

# Towards Transparent AI-Powered Cybersecurity in Financial Systems: The Deployment of Federated Learning and Explainable AI in the CaixaBank Pilot

Aikaterini Karampasi<sup>1</sup>, Panagiotis Radoglou-Grammatikis<sup>1</sup>, Marek Pawlicki<sup>2,3</sup>, Ryszard Choraś<sup>3</sup>, Ramon Martin de Pozuelo<sup>4</sup>, Panagiotis Sarigiannidis<sup>1</sup>, Damian Puchalski<sup>2</sup>, Aleksandra Pawlicka<sup>2,5</sup>, Rafał Kozik<sup>2,3</sup>, and Michał Choraś<sup>2,3</sup>

<sup>1</sup>Department of Electrical and Computer Engineering, University of Western Macedonia, Kozani, Greece

<sup>2</sup>ITTI sp. z o.o., Poznań, Poland

<sup>3</sup>PBS Bydgoszcz, University of Science and Technology, Bydgoszcz, Poland

<sup>4</sup>Technology and Information Security, CaixaBank, Barcelona, Spain

<sup>5</sup>University of Warsaw, Warsaw, Poland

## ABSTRACT

In the domain of financial cybersecurity, where trust and reliability is paramount, the advent of Artificial Intelligence is bringing novel tools for network intrusion detection. This paper introduces AI4FIDS, a novel AI-powered Intrusion Detection System leveraging Federated Learning (FL) to enhance data privacy while enabling decentralized model training across multiple financial entities. Concurrently, we present TRUST4AI.XAI, an explainability module designed to render AI decision-making transparent and interpretable, thereby aligning with the critical need for model accountability in financial applications. Our experimental results, conducted in the framework of the AI4CYBER project's financial sector pilot, demonstrate in detecting network intrusions in financial infrastructure while maintaining user privacy, while increasing trustworthiness via explainability methods. The integration of these technologies addresses the dual challenges of effective threat detection and regulatory compliance, offering a scalable solution for modern financial institutions. This work contributes to the ongoing dialogue on leveraging AI for financial security and sets a benchmark for the development of privacy-preserving, interpretable AI models in this sector.

Keywords: Federated Learning; Network Intrusion Detection; Fintech; AI explainability

## INTRODUCTION

Cybersecurity plays a crucial role in protecting sensitive information and assuring the continuity of critical sectors' operations in today's interconnected world. It is particularly important in sectors such as the financial domain and banking services, in which trustworthiness and reliability are essential to preserve social trust and economic stability, and which are targeted by a variety of malicious actors [25]. As evolving cyber threats target critical sectors with more and more sophisticated attack attempts, investing in advanced security measures, often assisted by Artificial Intelligence (AI) capabilities, becomes essential to maintain the robustness and reliability of key services. Digitalization of the financial sector, the growth of fintech as a whole, as a part of international critical infrastructure [12] and its vulnerability due to the sensitive nature of financial data are reasons why the financial sector became one of the prime targets of cyber criminals [32]. These are also reasons for increased frequency and impact of malicious attempts targeting financial systems over the last years [2].

At the same time, recent advances in AI and successful applications of machine learning (ML) in detecting and classifying intrusions at the network level have made AI widely recognized as a major tool for enhancing the cybersecurity of banks. However, the use of AI-assisted cybersecurity poses new challenges, such as issues of model transparency and inherent ML/AI vulnerabilities [2].

In this paper, a trustworthy network intrusion detection pipeline is proposed. It includes AI4FIDS – an AI-powered Intrusion Detection System (IDS) which leverages Federated Learning (FL), enabling the

34 training of federated models across multiple decentralized entities or environments. The TRUST4AI.XAI  
35 tool providing explainability mechanisms for AI models is also introduced. It is used in conjunction  
36 with AI4FIDS to make cybersecurity decisions proposed by the AI-powered system more transparent and  
37 interpretable. Both tools are part of a wider architecture, designed, developed and implemented in the  
38 AI4CYBER [AI4CYBER] project, which is a European Union (EU) project co-funded by the Horizon  
39 Europe research and innovation programme under the Grant Agreement (GA) No 101021936. AI4FIDS is  
40 one of the core components of the AI4CYBER project together with tools for root cause analysis, attack  
41 simulation, fixing and testing, vulnerability analysis and many more. The proposed explainability module is  
42 included in the framework providing trustworthiness for AI services developed in the project. AI4CYBER  
43 tools are validated in three real-world pilots. For the purpose of this work, the focus is on the financial  
44 sector.

45 This paper outlines several significant advancements in AI-powered cybersecurity deployments for  
46 financial systems through a detailed exploration of the CaixaBank pilot. Major contributions include:

- 47 • A novel AI-powered IDS that utilizes Federated Learning to ensure data privacy across multiple  
48 financial entities while enabling decentralized model training
- 49 • An explainability module that serves an array of xAI methods, giving the user a look at the AI model  
50 from different perspectives, contributing to the transparency and interpretability of AI decision-  
51 making, thus meeting the need for model accountability in financial applications
- 52 • Provides the Experimental Validation of the end-user-centric and sector-oriented AI4FIDS and  
53 TRUST4AI.XAI deployments in the CaixaBank pilot.
- 54 • Sets a precedent for the development of privacy-preserving, interpretable AI models in the financial  
55 sector

56 The paper is organized as follows: Section II provides the current state of the art related to federated  
57 learning in cybersecurity and to explainable AI (xAI). Section III details the approach and design of  
58 AI4FIDS and TRUST4AI.XAI solutions and CaixaBank use-case in which we validate the proposed tools.  
59 Section IV focuses on the experiments and experimental results obtained to prove the effectiveness of the  
60 proposed suite of tools in the banking pilot. Section V concludes this article.

