



HAL
open science

Federated Learning for Zero-Day Attack Detection in 5G and Beyond V2X Networks

Abdelaziz Amara Korba, Abdelwahab Boualouache, Bouziane Brik, Rabah Rahal, Yacine Ghamri-Doudane, Sidi Mohammed Senouci

► **To cite this version:**

Abdelaziz Amara Korba, Abdelwahab Boualouache, Bouziane Brik, Rabah Rahal, Yacine Ghamri-Doudane, et al.. Federated Learning for Zero-Day Attack Detection in 5G and Beyond V2X Networks. AlgoTel 2023 - 25èmes Rencontres Francophones sur les Aspects Algorithmiques des Télécommunications, May 2023, Cargese (Corse), France. hal-04087452

HAL Id: hal-04087452

<https://hal.science/hal-04087452>

Submitted on 3 May 2023

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Apprentissage fédéré pour la détection des attaques Zero-Day dans les réseaux V2X 5G et 5GB

Abdelaziz Amara korba ^{1,4 †}, Abdelwahab Boualouache ², Bouziane Brik ³, Rabah Rahal ⁴, Yacine Ghamri-Doudane ¹, Sidi Mohammed Senouci ³

¹L3I, University of La Rochelle, France, ²FSTM, University of Luxembourg, Luxembourg, ³LRS, Badji Mokhtar An-naba University, Algeria ⁴DRIVE, University of Bourgogne, France

Le déploiement de véhicules connectés et automatisés (CAV) utilisant les réseaux 5G et 5G-B les rend vulnérables à de nouveaux vecteurs d'attaques. Afin de détecter et de mitiger ces attaques, plusieurs systèmes de détection d'intrusions (IDS) basés sur l'apprentissage automatique ont été proposés, dont la grande majorité utilise l'apprentissage profond supervisé. Cependant, la principale limite de ce type de solution est son incapacité à détecter des attaques différentes de celles vues lors de l'apprentissage de l'IDS, ou de nouvelles attaques, également appelées attaques Zero-day. De plus, l'apprentissage du modèle de détection nécessite une collecte et un étiquetage importants des données, ce qui augmente le coût de communication et soulève des problèmes de confidentialité. Afin de pallier ces problèmes, nous proposons dans cet article un nouveau IDS basé sur l'apprentissage fédéré non supervisé préservant la confidentialité, et dont l'apprentissage du système ne nécessite que le trafic réseau légitime. L'expérimentation de la solution proposée sur un jeu de données récent montre que le système proposé présente un taux de détection élevé tout en minimisant le taux de faux positifs et le délai de détection.

Mots-clefs : 5G, V2X, IoV, Sécurité, Attaques, Détection, Apprentissage Fédéré

1 Introduction

Les réseaux de cinquième génération (5G) et au-delà de la 5G (5GB) promettent de révolutionner les systèmes de transport intelligents (STI) avec des communications à latence ultra-faible et à bande passante élevée. Cependant, les véhicules connectés et automatisés (CAV) à différents niveaux d'automatisation seront confrontés à de nouveaux vecteurs d'attaques provenant des technologies 5G/5GB, conduisant ainsi à des situations dangereuses pour les usagers de la route. L'intelligence artificielle, et l'apprentissage automatique en particulier apparaît comme un allié clé de la cybersécurité [BE22a]. Divers systèmes de détection d'intrusion (IDS) basés sur l'apprentissage profond ont été proposés pour protéger les STIs des attaques. La plupart de ces IDSs s'appuient sur un apprentissage supervisé et centralisé [RAKGZ⁺22]. L'apprentissage centralisé nécessite une collecte et un étiquetage importants des données, ce qui augmente le coût de communication, et soulèvent des problèmes de confidentialité. Pour pallier ces limitations, des recherches récentes [BE22b, LZZ⁺21] ont exploité le potentiel de l'apprentissage fédéré (FL), qui a montré des résultats prometteurs dans divers domaines. Par contre, tous les IDS à base de FL proposés pour les ITS reposent sur des techniques d'apprentissage supervisé. Une limitation importante de l'utilisation de ces techniques est leur incapacité à détecter des attaques différentes de celles vues lors de la phase d'apprentissage et des attaques de type « zero-day ». D'autant plus, la plupart des approches existantes supposent que les IDS déployés sur les CAVs (clients FL) détiennent des données préalablement étiquetées, ce qui n'est pas toujours vrai.

