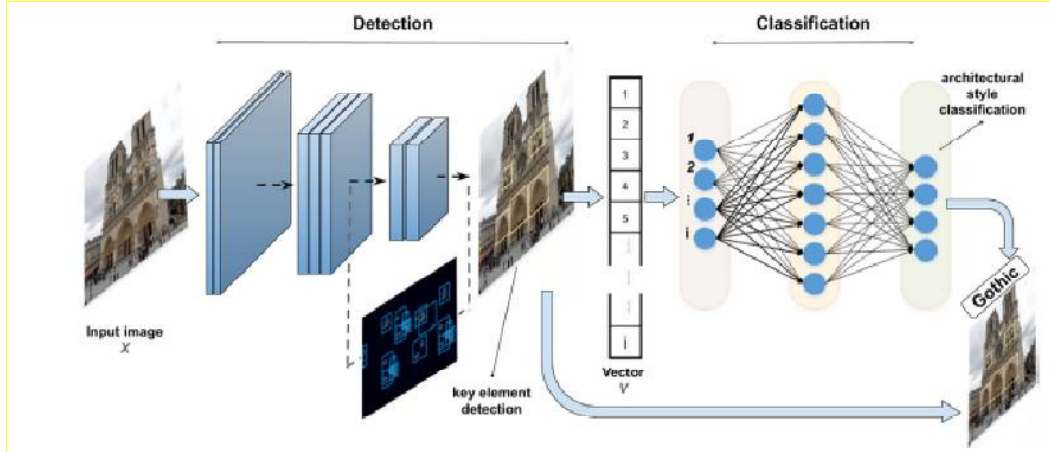


Example of a map visualizing behavior-based reason for recommendation.





EXPLANet architecture processing examples from MonuMAI dataset



Example CLEVR-XAI-complex data point

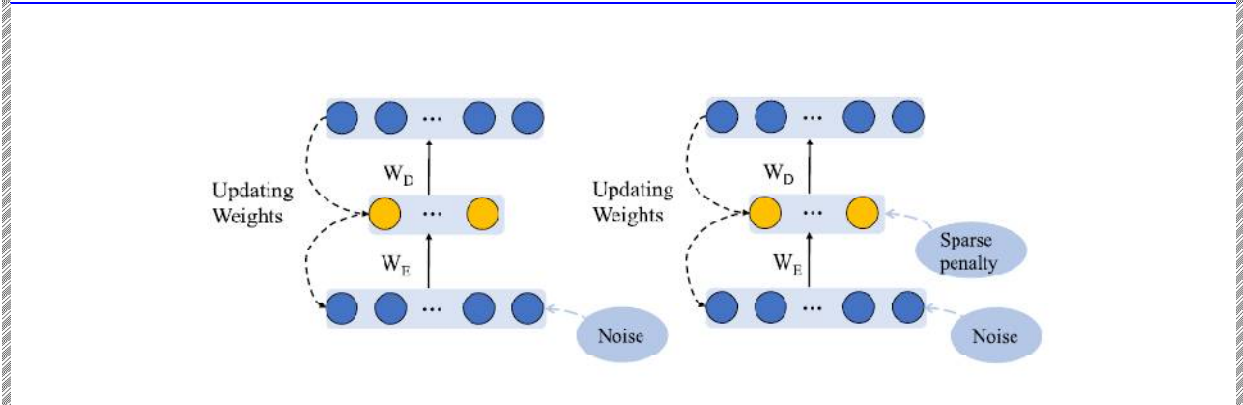
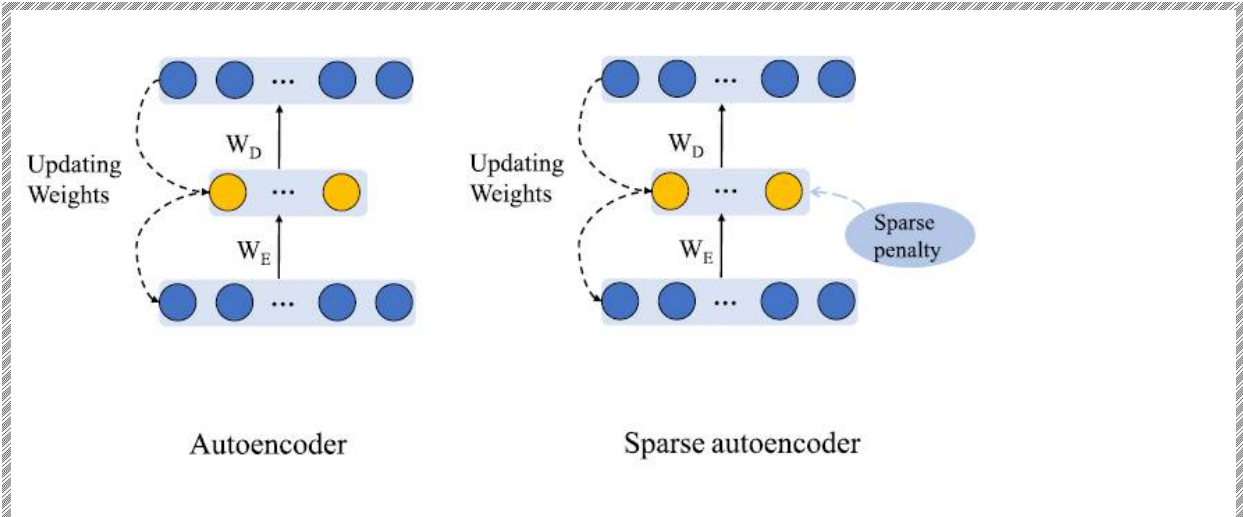
Image	Question/Answer	Program
	<p>What number of objects are large purple blocks or green metallic cubes?</p> <p>one</p>	<pre>branch1 = [scene, filter_size, filter_color, filter_shape] branch2 = [scene, filter_color, filter_material, filter_shape] program = [(branch1,branch2), union, count]</pre>
Ground Truth	GT Unique	GT Unique First-non-empty
	undefined	
	GT Union	GT All Objects