## 61 STATE-OF-THE-ART AND RELATED WORKS

### 62 Federated Learning for Network Intrusion Detection

63 The impact of Federated Learning in cybersecurity has been considered by a variety of research studies,  
64 with particular focus on its application in intrusion detection and prevention. In the following paragraphs,  
65 an overview of relevant survey papers in this field will be discussed. More notably, M. Alazab et al. [3]  
66 evaluated the manner in which FL operates and contributes in the context of cybersecurity, with a targeted  
67 analysis of selected use case scenarios, applications and confrontations. In the same manner, B. Ghimire  
68 and B. Rawat [15] explore the progression of FL and cybersecurity in a reciprocal fashion. The authors  
69 initially investigate the utilization of FL in cybersecurity applications, including but not limited to intrusion  
70 detection, with peculiar interest on Internet of Things (IoT) and Cyber-Physical Systems (CPS), while they  
71 subsequently discuss the impact of cybersecurity in FL. In an extensive review paper, E. M. Campos et  
72 al.[8] investigate the manner in which FL is employed for IoT environments, exploring the effect of FL in  
73 intrusion detection while also considering the progression of Machine Learning and Deep Learning (DL)  
74 approaches by reason of FL. Eventually, areas with potentials for additional exploration and avenues for  
75 future studies are described. A comprehensive survey is provided by L. Lavour et al. [19] regarding the  
76 evolution of federated IDS and Intrusion Prevention System (IPS). After establishing their methodological  
77 approach, the authors conduct a detailed analysis of the existing research, evaluating a range of criteria  
78 including the following: detection techniques, mitigation tactics, data sources and datasets, variations  
79 of FL, local models and aggregation methodologies, as well as communication protocols (e.g. overhead  
80 optimisation and encryption procedures). In light of the aforementioned criteria, a pertinent categorization  
81 of federated IDS and IPS is proposed, and a comparative analysis of the subject literature is undertaken.  
82 Subsequently, the authors present a discussion of the open issues and research directions that remain to  
83 be addressed. S. Arisdakessian et al. [7] conducted a survey regarding intrusion detection with respect  
84 to IoT applications. The authors combined and discussed a variety of technological and research areas,  
85 including FL, game theory, social psychology and Explainable Artificial Intelligence. A total of 19 criteria  
86 were taken into consideration in order to conduct an exhaustive study and analysis of several works which  
87 allowed the identification of significant research gaps pertaining to the aforementioned technological and  
88 research areas.

89 A federated DL framework is presented by S. I. Popoola et al. [26] which consists of multiple distinct  
90 nodes performing the training procedure of Deep Neural Networks (DNNs) based on the data provided  
91 by their local network traffic. Subsequently, a central server gathers the disparate parameters of each  
92 trained model and aggregates them using the Fed+ fusion technique, eventually distributing them back to  
93 the nodes. With respect to the DNNs' designation, these incorporate the input layer, two fully connected  
94 hidden layers and the output layer. The results that were attained, based on simulations that were executed,  
95 provided the authors with an accuracy of 99.27%, precision 97.03%, TPR 98.06%, and an F1-score of  
96 97.50%. This outcome indicates better performance of the federated DL over the local DNN models. With  
97 regard to the identification of the optimal fusion technique, a variety of methods were employed, namely  
98 Federated Averaging (FedAvg), Fed+, and Coordinate Median (CM). The experiments indicated that Fed+  
99 exceeded the performance of the other state-of-the-art (SOTA) methods, providing evidence for the overall  
100 superiority of the DNN-Fed+ model using FL for the intrusion detection assignment in heterogeneous  
101 wireless networks.

102 T. Dong et al. [11] proposed a novel intrusion detection system based on a learning-based methodology,  
103 namely FedForest, encompassing FL and Gradient Boosting Decision Trees (GBDT). The proposed  
104 technique is implemented by training a local encoder (GBDT classifier) on the distinct clients. The data  
105 based on which the clients were trained, were distinct private datasets of each client, while the parameters  
106 that were attained in each case were broadcast to the server. Consequently, the server decides for the finest  
107 encoders and transmits them to all clients. Eventually, the clients utilize the encoders to encode their data,  
108 train and deploy the new models. To further enhance data privacy, a random masking algorithm was utilized  
109 on the data. During the evaluation procedure, an illustration of the superiority of the proposed FedForest  
110 was performed with a Multi-layer Perceptron (MLP) composed of 3, 5, and 7 layers. The results that were  
111 attained indicated accuracy levels of 67.03% on the DDoS2019 dataset, 89.63% on MalDroid2020, 86.76%  
112 on Darknet2020, and 79.6% on DoHBrw2020, signifying the prevalence of the suggested methodology.

113 P. H. Mirzaee et al. [22] suggested a Federated Intrusion Detection System (FIDS) methodological  
114 scheme specifically implemented for 5G environments, the primary goal of which was to establish  
115 user privacy while simultaneously preserving a high detection rate. More notably, a federated DNN  
116 implementation was proposed, appropriate for ensuring privacy of the user's information. The algorithm  
117 encompassed a dedicated server for the aggregation of the updates from each respective local model, while  
118 the obtained results were sent back to the end nodes. Regarding the evaluation procedure, it was exhibited  
119 that the recommended implementation accomplished 99.5% in all metrics, namely accuracy, precision and  
120 F1-score, on the NSL-KDD dataset.

121 W. Schneble and G. Thamilarasu [29] proposed a widely distributed IDS based on ML methodologies,  
122 and more specifically they employed FL techniques for Medical Cyber-Physical Systems (MCPS), towards  
123 reducing communication and computation costs while increasing network security. The suggested model  
124 was evaluated on both real and simulated attacks such as Denial of Service (DoS), Data Modification, and  
125 Data Injection. The results that were attained showcased that the model under discussion outperformed  
126 SOTA methodologies by achieving 99% accuracy levels and an FPR of 1%, while simultaneously the  
127 communication costs were decreased.

128 O. Aouedi et al. [5] suggested the FLUIDS which describes a semi-supervised implementation for IDS,  
129 composed of encoders and trained on each end device with unlabeled data. The local models are afterwards  
130 aggregated to be trained on labeled data, which is located on a server, in a supervised manner, eventually  
131 providing an ameliorated classification of attacks. B. Li et al. [20] introduced a federated IDS especially  
132 trained on detecting DDoS attacks based on prototypical features extracted by GRU layers to eventually  
133 derive 97% accuracy. On the other hand, R. Zhao et al. [33] implemented an FL architecture based on  
134 BiLSTM towards identifying high-risk malicious behaviour, which had minor variations from a centralized  
135 model, to obtain 99.21% accuracy. The IoTDefender was proposed by Y. Fan et al. [13] for 5G IoT through  
136 a federated transfer learning architecture which surpassed the performance of traditional implementations  
137 achieving 91.93% accuracy and finer generalization abilities. O. Friha et al. [14] proposed the FELIDS  
138 framework, which was based on CNN and DNN architectures towards constructing an FL-based IDS  
139 model, outperforming other centralized architectures in maintaining privacy of the utilized data as well as  
140 high detection accuracy.

