

România  
Ministerul Apărării Naționale  
Academia Tehnică Militară "Ferdinand I"

Facultatea de Sisteme Informatică și Securitate Cibernetică  
CALCULATOARE ȘI SISTEME INFORMATICE PENTRU APĂRARE ȘI  
SECURITATE NAȚIONALĂ



Platformă de Analiză Automată a Atacurilor Data Poisoning asupra unei  
Infrastructuri de Învățare Federate

**Coordonator Științific**

Cpt. conf. dr. ing. Iulian Aciobăniței

**Absolvent**

Sd. Sg. Maj. Lepădatu Tudor

Conține \_\_\_\_\_ file  
Inventariat sub numărul \_\_\_\_\_  
Cu poziția din indicator \_\_\_\_\_  
Cu termen de păstrare \_\_\_\_\_

București  
An 2026

Mulțumiri pentru persoanele care au sprijinit procesul de realizare a acestei lucrări

# Referatul Coordonatorului Științific

Referatul reprezintă un text în care coordonatorul științific sumarizează, ulterior finalizării conținutului efectiv, efortul pe care l-ați depus și trage concluzii cu privire la gradul de realizare a obiectivelor propuse inițial (cele prezentate în cadrul detalierii). De regulă, acest referat nu depășește cele 4 pagini alocate în acest document.







# Tema Proiectului de Diplomă

Tema proiectului de diplomă (sau detalierea) reprezintă un text realizat de coordonatorul științific, în lunile următoare propunerii titlului proiectului de diplomă către facultate, în care detaliază nevoia reală a implementării unui astfel de proiect și modul în care se dorește a fi implementat. În plus, poate prezenta structura pe capitole a viitoarei lucrări, anexele ce vor fi incluse și sursele bibliografice din care studentul se va informa. De regulă, această detaliere nu depășește cele 4 pagini alocate în acest document.









# Abstract

The Artificial Intelligence integration with tools and applications has evolved since 2020s and the cybersecurity scene tries to adapt frequently. Information security and data integrity is more important than never before, being used by Machine Learning models or Neural Networks trained to perform specific tasks.

From a data science point of view, the quality of information is much important than securing it. The AI evolution has led to the creation of different attack boundaries, from changing the model parameters to perform data poisoning schemes. Confidentiality is the key in maintaining unique aspects for each entity involved in the federated learning process.

In this paper, data poisoning attacks are studied with different options in a segregated simulated infrastructure called federated learning, in which each client may change its scope intentionally or unintentionally. Each simulation has its own configuration providing data scientist with a dedicated environment for testing its machine learning algorithm against data poisoning attacks.

# Cuprins

<b>1</b>	<b>Introducere</b>	<b>1</b>
1.1	Context . . . . .	1
1.2	Motivatia lucrarii . . . . .	2
1.3	Obiectivele lucrarii . . . . .	2
1.4	Structura lucrarii . . . . .	2
<b>2</b>	<b>Notiuni Teoretice</b>	<b>3</b>
2.1	Notiuni introductive . . . . .	3
2.1.1	Diferenta dintre Machine Learning si Deep Learning . . . . .	3
2.1.2	Retea Neuronala . . . . .	3
2.2	Invatare automata federata . . . . .	3
2.2.1	Concept . . . . .	4
2.2.2	Arhitectura FL . . . . .	4
2.2.3	Procesul de antrenare FL . . . . .	6
2.2.4	Exemple in viata reala . . . . .	7
2.3	Atacuri de tip Data Poisoning . . . . .	8
2.3.1	Definirea tipurilor de atac . . . . .	8
2.3.2	Vectori de atac . . . . .	8
2.3.3	Atacul Data poisoning . . . . .	9
2.3.4	Impactul poisoning in Federated Learning . . . . .	9
2.4	Alte Notiuni . . . . .	10
2.4.1	Docker . . . . .	10
2.4.2	Python . . . . .	10
2.4.3	Rest API . . . . .	10
<b>3</b>	<b>Proiectare, Implementare si Testare</b>	<b>11</b>
3.1	Cerintele Software . . . . .	11
3.1.1	Cerintele functionale . . . . .	11
3.1.2	Cerintele nefunctionale . . . . .	11
3.2	Arhitectura platformei . . . . .	11
3.2.1	Containere . . . . .	11
3.2.2	Server . . . . .	11
3.3	Testare . . . . .	11
<b>4</b>	<b>Rezultate si Metrice Simulari</b>	<b>12</b>
4.1	Evaluare Performante . . . . .	12
4.1.1	Scalabilitatea Simularilor . . . . .	12
4.1.2	Scalabilitatea platformei . . . . .	12
4.2	Evaluare Rezultate . . . . .	12
4.2.1	Performante Gaussian Noise . . . . .	12
4.2.2	Performante Label-Flip . . . . .	12
4.2.3	Performante Backdoor . . . . .	12
<b>5</b>	<b>Concluzii si dezvoltare ulterioara</b>	<b>13</b>
5.1	Starea Curenta . . . . .	13
5.2	Dezvoltare Ulterioara . . . . .	13
5.3	Tabele . . . . .	13
5.4	Imagini . . . . .	15
5.5	Liste . . . . .	15
5.6	Formule Matematice . . . . .	15
5.7	Note de Subsol. Citări . . . . .	16
5.8	Etichete. Referințe . . . . .	16
	<b>Bibliografie</b>	<b>17</b>