Objective: Which objects in the scene are considered as ground truths

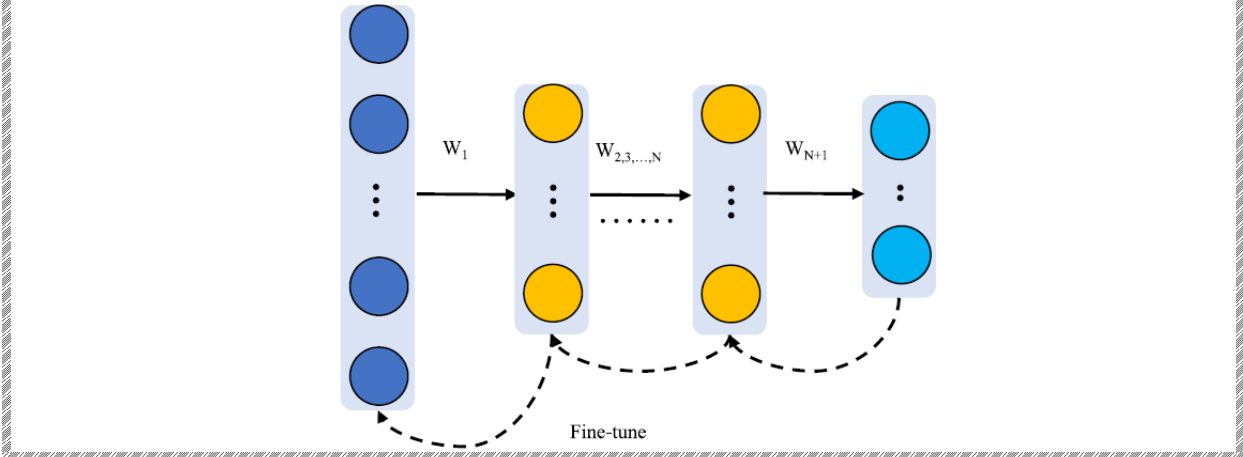
xAI.