### 141 **Explainable AI Methods for Network Intrusion Detection**

142 Many of the best-performing AI/ML methods function as black boxes, which presents significant ethical  
143 concerns for their use in various domains. This lack of transparency can undermine trust and become an  
144 obstacle for numerous potentially beneficial applications [9]. This creates a need for reliable interpretability  
145 methods, which would facilitate a way for the human operator to understand the decision-making process  
146 of the model. With this pressing necessity, xAI is now an intense area of research, with numerous emerging  
147 approaches [24].

148 In this work, the focus is on providing xAI-derived explanations, relying on the representation of AI  
149 models using methods that are easier to interpret. These methods can be derived from the original AI  
150 models or built from scratch using the available data. The techniques can be broadly classified into two  
151 categories: model-agnostic and model-specific methods [27]. The methods that are model-agnostic can be  
152 used for any ML model. The tool deployed in the project provides a plethora of xAI methods, including  
153 LIME, SHAP, DiCE and ProtoDash.

154 LIME (Local Interpretable Model-Agnostic Explanations) generates locally faithful explanations by  
155 fitting an interpretable model to the neighborhood of the input data, providing insights into complex models  
156 by creating simpler, locally interpretable linear models [27].

157 SHAP (SHapley Additive exPlanations) assigns a value to each input feature based on its contribution  
158 to the model's prediction using Shapley values to fairly allocate the payoff among features [21].

159 DiCE (Diverse Counterfactual Explanations) offers counterfactual explanations by synthesizing data  
160 points by perturbing features of a sample until the label flips, essentially presenting a 'what-if' scenario [23].

161 ProtoDash identifies 'Prototypical Samples' within a dataset to gain insights into the characteristics  
162 of a subset or specific class of data through its most representative samples. The samples are found by  
163 maximising similarity for a concise representation [18].

164 The ANCHORS explainer identifies "anchors," or rule-based conditions that reliably predict the same  
165 outcome when met. It reveals key factors influencing a model's decisions by testing feature combinations  
166 to find those that consistently lead to the same prediction. ANCHORS is model-agnostic [28].

167 PDP (Partial Dependence Plot) is a visualization tool that shows how a feature impacts the predicted  
168 outcome, on average, across a dataset. It highlights the global effect of a single feature, ignoring interactions  
169 with others, and helps interpret complex models by offering a clear graphical view of feature influence on  
170 predictions [17].

171 ICE (Individual Conditional Expectation) plots show how predictions change for each individual  
172 instance as a feature is varied, revealing interactions and variability in the model's behaviour for specific  
173 data points [16].

174 ALE (Accumulated Local Effects) plots calculate a model's predictions within intervals of a feature,  
175 accumulating these effects across the feature range. This offers a realistic view of how changes in the  
176 feature influence predictions [6].

177 Permutation Feature Importance (PFI) is a simple, model-agnostic method that measures the impact of  
178 each feature on model performance by shuffling its values. By disrupting each feature and observing the  
179 change in accuracy, PFI provides clear insights into feature relevance, making it broadly applicable and  
180 easy to interpret, regardless of the model structure [4].

## 181 **AI-EMPOWERED ANOMALY DETECTION IN BANKING SCENARIO**

### 182 **AI4FIDS - Anomaly Detection Tool**

183 AI4FIDS is illustrated through the C4 model which demonstrates the architecture of the proposed im-  
184 plementation. In general, the C4 model is a structural representation extensively employed in software  
185 engineering for conceptualizing and substantiating the architecture of the software systems. Context, Con-  
186 tainers, Components, and Code are the pillars of this model, proposed by Simon Brown, to substitute the  
187 distinguishable levels of abstraction in the model, each of which offers a distinct perspective of the system,  
188 thus making it more feasible to comprehend and explain the architecture to the involved stakeholders, both  
189 technical and non-technical. More specifically, Fig. 1 depicts the Context level of AI4FIDS, illustrating its  
190 communication with the other entities, interior and exterior.

191 To begin with, as an IDS based on multiple data sources, AI4FIDS gathers its input from captured  
192 network traffic, system logs and operational data which are captured in the Critical System under inspection  
193 by AI4FIDS. This data originates from the connections of the Critical System with the End Users and/or  
194 External Networks/Systems (i.e., the Internet). The purpose of AI4FIDS is to analyse this data and  
195 identify pertinent cyberattacks and anomalous behaviour. Upon producing its results, AI4FIDS spreads the  
196 equivalent security events to the Security Information and Event Management (SIEM), which is an external  
197 system. SIEM is mainly responsible for improving the security posture of an entity by provisioning for  
198 real-time perceptibility regarding security incidents and threats, hence assisting in the regulations' and  
199 policies' compliance actions of the entity. More notably, the core concept behind a SIEM system is to  
200 normalize, prioritize and correlate the events provided as input from AI4FIDS, while eventually the System  
201 Security Operator is able to supervise, evaluate and determine the AI4FIDS security events analysed by  
202 SIEM.

203 Even though the Context level offers a comprehensive examination of the system's engagement with  
204 exogenous entities, the Container level provides a more detailed analysis of the architectural underpinnings  
205 of AI4FIDS. More specifically, in the Container level, the communication of the distinct entities, identified

206 as logical, that comprise AI4FIDS are described, while the interaction with extrinsic components is  
 207 depicted. To that end, the following containers might be identified while constructing the FIDS in AI4FIDS:  
 208 (a) Log-based (L-FIDS), (b) Operational Data-based (O-FIDS), (c) Network Flows-based (N-FIDS), (d)  
 209 Visual-based (V-FIDS) and (f) Training for FIDS (T4FIDS). Then, in the architectural representation of the  
 210 system, on top of the Context and Container levels, the Component level is able to specifically define the  
 211 architecture of the distinct components along with their communications.