# Listă de figuri

2.1	Imaginea 2.2.2: Modele Federated Learning . . . . .	4
2.2	Imagine 2.2.3.1: Arhitectura interna a unui dispozitiv . . . . .	6
2.3	Imagine 2.2.3.2: Procedee in invatare automata federata . . . . .	7
5.1	Arhitectura unui calculator . . . . .	15

# Listă de Abrevieri

UE ..... Uniunea Europeană

EU ..... *European Union*







# Capitolul 1:

## Introducere

### 1.1 Context

Odata cu dezvoltarea sistemelor de calcul moderne si a componentelor Hardware, s-au putut realiza produse software complexe cu capacitati de stocare net superioare. Revolutia tehnologica a permis nu doar realizarea unor sarcini simple, precum calcule matematice, sau automatizarea unor dispozitive ( de exemplu aprinderea automata a unui bec printr-un microcontroler), ci si posibilitatea gestionarii mai eficiente a informatiilor digitale (de la date bancare la fisiere media).

Aceasta a devenit treptat principala sursa legitima de inregistrare a oricarui tip de date (text, imagini, video, audio). Pentru a accesa si actualiza informatia digitala, s-au dezvoltat diferite versiuni de baze de date centralizate si distribuite.

Bazele de date centralizate sunt aplicatii software specializate ce folosesc resursele sistemului (a statiei) pentru a raspunde cat mai rapid interogarilor. Statiile trebuie sa detina multa putere de stocare si de procesare in comparatie cu un sistem de calcul normal destinat utilizatorilor casnici. In alta ordine de idei, s-au dezvoltat si baze de date distribuite, menite sa reduca din capacitatile tehnice ale serverului si sa stocheze informatia sub forma descentralizata. Cautarea resursei in acest context ar presupune interogarea recursiva a fiecarei entitati pana la gasirea sa. Prin acest mod, nu doar ca statiile pot avea si capacitati tehnice mai reduse, dar si pot pastra copii de rezerva (backup) locale pentru fiecare segment de informatie in parte.

Această evoluție naturală spre descentralizare a deschis drumul unor concepte moderne precum învățarea automată federată (federated learning), unde datele nu mai sunt transferate către un server central. În schimb, modelele sunt antrenate local, iar parametrii sunt ulterior agregați global. Astfel, se menține confidențialitatea datelor, fără a compromite performanța modelului.

Evolutia tehnologica continua a dat nastere la o serie de atacuri cibernetice menite sa destabilizeze securitatea aplicatiilor si totodata sustragerea a cat mai multe date sau identitati private in contradictie cu normele legale. Cele mai populare atacuri raportate la scara globala pentru anul curent 2025 sunt ransomware (conform <sup>1</sup>, in SUA s-au raportat cresteri de 149%), furtul de identitate prin exfiltrarea de credentiale, si phishing. Dezvoltarea modelelor de inteligentă artificială a amplificat aceste riscuri, oferind atacatorilor instrumente automatizate pentru generarea și adaptarea atacurilor.

Pentru a limita utilizarea abuzivă a tehnologiilor bazate pe AI, Uniunea Europeană a adoptat în 2024 un set de reglementări stricte privind integrarea acestor module în aplicațiile software, prin AI Act <sup>2</sup>.