Puisque un CAV exécute un ensemble d'applications bien connues (sécurité, commodité, commerciales, etc.), sont comportement réseau doit présenter un degré élevé de régularité tant qu'il n'est pas attaqué

[†]5G-INSIGHT (ANR-20-CE25-0015-16)

ou défectueux. De même, une attaque doit altérer son modèle de communication. Dans cet article, nous proposons un IDS à base d'apprentissage fédéré non-supervisé pour modéliser le comportement réseau légitime d'un CAV (ou attendu), et ainsi détecter les attaques comme des anomalies. L'apprentissage fédéré de l'IDS est orchestré par le serveur Multi-access Edge Computing (MEC) pour améliorer l'efficacité de l'apprentissage et minimiser la latence. L'expérimentation de la solution proposée sur un jeu de données récent montre que le système proposé présente un taux de détection élevé tout en minimisant le taux de faux positifs et le délai de détection.

2 Solution proposée

Nous avons développé un système de détection d'intrusion (IDS) basé sur un auto-encodeur (AE) entraîné de manière fédérée, exclusivement sur des flux réseau légitimes. Pour cela, nous avons utilisé un serveur MEC qui a permis d'orchestrer l'apprentissage en distribuant les modèles et en agrégeant les gradients. Le modèle initial de l'AE a été distribué par le serveur MEC à chaque CAV, qui l'a ensuite entraîné sur ses propres données locales. Ces données ont été constituées d'un ensemble de flux extraits à partir des paquets capturés localement par le CAV. Pour chaque flux, un ensemble d'attributs a été calculé. Les cycles d'apprentissage (round) ont été orchestrés par le serveur MEC et exécutés par l'AE sur l'ensemble des flux locaux.

2.1 Extraction des flux et calcul des attributs

Tout d'abord, pour identifier un flux réseau, nous utilisons une combinaison de cinq informations de l'en-tête du paquet, à savoir : l'adresse IP source, l'adresse IP de destination, le numéro de port source, le numéro de port de destination et le type de protocole de couche transport. Pour chaque flux extrait, un ensemble d'attributs est calculé selon une fenêtre temporelle donnée (ex. 100 secondes). Les attributs sont constitués principalement de statistiques (moyenne, médiane, écart-type, etc.) calculées à partir des informations des en-têtes des couches réseau et transport de la pile TCP/IP, tels que le nombre d'octets envoyés/reçus par seconde, l'intervalle moyen entre deux paquets successifs, la taille maximale des paquets envoyés/reçus, etc.

2.2 Modélisation du trafic légitime et détection d'anomalies

Le modèle Auto-Encoder (AE) [HZ93] est un modèle d'apprentissage profond non supervisé qui comprime les vecteurs d'entrée sous forme de vecteurs de code à l'aide d'un ensemble de poids de reconnaissance, puis les reconstruit en m ($m < d$) nombre de neurones reconstruits en utilisant un ensemble de poids génératifs. Il y a deux parties principales dans l'architecture de l'AE : l'encodeur et le décodeur. L'encodeur réduit la dimension des vecteurs d'entrée ($x_i \in R^d$) au nombre de neurones qui forment la couche cachée. L'activation du neurone i dans la couche cachée est donnée par :

$$h_i = f^\theta(x) = s\left(\sum_{j=1}^n W_{ij}^{input} x_j + b_i^{input}\right) \quad (1)$$

Où x est le vecteur d'entrée, θ les paramètres $\{W^{input}, b^{input}\}$, W est une matrice de poids du codeur de dimension $m \times d$, ou b est un vecteur de biais de dimension m . Ainsi, le vecteur d'entrée est codé en un vecteur de faible dimension. Le décodeur fait le mappage de la représentation cachée de faible dimension h_i à l'espace d'entrée initial R^d par la même transformation que le codeur.