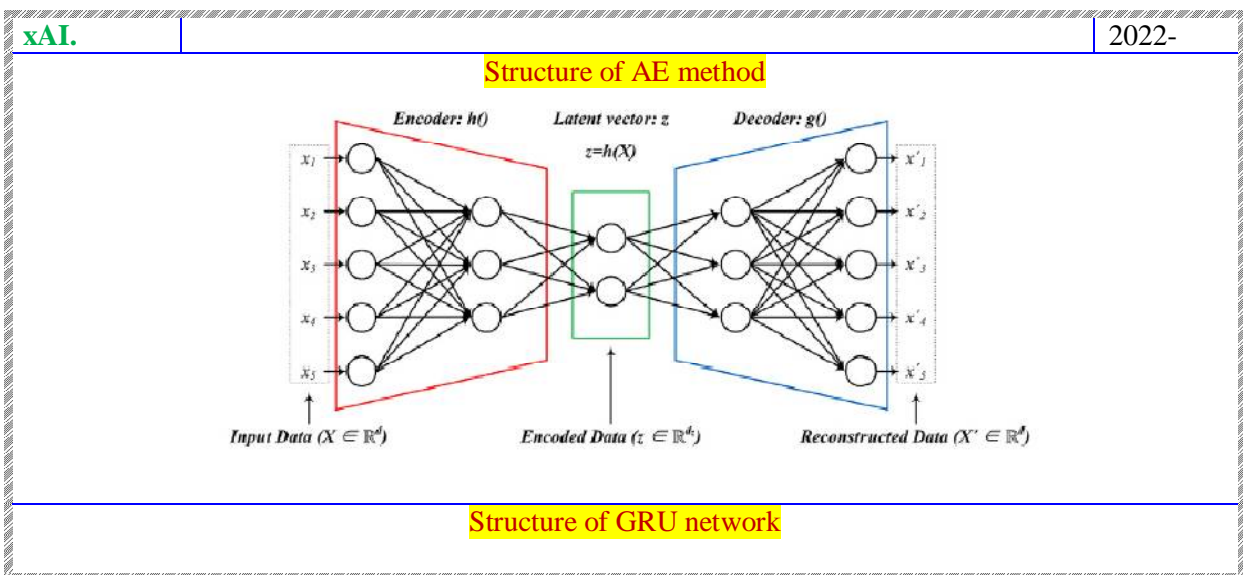
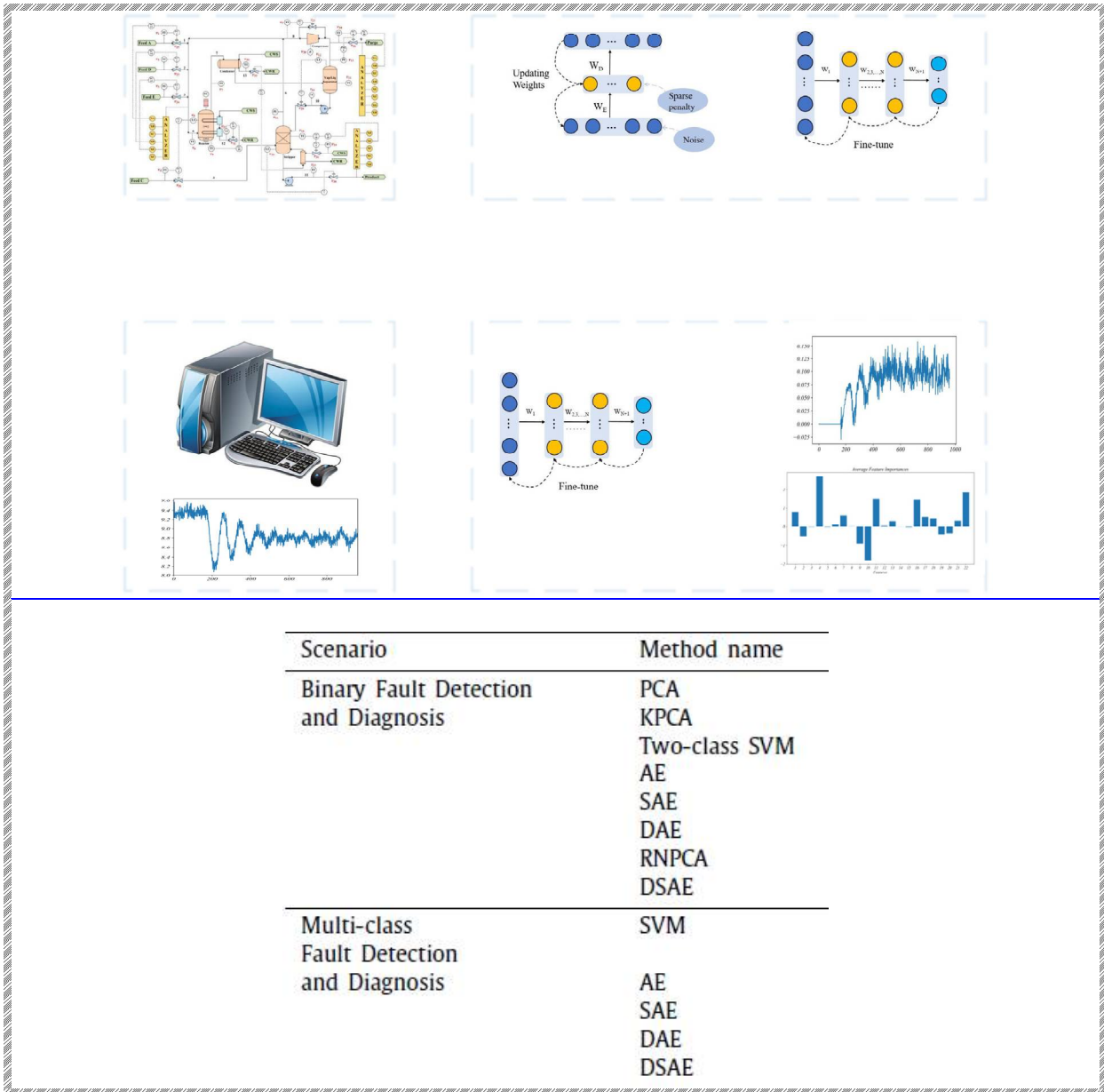
2022-154