Progresul din domeniul machine learning si a Large Language Models a fost posibil ca urmare a unui volum masiv de date disponibile si a nevoii tot mai mari de analiza. Acest lucru a determinat aparitia unei noi categorii de specialisti, data scientists, dedicati colectarii si prelucrarii minutioase a datelor pentru antrenarea modelelor.

Totuși, pe măsură ce investițiile în tehnologii AI au crescut, au apărut și actori rău intenționați care încearcă să exploateze vulnerabilitățile din procesul de antrenare. Întrucât modelele moderne depind de calitatea datelor folosite, acestea au devenit o țintă principală a atacurilor. Atacatorii se regasesc si ei intr-o pozitie constanta de adaptare la noile formalitati de securitate si incearca sa contracareze fiecare element nou. Astfel, avand in vedere complexitatea dezvoltarii unui modul de inteligenta artificiala specializat pe diferite domenii, tinta s-a redirectionat spre volumul de date pe care acestea le folosesc si care pot determina starea finala a aplicatiei.

În contextul învățării automate distribuite, literatura de specialitate identifică trei categorii majore de atacuri:

- Atacuri asupra datelor, precum data poisoning, unde setul de antrenare este manipulat pentru a altera comportamentul modelului;
- Atacuri asupra modelului, prin modificarea parametrilor sau a gradientului (de exemplu, model poisoning);
- Atacuri asupra canalului de comunicare, care vizează interceptarea sau modificarea mesajelor dintre entitățile participante.

Lucrarea de față se concentrează pe prima categorie, data poisoning, în cadrul unei infrastructuri de

invatare federate.

## 1.2 Motivatia lucrarii

Avand in vedere aspectele legate de posibilitatea unei interventii asupra setului de date de antrenare, atac denumit otravire a datelor (data poisoning), munca cercetatorilor s-a ingreunat. Preocuparea nu mai este primordial asupra analizei setului de date de antrenare, cat despre mentinerea integritatii si a confidentialitatii lor. Pentru a raspunde acestor nevoi, colaborarea dintre cercetatori s-a orientat către modele distribuite de lucru, iar învățarea federată (federated learning) a devenit una dintre principalele direcții. Aceasta permite colaborarea între participanți fără a partaja direct seturile lor de date, menținând o barieră naturală împotriva accesului neautorizat. Totuși, deși infrastructura este diferită față de abordările centralizate, vulnerabilitățile rămân, iar atacurile asupra datelor utilizate local pot afecta modelul global.

In urma unei analize proprii, am putut observa diferite solutii/frameworks de simulare a procesului de invatare automata federata, dar fara o integrare cu mecanisme moderne de testare pentru atacuri precum otravirea datelor (data poisoning) amintite anterior <sup>3</sup>. Unele dintre aceste framework-uri sunt poate dificil de gestionat si configurat <sup>4</sup>, si nu permit extinderea usoara prin integrare altor componente. In acelasi timp, gandindu-ne la multitudinea de atacuri malware si la platformele de detectie a lor, devine clar ca in domeniul inteligentei artificiale lipseste o platforma centralizata, flexibila, dedicata testarii si evaluarii cu diferite tipuri de atacuri asupra modelelor distribuite.

Aceste limitări justifica realizarea prezentei lucrari, care își propune dezvoltarea unei platforme de simulare capabile sa testeze atacuri de tip data poisoning într-o infrastructura de învățare federată.

## 1.3 Obiectivele lucrării

Plecand de la neajunsurile prezentate, ne propunem in aceasta lucrare sa venim in sprijinul comunitatii de cercetare stiintifica in domeniul securizarii solutiilor cu AI cu o platforma de simulare cu sursa deschisa ("open source"), a acestei clase de atacuri pe mai multe directii. Astfel, oferim cercetatorilor posibilitatea analizei algoritmului de antrenare propriu dezvoltat, plecand de la o retea neuronală de baza si un set de date uzual (imagini), si testarea sa prin antrenare in diferite conditii. Platforma in sine respecta toate normele unei aplicatii software de productie, in care fiecare actiune are propria sa logica de implementare. Serviciile sunt segregate suficient de mult incat sa permita o dezvoltare ulterioara prin integrarea lor cu alte sisteme.