## 212 TRUST4AI.XAI - AI Explainer

213 For the AI explainability purposes, the modular architecture leveraging microservices is proposed for  
 214 the TRUST4AI.XAI Explainer. The modularity and microservice approach ensures the scalability, main-  
 215 tainability, and flexibility of the solution. The system allows the end user to perform analyses of any  
 216 supervised learning model with minimal setup. This process is based on communication between the  
 217 xAI components of the TRUST4AI.XAI system and the AI models, integrated via REST API or via  
 218 Apache Kafka, depending on use case needs. The entire TRUST4AI.XAI system is an arrangement of  
 219 microservices, as illustrated in Fig. 2. The xAI serving component features a user-friendly, React-based  
 220 front-end. The main function of the component, as it is provided by the frontend, is to easily allow the  
 221 user to perform analyses and visualize the decision-making process undertaken by the AI classification  
 222 tools in graphical/chart form. It is important to note that the AI models are external to the xAI system.  
 223 Another component, as seen in Fig. 2, is API Gateway, written using the Spring framework in Java. It  
 224 serves as a gateway between the frontend and microservices. The API Gateway orchestrates the local  
 225 and global microservices, the preprocessing microservices, and the data sink. The user, using the web  
 226 application/frontend can issue requests to explain particular samples for a particular model, choosing from  
 227 the samples visible on an APACHE Kafka topic. Those samples are the concatenation of the feature vector  
 228 and the classification results coming from the AI4FIDS component, which are pushed to the xAI Kafka  
 229 topic.

230 The microservices are responsible for creating respectively both local and global explanation objects  
 231 and visualizing their analysis. The interdependencies of these components have been illustrated in Fig. 2.  
 232 The architecture also includes components such as a service that collects logs from individual microservices,  
 233 and a configuration centre.

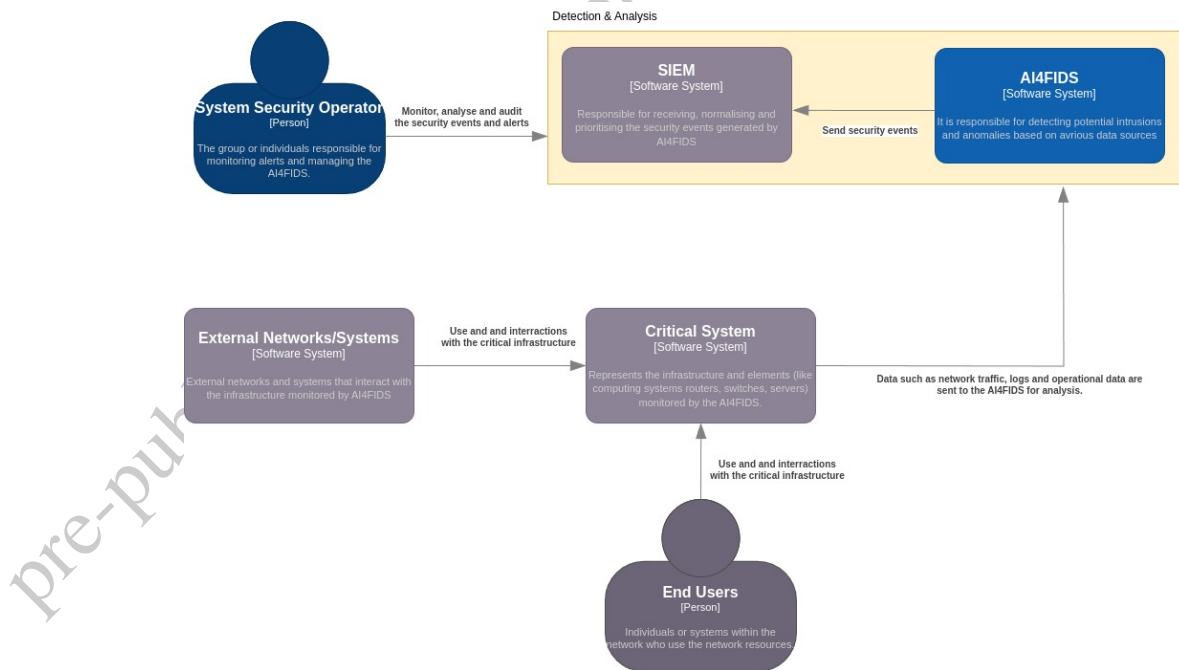
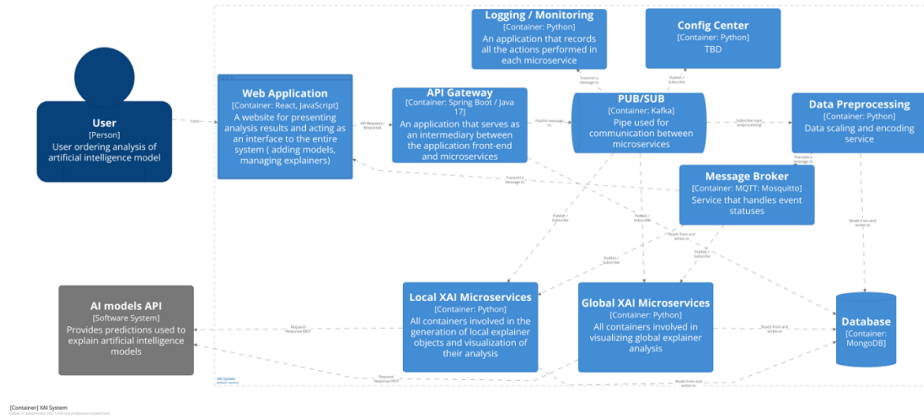


Figure 1. AI4FIDS system context model.

## 234 CaixaBank Scenario

235 The scenario defined for the validation of the proposed pipeline of AI4FIDS and TRUST4AI.XAI solutions  
 236 is rooted in the AI4CYBER project finance sector cybersecurity use case. The piloting partner is CaixaBank,  
 237 one of the leading financial institutions in Spain, chosen as the "Best Bank in Western Europe 2024" [Daly].  
 238 CaixaBank employs a set of both proprietary and commercial cybersec solutions. Tools for vulnerability



**Figure 2.** The architecture of the TRUST4AI.XAI subcomponent in AI4CYBER.

239 management, static and dynamic application security testing or pen-testing analysis are used to protect  
 240 the bank's infrastructure and for monitoring of security controls in the offered and used applications.  
 241 CaixaBank's main focus is strengthening its defence mechanisms to be used against sophisticated network  
 242 cyberattacks, as those attacks could result in the exposure of sensitive and personal data. Additionally,  
 243 the bank aims at the detection and management of software vulnerabilities in a dynamic and evolving  
 244 environment.