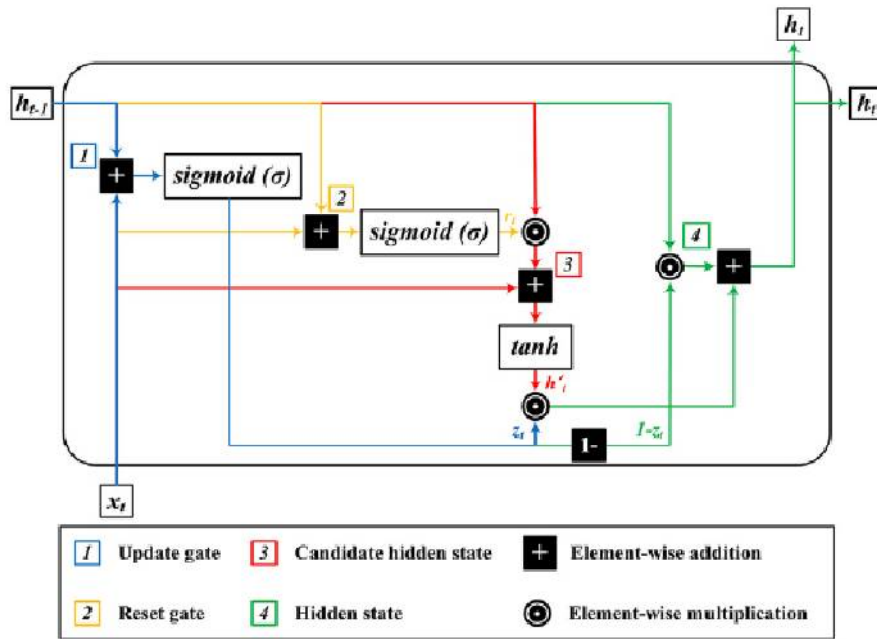
Structure of autoencoder, sparse autoencoder, denoising autoencoder, denoising sparse autoencoder (DSAE)