Rezultatele pot fi utile in contextul securizarii procesului de antrenare al algoritmului, dar si pentru analiza factorilor de risc la care e expus in acest mediu.

Cercetatorul este cel care furnizeaza algoritmul python de antrenare a propriei retele neuronale sau algoritmul de Machine Learning. El seteaza parametrii simulării atat pentru procesul de antrenare, cat si pentru tipul de atac de otravire a datelor. Platforma isi propune sa simuleze acest tip de atac cu ajutorul acestor setari de inceput intr-un mediu de invatare federata, furnizand la final o comparatie intre modelul antrenat folosind datele normale de antrenare si cel antrenat cu datele otravite. Aceste rezultate pot fi utile in semnalarea unui posibil risc la nivelul modelului dezvoltat, oferind mai apoi solutii de imbunatatire a implementarii sale.

## 1.4 Structura lucrării

---

<sup>1</sup><https://www.dnsc.ro/vezi/document/buletin-de-indicatori-statistici-si-tendinte-de-securitate-cibernetica-h1-2025>

<sup>2</sup><https://artificialintelligenceact.eu/wp-content/uploads/2024/11/Future-of-Life-InstituteAI-Act-overview-30-May-2024.pdf>

<sup>3</sup><https://ibmfl-api-docs.res.ibm.com/index.html>

<sup>4</sup><https://github.com/IBM/federated-learning-lib/tree/main>

# Capitolul 2:

## Notiuni Teoretice

În acest capitol, vor fi prezentate notiunile teoretice specifice înțelegerii procesului de dezvoltare a platformei de simulare. Vom începe cu Notiunile introductive despre conceptele de Machine Learning în antiteza cu Deep Learning. În continuare, vom discuta despre învățarea federată și arhitectura unei infrastructuri federate de învățare automată, tipurile de atacuri data poisoning implementate în procesul de simulare a atacurilor, precum și alte notiuni specifice implementării.

### 2.1 Notiuni introductive

Machine Learning și Deep Learning sunt două ramuri importante ale Inteligenței Artificiale care au rolul dezvoltării unor modele specifice rezolvării unor anumite acțiuni. Pornind de la antrenarea de rețele neuronale, ne orientăm atenția spre setul de date de antrenare și spre actorii ce pot interveni în acest proces. Mediul în care testăm oferă o perspectivă reală asupra impactului pe care îl pot avea aceste atacuri la nivelul unei organizații sau aplicații.

#### 2.1.1 Diferența dintre Machine Learning și Deep Learning

Inteligența Artificială (AI) este domeniul vast care înglobează orice tehnică ce permite calculatoarelor să imite comportamentul uman. Informația a evoluat treptat odată cu îmbunătățirea capacităților de stocare ale dispozitivelor și apariția programelor software complexe. De la simplul format de text, înregistrări audio, până la imagini și video în rezoluții 4K, modul de lucru s-a diversificat constant.

La fel au evoluat și cerințele utilizatorilor, care tind să acceadă către soluții automate care să le rezolve problemele uzuale, precum identificarea de patterns în imagini sau chiar din video, sau generarea de text.

IA vine să rezolve aceste probleme și să introducă algoritmi de rezolvare specifici pentru fiecare tip de informație furnizată.

Machine Learning este o componentă importantă din domeniul IA care se diferențiază de alte metode de antrenare prin optimizările pe care le aduce erorilor ce apar din predicția rezultatului corect. Modelele de ML clasice se bazează pe intervenția umană în factorul de decizie (supervised learning), mai precis datele de intrare sunt etichetate pentru a oferi un context de predicție stabil.

Deep Learning este o subcategorie a Machine Learning, care are rolul de a minimiza intervenția umană și a automatiza procesul de decizie. Prin această metodă se automatizează mare parte din extragerea caracteristicilor pe setul de date, eliminând nevoia de a defini etichete pentru fiecare valoare de intrare (unsupervised learning).

Diferența dintre aceste două concepte este în modul în care acești algoritmi învață și procentul de utilizare a datelor [1]. Scopul principal al învățării automate este predicția. Pe baza unui set de date de antrenare și de testare, se determină o anumită categorie de ieșire predefinită.