245 At CaixaBank, Identity and Access Management (IAM) is carried out with the suite of AIM/PAM  
 246 (Access Identity Management/Privileged Access Management) tools designed to manage user access to  
 247 different applications and services. From a security perspective, it is also important to note that bank  
 248 applications, systems, and data are accessed not only by the bank's employees, but some of the environments  
 249 are also accessible by third-party providers. Another important issue related to security is remote work,  
 250 which impacts access control configuration to ensure employees can securely reach internal applications  
 251 from their homes. The bank customizes different layers of IAM depending on the user type. As for  
 252 log collection and monitoring, and incident response mechanisms, the bank utilizes SIEM and SOAR  
 253 solutions. The Security Information and Event Management system provides monitoring and data analytics  
 254 based on the logs collected from network assets (devices, applications), while the Security Orchestration,  
 255 Automation, and Response (SOAR) system is used to automate response after detection and alerting of  
 256 potential attacks.

257 For the purposes of the banking cybersecurity scenario, the proposed tools, i.e. AI4FIDS and  
 258 TRUST4AI.XAI are sandboxed using the bank's Innovation Sandbox. The main source of the infor-  
 259 mation is an interconnected SIEM solution feeding the isolated tools, while the abovementioned corporate  
 260 SOAR solution is used to trigger tailored playbooks. The tools can also access internal development tools  
 261 and repositories to analyse the code and development processes.

262 From the end-user viewpoint, CaixaBank adopts AI4CYBER tools for ensuring robust security in critical  
 263 environments like the SWIFT client and the Financial Terminal. AI4FIDS enhanced by xAI capabilities,  
 264 delivers AI services that can detect abnormal actions, identify impersonations of privileged users, and  
 265 prevent intrusions and AI-driven attacks in real-time. In addition, the tools facilitate comprehensive  
 266 monitoring of user behaviours and activities across different bank services. The xAI enhancement, provides  
 267 insight into the model's output, making cybersecurity-related decisions transparent and understandable for  
 268 the bank's security staff, saving time and cost of security operations.

## 269 EXPERIMENTS AND RESULTS

### 270 AI4FIDS – Initial Evaluation Results with Existing Cybersecurity Datasets

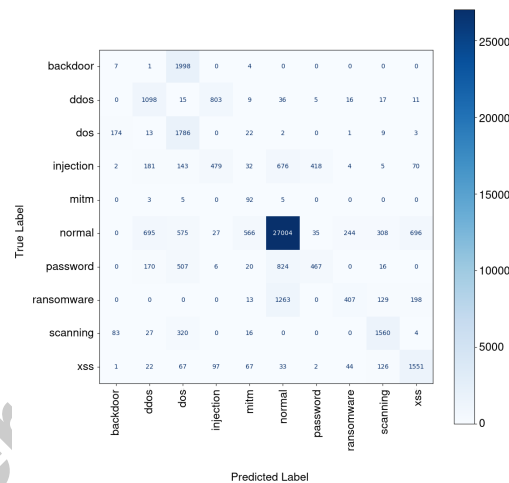
271 In order to demonstrate the efficacy and soundness of the initial version of AI4FIDS in conjunction  
 272 with benchmark cybersecurity datasets, this section describes the system prerequisites and technical  
 273 specifications of AI4FIDS. Multiple datasets were employed for evaluating the proposed implementation,  
 274 yet in the context of this study, the distinct results for the aggregation techniques are provided based on the  
 275 CSE CIC-IDS-2018 Dataset [31].

276 To that end, the initial results of the evaluation procedure of the respective containers, namely L-FIDS,  
 277 O-FIDS, N-FIDS and V-FIDS, are presented. More specifically, T4FIDS is responsible for generating the  
 278 federated models that will be utilized by the detection engines of the preceding containers. The testbed

279 onto which the proposed implementation is trained on, is composed of three Federated Clients and one  
 280 Federated Server. In the current study, preliminary results are provided for the N-FIDS container along with  
 281 the CSE CIC-IDS-2018 Dataset [30]. In Table 1, the respective results from the evaluation procedure are  
 282 shown when employing the network flow statistics calculated by the CICFlowMeter. As one may observe, a  
 283 thorough inspection of a variety of aggregation methods is performed in the context of evaluating the initial  
 284 version of AI4FIDS, where the N-FIDS container must detect attacks when the aforementioned features are  
 285 considered. More notably, the detection engine, namely T4FIDS, is trained with the TCP/IP network flow  
 286 statistics, where multiple cyberattacks are considered and five aggregation techniques are scrutinized. The  
 287 attained results indicate that the finest performance was achieved from the FedProx technique. Additionally,  
 288 in Fig. 3 the confusion matrix of the model that uses FedAvg with the TON IoT Dataset is illustrated.

**Table 1.** Performance comparison of different aggregation methods.

Aggregation	ACC	TPR	FPR	F1	AUC
FedAvg	84.97%	80.96%	1.13%	85.80%	98.60%
FedProx	86.73%	78.68%	1.01%	87.42%	98.17%
FedAdam	27.80%	35.66%	5.52%	28.11%	77.83%
FedAdagrad	85.66%	74.28%	1.10%	86.19%	97.86%
FedYogi	74.93%	71.01%	1.86%	77.22%	95.41%



**Figure 3.** Confusion Matrix of the FL model that uses FedAvg with the TON IoT Dataset – Network Flow Statistics.

289 **xAI Interpretation – Experiments and Results with Existing Cybersecurity Datasets**

290 With the detection model established and evaluated, the experiment proceeded with attempting to gain  
 291 insight into the decision-making process of the classifier. To this end, the xAI methods were employed.  
 292 Following a scenario where a security operative needs to justify a decisive action to ban or not to ban a user  
 293 based on the detection result, the xAI methods aim to provide reasoning as to why the samples in question  
 294 were classified as an attack.