Supervised fine-tuning process of DSAE







$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t] + b_z)$$

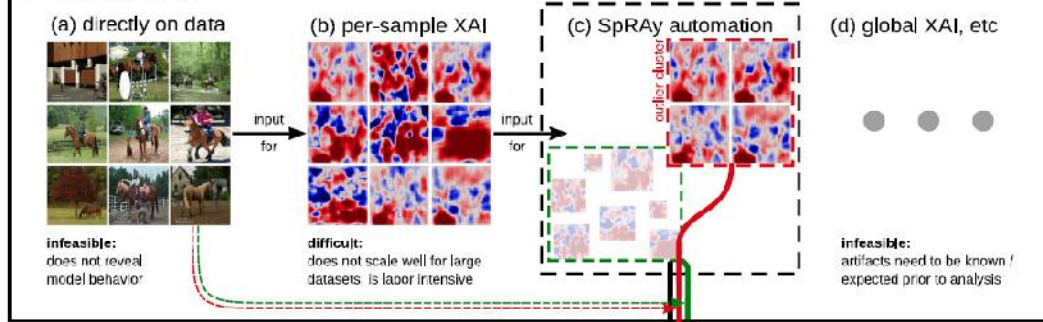
$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t] + b_r)$$

$$\hat{h}_t = \tanh(W_h \cdot [r_t \odot h_{t-1}, x_t] + b_h)$$

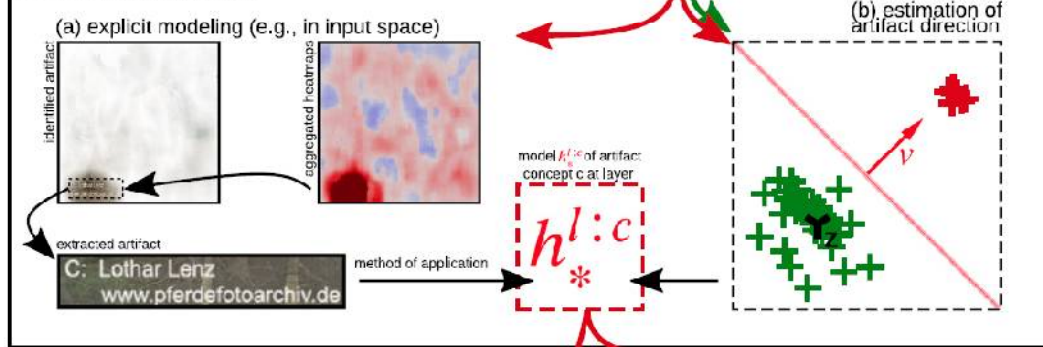
$$h_t = (1 - z_t) \odot h_{t-1} + z_t \odot \hat{h}_t$$