#### 2.1.2 Rețea Neuronala

Rețelele Neuronale sunt un subset al Machine Learning și se identifică drept infrastructura de bază din cadrul algoritmilor de Deep Learning. Denumirea de "neural" se referă la structura lor internă, în care fiecare caracteristică (feature) este un neuron ce interacționează unii cu alții. Ele sunt compuse din 3 straturi/layers: primul strat îl reprezintă stratul nodurilor de intrare, al doilea strat este denumit "stratul ascuns" (hidden layer) pt că încapsulează mai multe straturi, iar ultimul strat este cel de ieșire în care se face predicția propriu-zisă. Straturile ascunse sunt concepute pentru a procesa iterativ datele pornind de la starea lor din nodurile de intrare până la stratul de ieșire.

### 2.2 Învățare automată federată

Evoluția hardware în tehnologie a condus la creșterea numărului dispozitivelor mobile (telefoane, tablete), denumite gadgets din faptul că sunt mici, portabile și moderne. Ele au fost mai departe adoptate la scară largă, devenind obiecte indispensabile în era tehnologică ce avea să vină.

### 2.2.1 Concept

Invatarea automata federata permite lucrul cu modele de ML sau chiar retele neuronale, antrenate distribuit, pe un numar mare de dispozitive in scopul rezolvarii unei probleme de IA. Distribuirea sarcinilor a fost adoptata si in contextul opozitiei lucrului centralizat, pe servere ce detin capabilitati Hardware performante (placi grafice de ultima generatie), dar care genereaza costuri mari si care pot fi predispuse la amenintari de securitate cibernetica, fiind considerate SPOF(Single Point of Failure).

### 2.2.2 Arhitectura FL

In literatura, exista mai multe categorii de arhitecturi de invatare automata federata. In aceasta sectiune ne vom concentra pe clasificarea generala a arhitecturii unei aplicatii folosind federated learning, si vom enumera pentru o anumita categorie cum se clasifica dispozitivele utilizate.

Federated Learning, asa cum a mai fost mentionat, este organizat dintr-un server (agregator) si multiple dispozitive client. Modul in care aceste entitati comunica este fundamentul principal in modului de imbunatatire al invatarii.

In modul clasic al federated learning, dispozitivele client transmit actualizari ale modelului de baza la un server central care aplica asupra lor o functie de agregare, reconstruind intreg modelul de baza. Aceasta setare/model, presupune de fapt o delegare a sarcinii de invatare, insa pastreaza entitate centrala necesara imbunatatirii solutiei. Acest fapt, nu tine sa evita posibilitatea amenintarilor cibernetice (Single Point of Failure), ci doar sa usureze costurile centralizatorului in a procesa local problema, distribuind sarcinile.

Modelul Fully decentralized (peer-to-peer) learning, ofera o noua abordare si rezolva problema cibernetica amintita. In aceasta setare nu exista agregator, imbunatatirile fiind comunicate intre clienti interconectati. Ideea principala se bazeaza pe inlocuirea comunicarii cu agregatorul cu cea intre dispozitive individuale printr-un protocol prestabilit. In functie de numarul de dispozitive, se concepe un graf de conexiuni in care fiecare nod reprezinta un client, iar fiecare muchie un canal de comunicatie. Restrictia principala este ca un dispozitiv sa fie conectat la un numar maxim limitat de dispozitive adiacente, prestabilit, in contradictie cu un graf complet (stea) specific arhitecturii clasice client-server. Nodurile isi imbunatatesc propriile variante ale retelei, si isi comunica rezultatele pe care le agrega local, realizand o medie a ponderilor. In comparatie cu modelul federated learning clasic, modelul fully decentralized nu specifica de la inceput dispozitivelor un model de baza global de la care sa porneasca in procesul de rezolvare a problemei.

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

Figura 2.1: Imaginea 2.2.2: Modele Federated Learning

Imaginea de mai sus ofera o privire de ansamblu asupra celor doua modele de arhitecturi si caracteristicile acestora.