295 Fig. 4 demonstrates the application of the LIME technique within the TRUST4AI.XAI component  
 296 of AI4CYBER. It shows how LIME decomposes a cybersecurity model’s decision-making process for  
 297 individual predictions, highlighting the contribution of each feature towards the predicted outcome of  
 298 identifying network threats. Fig. 5 illustrates a decision tree that approximates the complex decision  
 299 boundaries of the explained model, providing a simplified view of how various network statistics influence  
 300 the classification of traffic. Fig. 6 showcases the SHAP explanations, each bar in the visualization represents  
 301 the impact of an individual feature on the model’s prediction, giving insight into which attributes are most  
 302 influential for detecting cybersecurity threats. Fig. 7 shows ICE and PDP plots that analyze the impact of  
 303 specific features on the predictions of the AI4FIDS cybersecurity model. These plots help in understanding  
 304 the relationships between the feature values and the likelihood of an event being flagged as a security threat,  
 305 across a range of values for the selected features.



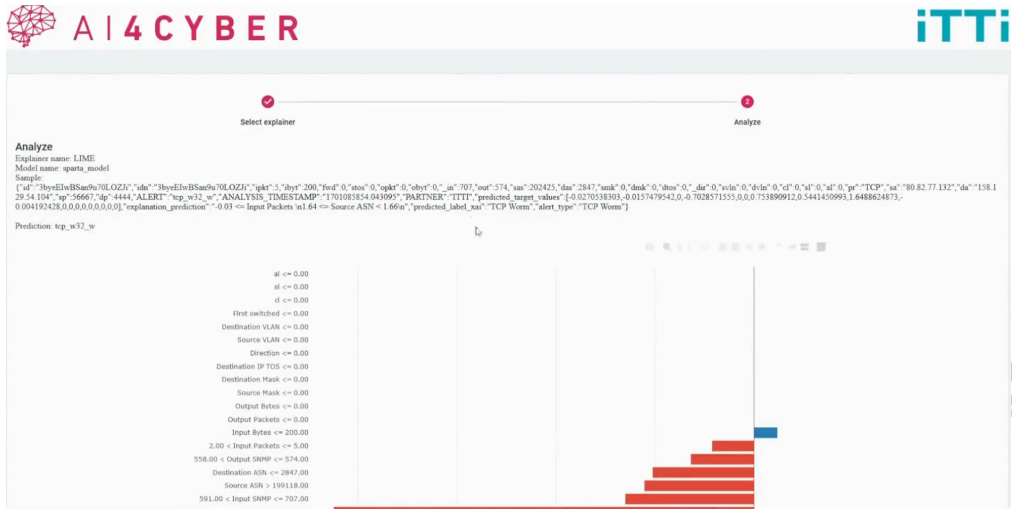


Figure 4. Examples of xAI methods: LIME explanations in the TRUST4AI.XAI component of AI4CYBER.

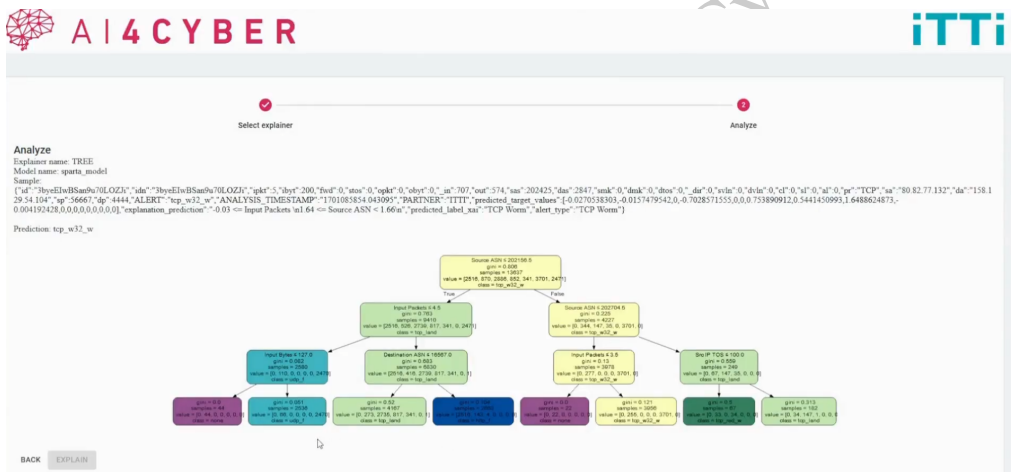


Figure 5. Examples of xAI methods: Surrogate Tree Aggregations in the TRUST4AI.XAI component of AI4CYBER.