L'erreur de reconstruction de l'AE est utilisée comme score d'anomalie. Les flux réseau avec des erreurs de reconstruction importantes sont considérés comme des flux malveillants (anomalies). Après l'apprentissage, le modèle reconstruira exceptionnellement bien les flux légitimes, mais pas les flux malveillants qu'il n'a jamais vus. **Algorithm 1** décrit le processus de détection d'anomalies en utilisant les erreurs de reconstruction de l'auto-encodeur. L'écart absolu médian (MAD) a été utilisé pour le calcul du seuil α , parce qu'il est moins susceptible d'être faussé par les valeurs aberrantes.

Algorithm 1: Détection du trafic malveillant

```

1 BEGIN
2 INPUT : dataset légitime  $X$ , dataset malveillant  $x^{(i)} \ i \in \{1, \dots, N\}$ , Seuil :
    $\alpha = \widehat{MSE}_{benign\_val} + 5 \times MAD(MSE_{benign\_val})$ 
3 OUTPUT : RE (Erreur de reconstruction)  $\|x - \hat{x}\|$ 
4  $\phi, \theta \leftarrow$  entraîner l'AE avec le dataset légitime  $X$ 
5 for  $i \in \{1, \dots, N\}$  do
6    $RE(i) = \|x^{(i)} - g_{\theta}(f_{\phi}(x^{(i)}))\|$ 
7   if  $RE(i) > \alpha$  then
8      $x^{(i)}$  est un flux malveillant
9   else
10     $x^{(i)}$  est un flux légitime
11  end if
12 end for
13 END

```

2.3 Processus d'apprentissage fédéré assisté par MEC

Le problème d'apprentissage fédéré est formulé comme un problème d'optimisation fédérée et résolu à l'aide de l'algorithme FedAvg [MMR⁺17]. En effet, en utilisant ses données locales, chaque CAV calcule le gradient moyen pour un cycle d'apprentissage. D'autre part, le serveur MEC agrège les mises à jour locales reçus et retransmet le modèle global aux CAVs participants. Ce processus est répété pendant un certain nombre de rounds, défini initialement par le serveur MEC.

3 Évaluation des performances

Nous avons testé notre solution sur le jeu de données [RAKGZ20], soixante-quatre attributs ont été calculés pour les flux légitimes et malveillants extraits à partir du trafic brut (fichiers PCAP). Tenant compte du fait qu'une fenêtre temporelle d'échantillonnage plus petite permet une détection plus rapide, nous avons fixé la fenêtre temporelle à 1 seconde. Nous avons implémenté un auto-encodeur profond avec trois couches cachées. La fonction de loss d'erreur quadratique moyenne (MSE) binaire a été utilisée. Pour évaluer les performances de détection de notre solution, en plus de l'accuracy, nous avons aussi considéré le F1-Score et le taux de faux positifs (FPR).

Globalement, notre solution présente un taux de détection élevé avec un faible taux de faux positifs. Le tableau 1 montre comment la l'accuracy moyenne, le F1-score et le taux de faux positifs varient en changeant le nombre de CAVs participants à l'apprentissage fédéré. L'apprentissage avec trois CAV a permis d'obtenir à peu près les mêmes performances tout en réduisant le temps d'apprentissage d'environ 30%. D'après la figure 1a, notre solution détecte légèrement mieux l'attaque SynFlood. Cela peut s'expliquer par le nombre relativement important d'attributs liées au protocole TCP (comparé au protocole UDP). Bien que 99,99 % de F1-Score et 0,01 % de FPR aient été obtenus en utilisant l'apprentissage supervisé dans [RAKGZ20], notre solution présente des performances comparables en utilisant uniquement le trafic légitime. Nous avons aussi comparé notre modèle fédéré avec le même modèle entraîné de façon centralisée. La comparaison des résultats illustrée dans la figure 1b montre que le modèle fédéré a donné des performances de détection remarquablement similaires à celles du modèle centralisé.