Model	Avantaje	Dezavantaje
<b>Federated learning</b> (centralized coordination)	<ul style="list-style-type: none"> <li>• mai simplu de configurat (topologie hub-and-spoke)</li> <li>• pornește de la un model global de bază</li> <li>• agregarea centralizată reduce sarcina de calcul pe clienți</li> <li>• necesită mai puține conexiuni (doar client <math>\rightarrow</math> server)</li> <li>• gestiunea și monitorizarea sunt mai simple</li> </ul>	<ul style="list-style-type: none"> <li>• SPOF (Single Point of Failure) – serverul central</li> <li>• serverul poate deveni țintă pentru atacuri</li> <li>• nu elimină riscurile cibernetice, doar distribuie munca</li> <li>• dependența de coordonator pentru progresul antrenării</li> <li>• necesită infrastructură centralizată permanent disponibilă</li> </ul>
<b>Fully decentralized / peer-to-peer learning</b>	<ul style="list-style-type: none"> <li>• previne atacurile specifice unui server central (evită SPOF)</li> <li>• reziliență crescută – compromiterea unui nod nu debilizăază întregul sistem</li> <li>• îmbunătățirile se propagă prin graf, fără entitate centrală</li> <li>• agregare locală (fiecare nod mediază ponderile)</li> <li>• poate scala natural dacă graful este bine proiectat</li> </ul>	<ul style="list-style-type: none"> <li>• necesită conexiuni suplimentare între clienți</li> <li>• topologie complexă, dificil de administrat</li> <li>• nu există model global inițial furnizat tuturor</li> <li>• performanța depinde de calitatea grafului de conexiuni</li> <li>• nodurile malițioase pot influența direct vecinii</li> <li>• necesită protocoale suplimentare pentru consistența actualizărilor</li> </ul>

Tabel 2.1: Compararea modelelor Federated Learning și Fully Decentralized Learning

În tabelul de mai sus, sunt evidențiate avantajele și dezavantajele utilizării celor două tipuri de modele federated learning.

În continuarea acestei lucrări, se va discuta preponderent despre modelul clasic federated learning, fiind unul adoptat la scară largă și care oferă performanțe bune raportat la costurile de producție.

Modelul clasic la rândul său, se cataloghează în literatură în funcție de tipul dispozitivelor care iau parte la procesul de antrenare. Sub acest filtru, există Cross-Device Federated Learning și Cross-Silo Federated Learning.

Cross-Device Federated Learning sunt dispozitivele client IOT uzuale, individuale, care comunică orchestratorului printr-un protocol prestabilit.

Cross-Silo Federated Learning sunt dispozitive din instituții guvernamentale, companii, sau centre de date distribuite geografic. Instituțiile nu doresc să schimbe informații între ele sau cu un furnizor de servicii central, păstrându-și confidențialitatea, folosind federated learning pentru a antrena propriul model pe datele private ale fiecăruia.

### 2.2.3 Procesul de antrenare FL

Primul pas este stabilirea conexiunii dintre dispozitive și un server de agregare ce permite antrenare distribuită a tipului de rețea neuronală sau model ML specific problemei. Odată stabilit canalul, în faza de configurare inițială, serverul trimite dispozitivelor starea de bază a rețelei neuronale, ponderile, în vederea antrenării individuale. Fiecare rețea se antrenează cu datele extrase local (on device) și își îmbunătățește configurația internă la fiecare epocă pentru o perioadă de timp bine determinată.