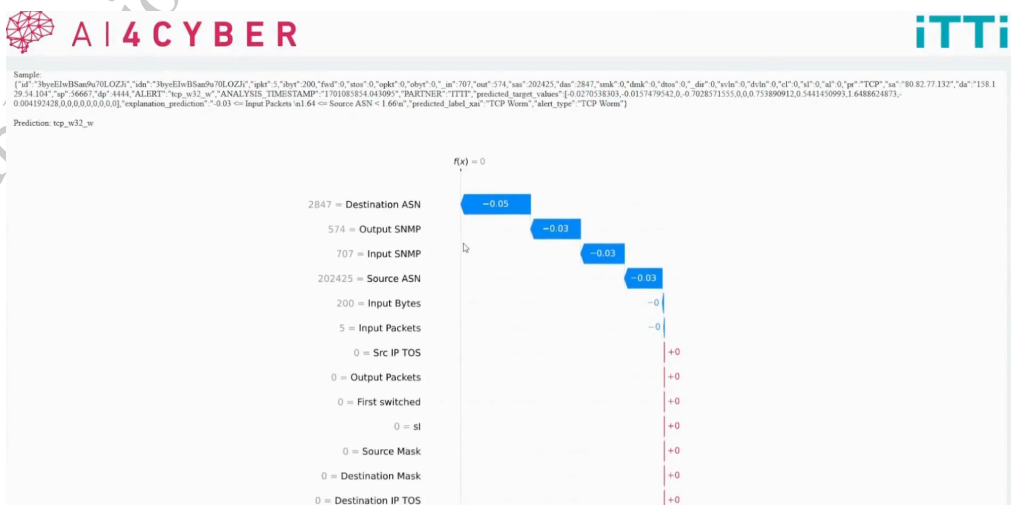
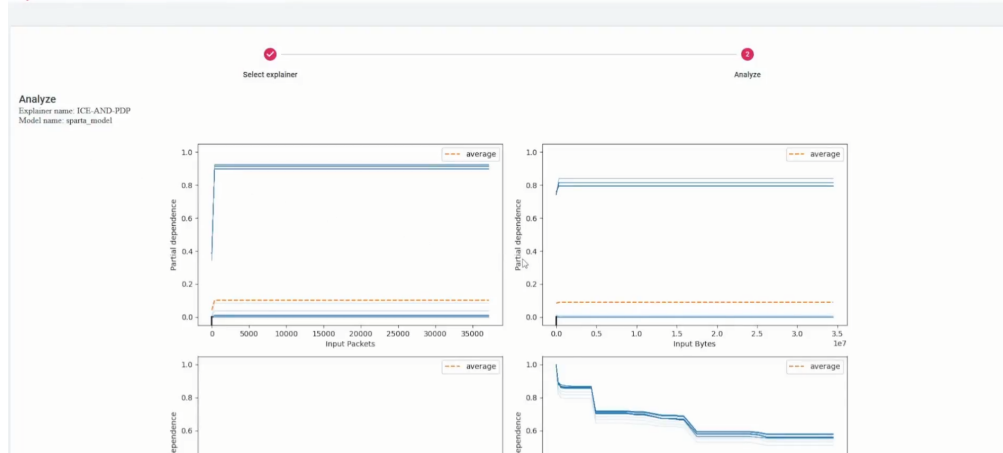


Figure 6. Examples of xAI methods: SHAP explanations in the TRUST4AI.XAI component of AI4CYBER.





**Figure 7.** Examples of xAI methods: ICE and PDP plots, explanations in the TRUST4AI.XAI component of AI4CYBER.

## CONCLUSION

In this paper, AI4FIDS and TRUST4AI.XAI were introduced, two pivotal components of the AI4CYBER project aimed at enhancing cybersecurity in the financial sector through the utilization of advanced AI techniques and end-user-centric deployment in the CaixaBank pilot. AI4FIDS, empowered by Federated Learning, offers a privacy-preserving approach that enables robust intrusion detection across decentralized networks without compromising sensitive data. The TRUST4AI.XAI module provides critical explainability, ensuring that AI-driven decisions are transparent and interpretable to end-users. The sector-oriented experiments conducted with these systems within the financial pilot of CaixaBank demonstrate their efficacy in both detecting a range of cyberthreats and in aligning with requirements for transparency.

The implementation demonstrated the practical viability and effectiveness of combining Federated Learning with advanced explainability to secure financial infrastructures. The success of this integration not only confirms the potential of these technologies to improve cybersecurity practices but also sets a precedent for the development of privacy-preserving, interpretable AI models in the financial sector.

## ACKNOWLEDGMENTS

This research is funded under the Horizon Europe AI4CYBER Project, which has received funding from the European Union's Horizon Europe research and innovation programme under grant agreement No. 101070450.

## CONFLICTS OF INTEREST

The authors declare no conflict of interest.

## REFERENCES

- [AI4CYBER] AI4CYBER. Ai4cyber. (Accessed on 07/22/2024).
- [2] AL-Dosari, K., Fetais, N., and Kucukvar, M. (2024). Artificial intelligence and cyber defense system for banking industry: A qualitative study of ai applications and challenges. *Cybernetics and systems*, 55(2):302–330.
  - [3] Alazab, M., RM, S. P., Parimala, M., Maddikunta, P. K. R., Gadekallu, T. R., and Pham, Q.-V. (2021). Federated learning for cybersecurity: Concepts, challenges, and future directions. *IEEE Transactions on Industrial Informatics*, 18(5):3501–3509.
  - [4] Altmann, A., Toloşi, L., Sander, O., and Lengauer, T. (2010). Permutation importance: a corrected feature importance measure. *Bioinformatics*, 26(10):1340–1347.
  - [5] Aouedi, O., Piamrat, K., Muller, G., and Singh, K. (2022). Fluids: Federated learning with semi-supervised approach for intrusion detection system. In *2022 IEEE 19th Annual Consumer Communications & Networking Conference (CCNC)*, pages 523–524. IEEE.