4 Conclusion

Dans cet article, nous avons proposé un nouveau IDS basé sur un modèle d'auto-encodeur profond, qui exploite la prédictibilité du trafic réseau légitime des CAVs pour détecter les attaques. En utilisant un apprentissage fédéré orchestré par le serveur MEC, l'IDS que nous avons proposé ne nécessite ni étiquetage ni partage de données entre les CAV. Nous avons démontré que notre IDS présente des performances de

TABLE 1 – Évaluation des performances de détection

Nb. Clients	Accuracy	F1-Score	FPR	Temps (mn)
10	87.94%	91.21%	6.95%	14.80
8	87.94%	91.22%	6.98%	11.53
6	87.95%	91.23%	7.06%	9.32
3	87.94%	91.21%	6.92%	4.43

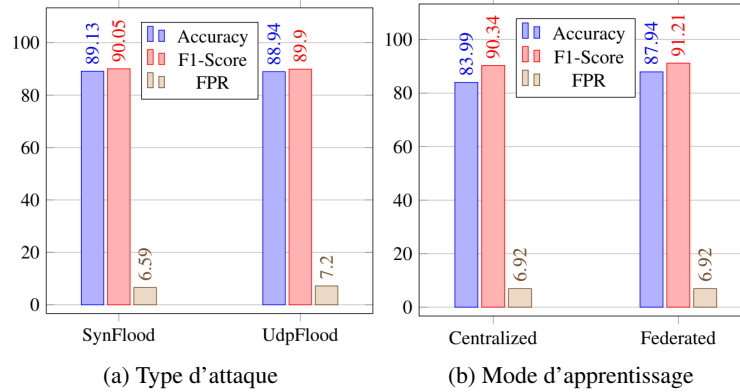


FIGURE 1 – Comparaison des performances prédictives

détection élevées, même avec un faible nombre de cycles de communication et un échantillonnage court (1 seconde). Cela permet un apprentissage rapide et un délai de détection court. À l'avenir, nous envisageons d'évaluer notre IDS sur d'autres ensembles de données de distribution indépendante non identiques (non-IID), y compris des attaques plus sophistiquées et récentes.

Références

- [BE22a] Abdelwahab Boualouache and Thomas Engel. A Survey on Machine Learning-based Misbehavior Detection Systems for 5G and Beyond Vehicular Networks. *arXiv preprint arXiv :2201.10500*, 2022.
- [BE22b] Abdelwahab Boualouache and Thomas Engel. Federated learning-based scheme for detecting passive mobile attackers in 5G vehicular edge computing. *Annals of Telecommunications*, 77(3) :201–220, 2022.
- [HZ93] Geoffrey E Hinton and Richard Zemel. Autoencoders, minimum description length and helmholtz free energy. *Advances in neural information processing systems*, 6, 1993.
- [LZZ⁺21] Hong Liu, Shuaipeng Zhang, Pengfei Zhang, Xinqiang Zhou, Xuebin Shao, Geguang Pu, and Yan Zhang. Blockchain and Federated Learning for Collaborative Intrusion Detection in Vehicular Edge Computing. *IEEE Transactions on Vehicular Technology*, 2021.
- [MMR⁺17] H. B. McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *AISTATS*, 2017.
- [RAKGZ20] Rabah Rahal, Abdelaziz Amara Korba, and Nacira Ghoualmi-Zine. Towards the development of realistic dos dataset for intelligent transportation systems. *Wireless Personal Communications*, 115(2) :1415–1444, 2020.
- [RAKGZ⁺22] Rabah Rahal, Abdelaziz Amara Korba, Nacira Ghoualmi-Zine, Yacine Challal, and Mohamed Yacine Ghamri-Doudane. Antibotv : A multilevel behaviour-based framework for botnets detection in vehicular networks. *Journal of Network and Systems Management*, 30(1) :1–40, 2022.