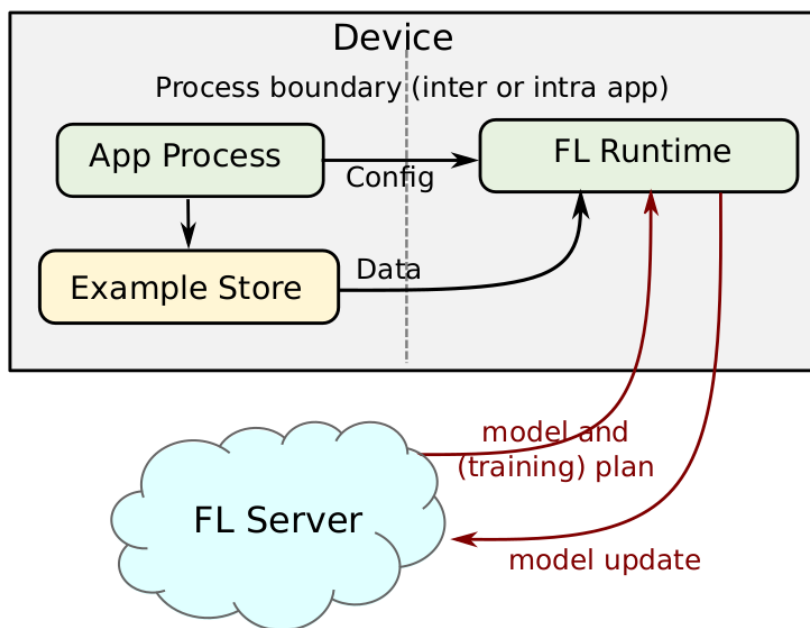


Figura 2.2: Imagine 2.2.3.1: Arhitectura internă a unui dispozitiv

Figura de mai sus descrie operațiile specifice programului Software care se ocupă de antrenarea rețelei/modelului. Putem observa cum dispozitivele primesc un plan de antrenare de bază pe care îl vor antrena local pe un set de date limitat.

Deși acest pas nu aduce un procent de îmbunătățiri foarte mari, în faza următoare, dispozitivele vor transmite configurațiile curente ale rețelelor lor la orchestrator (server). Rutina FL\_Runtime extrage configurația nouă locală și îi comunică serverului pentru o posibilă actualizare a sa.

În cele din urmă, entitatea centrală combină toate aceste ponderi aplicând o funcție de agregare și în cazul îmbunătățirii setului de ponderi, modifică configurația de bază și o retransmite dispozitivelor pereche. Dacă ponderile noi nu se îmbunătățesc semnificativ față de configurația de bază, atunci se păstrează aceasta din urmă, iar în caz contrar se actualizează cu noile ponderi.

În figura de mai sus, se pot observa într-o manieră continuă, fluxul de comunicație dintre dispozitive și serverul agregator, precum și operațiile specifice fiecărei entități dintr-o rundă de mesaje.

Securitatea protocoalelor de agregare, utilizate în comunicații dintre clienți și orchestrator, este o componentă importantă în procesul federated learning. De menționat este faptul că, în această topologie,

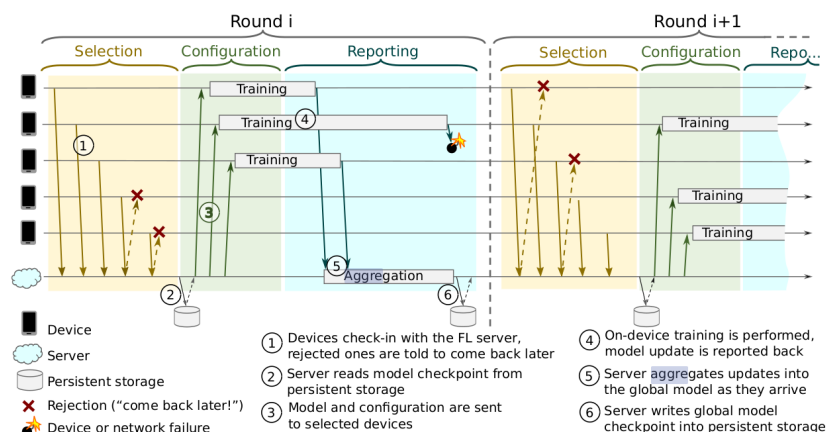


Figura 2.3: Imagine 2.2.3.2: Procedee in invatare automata federata

comunicatiile au loc criptat, folosind metode specifice precum criptare homomorfica, sau chiar OTP, insa securitatea datelor de pe dispozitive ramane la latitudinea acestuia.

## 2.2.4 Exemple in viata reala

Federated Learning s-a extins rapid în numeroase domenii datorită capacității sale de a antrena modele performante fără a colecta sau centraliza date sensibile. Prin păstrarea informațiilor la nivelul fiecărui dispozitiv sau instituții, FL reduce riscurile asociate scurgerilor de date și permite colaborarea între entități care altfel nu ar putea împărtăși date brute. În continuare sunt prezentate câteva exemple reprezentative ale utilizării sale în aplicații din lumea reală.

### Industrie și IoT

- **Mentenanță predictivă:** Vehiculele moderne, utilajele industriale și echipamentele IoT generează constant date despre starea componentelor. FL permite antrenarea unui model comun care poate prezice momentul oportun pentru realizarea mentenanței fără a colecta date brute de la fiecare dispozitiv.
- **Dispozitive de monitorizare:** Senzori portabili și dispozitive smart home pot furniza statistici privind activitatea sau consumul energetic, păstrând datele utilizatorilor la sursă.