Class Artifact Compensation Workflow

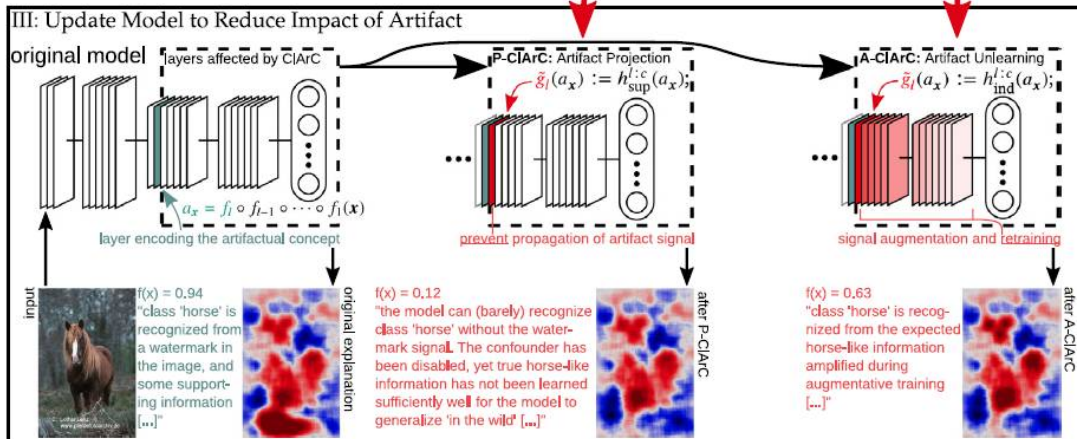
I: Identify Artifacts

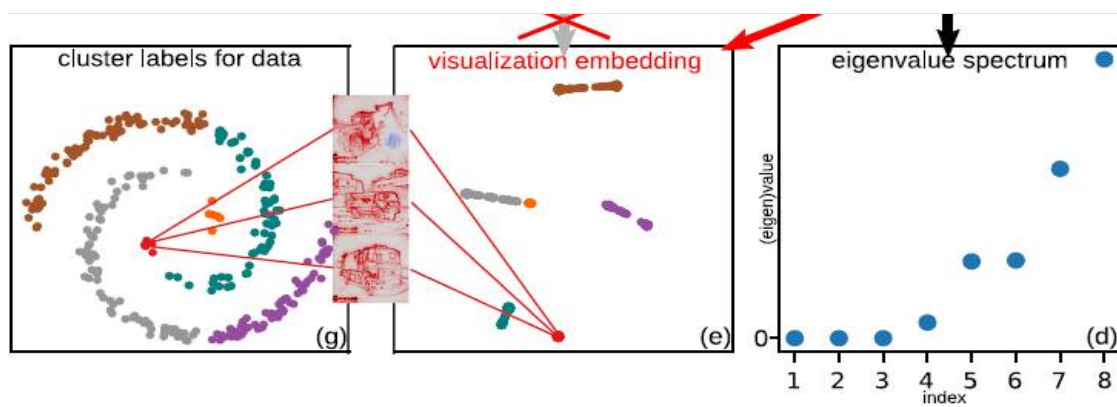
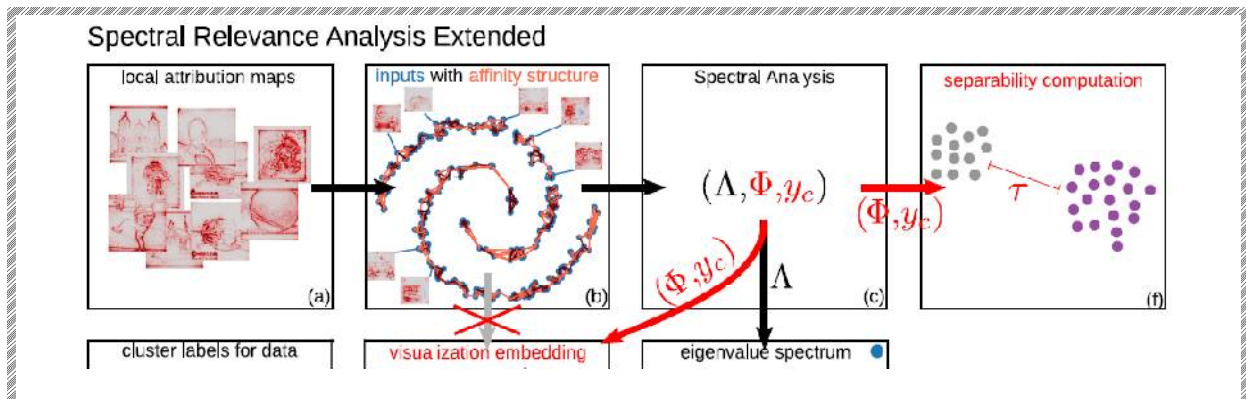