- 338 [6] Apley, D. W. and Zhu, J. (2020). Visualizing the effects of predictor variables in black box supervised  
339 learning models. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 82(4):1059–  
340 1086.
- 341 [7] Arisdakessian, S., Wahab, O. A., Mourad, A., Otrók, H., and Guizani, M. (2022). A survey on iot  
342 intrusion detection: Federated learning, game theory, social psychology, and explainable ai as future  
343 directions. *IEEE Internet of Things Journal*, 10(5):4059–4092.
- 344 [8] Campos, E. M., Saura, P. F., González-Vidal, A., Hernández-Ramos, J. L., Bernabe, J. B., Baldini, G.,  
345 and Skarmeta, A. (2022). Evaluating federated learning for intrusion detection in internet of things:  
346 Review and challenges. *Computer Networks*, 203:108661.
- 347 [9] Choraś, M., Pawlicki, M., Puchalski, D., and Kozik, R. (2020). Machine learning—the results are  
348 not the only thing that matters! what about security, explainability and fairness? In *Computational  
349 Science—ICCS 2020: 20th International Conference, Amsterdam, The Netherlands, June 3–5, 2020,  
350 Proceedings, Part IV 20*, pages 615–628. Springer.
- 351 [Daly] Daly, R. A standout performance: Q&a with caixabank ceo gonzalo gortázar — global finance  
352 magazine. (Accessed on 09/06/2024).
- 353 [11] Dong, T., Qiu, H., Lu, J., Qiu, M., and Fan, C. (2021). Towards fast network intrusion detection based  
354 on efficiency-preserving federated learning. In *2021 IEEE Intl Conf on Parallel & Distributed Processing  
355 with Applications, Big Data & Cloud Computing, Sustainable Computing & Communications, Social  
356 Computing & Networking (ISPA/BDCLOUD/SocialCom/SustainCom)*, pages 468–475. IEEE.
- 357 [12] Familoni, B. T. and Shoetan, P. O. (2024). Cybersecurity in the financial sector: a comparative analysis  
358 of the usa and nigeria. *Computer Science & IT Research Journal*, 5(4):850–877.
- 359 [13] Fan, Y., Li, Y., Zhan, M., Cui, H., and Zhang, Y. (2020). Iotdefender: A federated transfer learning  
360 intrusion detection framework for 5g iot. In *2020 IEEE 14th international conference on big data  
361 science and engineering (BigDataSE)*, pages 88–95. IEEE.
- 362 [14] Friha, O., Ferrag, M. A., Shu, L., Maglaras, L., Choo, K.-K. R., and Nafaa, M. (2022). Felids:  
363 Federated learning-based intrusion detection system for agricultural internet of things. *Journal of  
364 Parallel and Distributed Computing*, 165:17–31.
- 365 [15] Ghimire, B. and Rawat, D. B. (2022). Recent advances on federated learning for cybersecurity and  
366 cybersecurity for federated learning for internet of things. *IEEE Internet of Things Journal*, 9(11):8229–  
367 8249.
- 368 [16] Goldstein, A., Kapelner, A., Bleich, J., and Pitkin, E. (2015). Peeking inside the black box: Visualizing  
369 statistical learning with plots of individual conditional expectation. *Journal of Computational and  
370 Graphical Statistics*, 24(1):44–65.
- 371 [17] Greenwell, B. M., Boehmke, B. C., and McCarthy, A. J. (2018). A simple and effective model-based  
372 variable importance measure. *arXiv preprint arXiv:1805.04755*.
- 373 [18] Gurumoorthy, K. S., Dhurandhar, A., Cecchi, G., and Aggarwal, C. (2019). Efficient data representation  
374 by selecting prototypes with importance weights. In *2019 IEEE International Conference on Data  
375 Mining (ICDM)*, pages 260–269. IEEE.
- 376 [19] Lavaur, L., Pahl, M.-O., Busnel, Y., and Autrel, F. (2022). The evolution of federated learning-based  
377 intrusion detection and mitigation: a survey. *IEEE Transactions on Network and Service Management*,  
378 19(3):2309–2332.
- 379 [20] Li, B., Wu, Y., Song, J., Lu, R., Li, T., and Zhao, L. (2020). Deepfed: Federated deep learning for  
380 intrusion detection in industrial cyber-physical systems. *IEEE Transactions on Industrial Informatics*,  
381 17(8):5615–5624.
- 382 [21] Lundberg, S. M. and Lee, S.-I. (2017). A unified approach to interpreting model predictions. *Advances  
383 in neural information processing systems*, 30.
- 384 [22] Mirzaee, P. H., Shojafar, M., Pooranian, Z., Asefy, P., Cruickshank, H., and Tafazolli, R. (2021). Fids:  
385 A federated intrusion detection system for 5g smart metering network. In *2021 17th International  
386 Conference on Mobility, Sensing and Networking (MSN)*, pages 215–222. IEEE.
- 387 [23] Mothilal, R. K., Sharma, A., and Tan, C. (2020). Explaining machine learning classifiers through  
388 diverse counterfactual explanations. In *Proceedings of the 2020 conference on fairness, accountability,  
389 and transparency*, pages 607–617.
- 390 [24] Nauta, M., Trienes, J., Pathak, S., Nguyen, E., Peters, M., Schmitt, Y., Schlötterer, J., Van Keulen,  
391 M., and Seifert, C. (2023). From anecdotal evidence to quantitative evaluation methods: A systematic  
392 review on evaluating explainable ai. *ACM Computing Surveys*, 55(13s):1–42.
- 393 [25] Pawlicka, A., Choraś, M., and Pawlicki, M. (2020). Cyberspace threats: not only hackers and criminals.  
394 raising the awareness of selected unusual cyberspace actors—cybersecurity researchers’ perspective. In  
395 *Proceedings of the 15th International Conference on Availability, Reliability and Security*, pages 1–11.
- 396 [26] Popoola, S. I., Gui, G., Adebisi, B., Hammoudeh, M., and Gacanin, H. (2021). Federated deep

- 397 learning for collaborative intrusion detection in heterogeneous networks. In *2021 IEEE 94th Vehicular*  
398 *Technology Conference (VTC2021-Fall)*, pages 1–6. IEEE.
- 399 [27] Ribeiro, M. T., Singh, S., and Guestrin, C. (2016). ” why should i trust you?” explaining the predictions  
400 of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge*  
401 *discovery and data mining*, pages 1135–1144.
- 402 [28] Ribeiro, M. T., Singh, S., and Guestrin, C. (2018). Anchors: High-precision model-agnostic explana-  
403 tions. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32.
- 404 [29] Schneble, W. and Thamarasu, G. (2019). Attack detection using federated learning in medical  
405 cyber-physical systems. In *Proc. 28th Int. Conf. Comput. Commun. Netw.(ICCCN)*, volume 29, pages  
406 1–8.
- 407 [30] Sharafaldin, I., Gharib, A., Lashkari, A. H., Ghorbani, A. A., et al. (2018a). Towards a reliable  
408 intrusion detection benchmark dataset. *Software Networking*, 2018(1):177–200.
- 409 [31] Sharafaldin, I., Lashkari, A. H., Ghorbani, A. A., et al. (2018b). Toward generating a new intrusion  
410 detection dataset and intrusion traffic characterization. *ICISSp*, 1:108–116.
- 411 [32] Uzougbo, N. S., Ikegwu, C. G., Adewusi, A. O., et al. (2024). Cybersecurity compliance in financial  
412 institutions: a comparative analysis of global standards and regulations. *International Journal of Science*  
413 *and Research Archive*, 12(1):533–548.
- 414 [33] Zhao, R., Wang, Y., Xue, Z., Ohtsuki, T., Adebisi, B., and Gui, G. (2022). Semisupervised federated-  
415 learning-based intrusion detection method for internet of things. *IEEE Internet of Things Journal*,  
416 10(10):8645–8657.

pre-publication version only for AI4Cyber Zenodo