II: Estimate Artifact Model

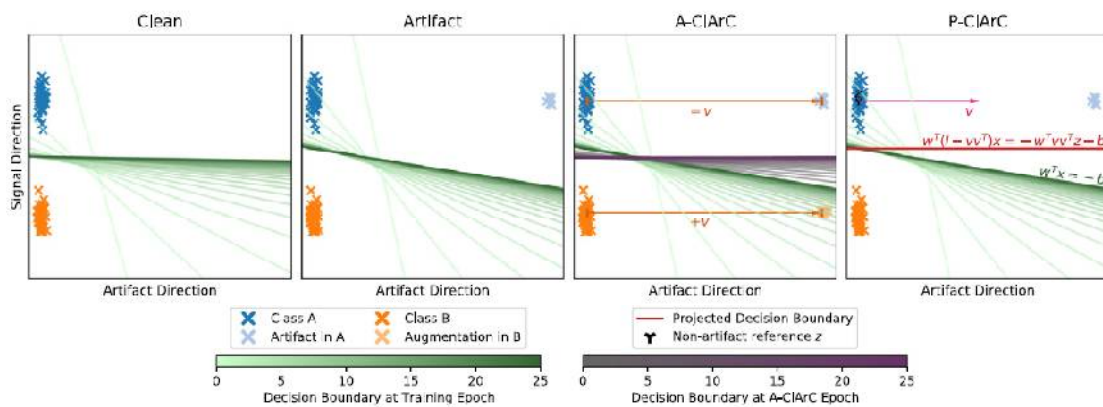


III: Update Model to Reduce Impact of Artifact

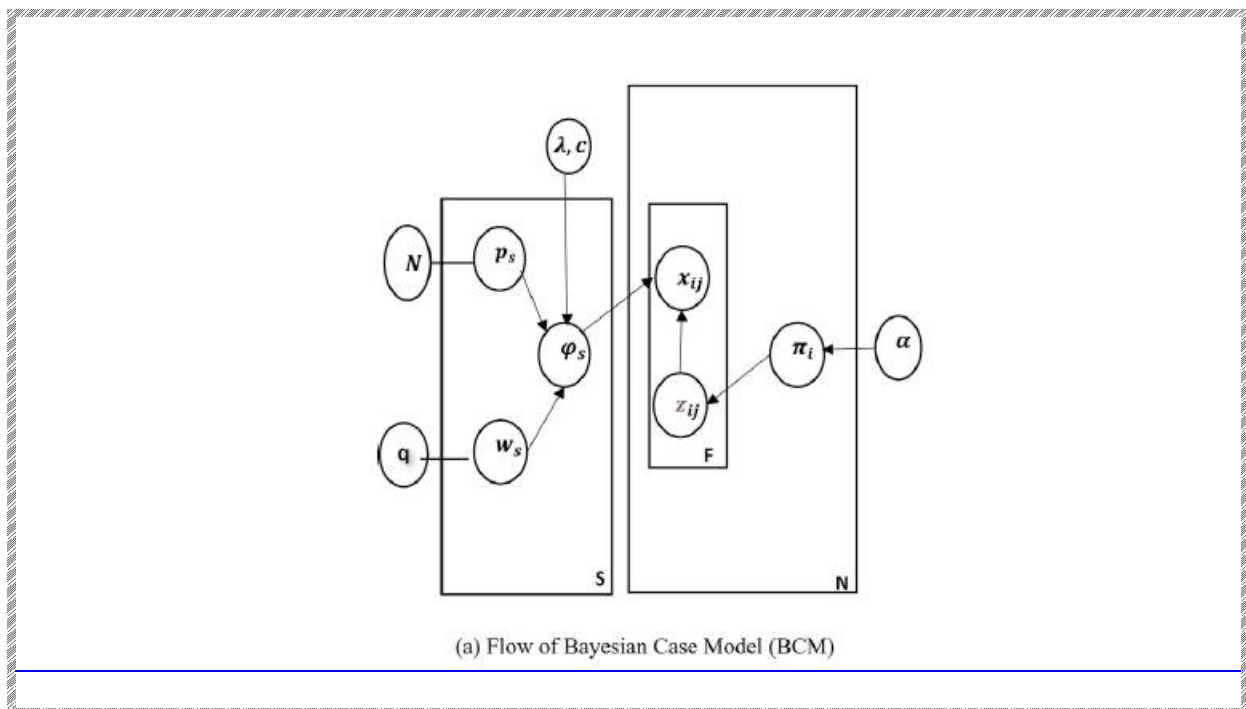
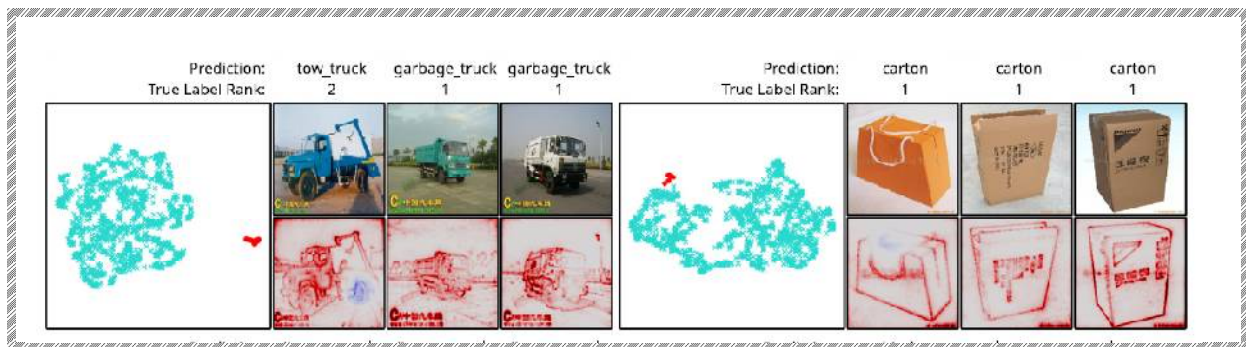
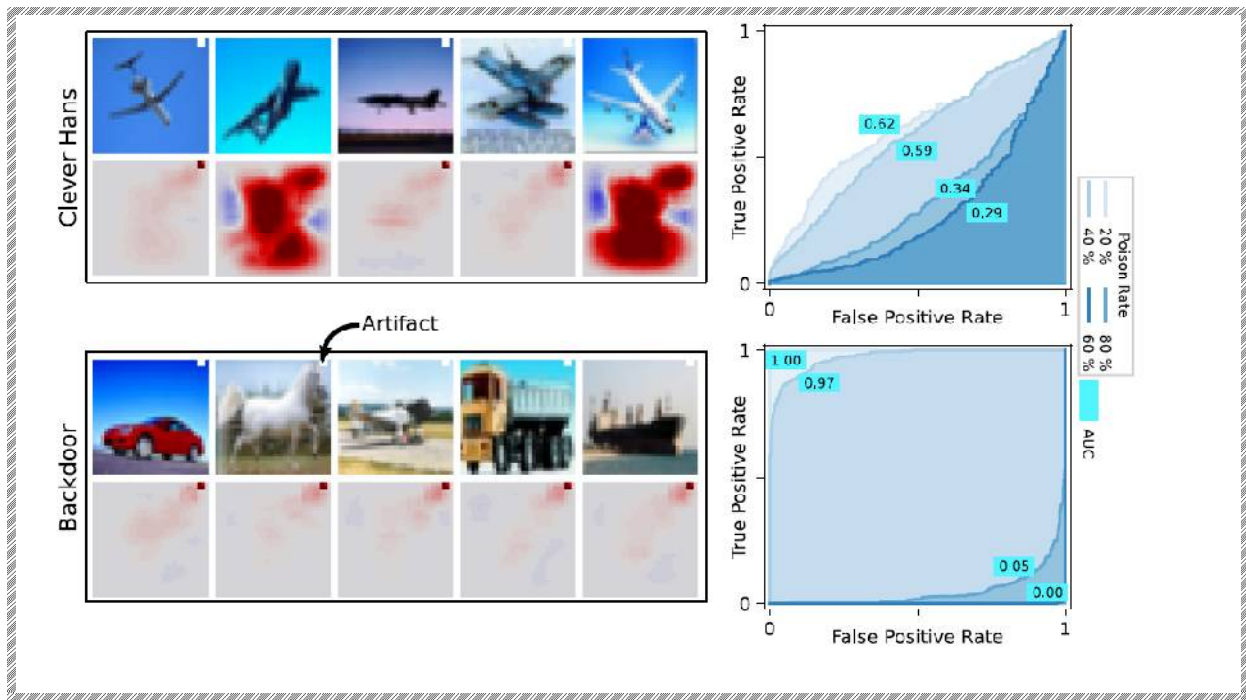




Logistic regression



Differences in detection of CH artifacts (top) and BDs (bottom)



	Data in assigned to cluster	LDA			BCM	
		Top 3 words and probabilities			Prototype	Subspaces
1		0.26	0.23	0.12		color () and shape () are important.
2		0.26	0.24	0.16		color () and eye () are important.
3		0.35	0.27	0.15		eye () and mouth () are important.

(b) BCM and LDA interpretation for dataset of smiley faces [77]