### Medical

- **Diagnostic, prognoză și imagistică:** FL este folosit în spitale și clinici pentru detectarea celulelor canceroase din imagini RMN, CT sau radiografii, fără transferul imaginilor către un server central.
- **Confidențialitate menținută la sursă:** Fiecare instituție medicală antrenează local o parte din model, partajând doar actualizările, ceea ce permite colaborarea fără a încălca regulile privind datele pacienților.

### Financiar

- **Detectarea fraudelor:** Instituțiile financiare pot îmbunătăți detectarea tranzacțiilor suspecte analizând tipare comune fără a expune date sensibile despre clienți.

### Servicii și experiență utilizator

- **Recomandări personalizate:** Platformele de streaming și aplicațiile mobile generează recomandări local, pe dispozitiv, fără a trimite istoricul complet al utilizatorului către server.
- **Analiză comportamentală:** FL poate analiza activitatea utilizatorilor pentru a sugera rutine sănătoase sau îmbunătățiri ale stilului de viață, păstrând confidențialitatea datelor.

## Securitate și privacy

- **Supraveghere fără expunerea datelor sensibile:** Modelele de recunoaștere facială pot fi antrenate fără a transmite imagini reale, doar parametrii aferenți.
- **Analiză a sentimentelor:** FL poate analiza reacțiile utilizatorilor la evenimente sociale (like-uri, share-uri, comentarii) fără colectarea directă a acestor date de către platformă.

## 2.3 Atacuri de tip Data Poisoning

În acest capitol se va discuta una dintre avantajele pe care le ofera mediul federated learning atacatorilor și ce înseamnă acest lucru pentru fluxul configurației modelului. Vom începe cu definirea vectorilor de atac și concentrarea pe una dintre categorii, data poisoning. În continuarea lucrării, vom analiza impactul pe care îl are acest tip de atac asupra infrastructurii federate de învățare și riscurile pe care le introduce, precum și câteva propuneri de identificare și constientizare a existenței sale.

### 2.3.1 Definirea tipurilor de atac

În literatura de specialitate, atacurile asupra unui model de ML sau asupra unei rețele neuronale sunt definite drept atacuri adversariale (adversarial attacks). Această clasă de atacuri are rolul de a produce modificări în comportamentul normal al modelului într-un mod indizibil. Conform lucrării [1], în funcție de nivelul de scop al unui atac, putem avea:

- **Atacuri fără țintă (untargeted attacks)** Scopul este reducerea acurateții generale a modelului sau chiar destabilizarea completă a acestuia. Un exemplu pentru clasificarea imaginilor este introducerea unui zgomot care degradează calitatea setului de date. O altă metodă este modificarea etichetelor din setul de antrenare, de exemplu atribuirea etichetei *leu* unor imagini care în mod normal ar trebui etichetate drept *pisică*.
- **Atacuri țintite (targeted attacks)** Denumite și *backdoor attacks*, deoarece urmăresc modificarea comportamentului modelului doar pentru un anumit subset de date, fără a afecta vizibil acuratețea globală. Continuând exemplul anterior, pentru un anumit tip de imagini ce reprezintă pisici într-o anumită poziție, se poate introduce un artefact vizual menit să păcălească modelul și să clasifice pisica drept *leu*. Astfel apare un backdoor activ doar pentru acele imagini precise.

Plecând de la această categorisire, există o multitudine de modalități prin care un atacator poate submina capacitatea de predicție a modelului. Mediul federated learning introduce prin construcție o serie de întrebări la care trebuie găsit un răspuns pentru a determina prin ce moduri un adversar se poate infiltra și poate profita de anumite drepturi pentru a introduce incertitudine în antrenarea sau reglarea (fine-tuning) a modelelor.

### 2.3.2 Vectori de atac

Analiza amenințărilor asupra modelelor de ML, a introdus o serie de posibile vulnerabilități asupra componentelor ce alcătuiesc fluxul de învățare. Un atacator își poate alege zona de interes, pe baza posibilităților de exploatare a sistemului respectiv.

Vectorii de atac cunoscuți sunt:

- **Data Poisoning** Când adversarul încearcă să corupă setul de date antrenare cu scopul defectării modelului încă de la început.
- **Model Update Poisoning** Când adversarul se folosește de o vulnerabilitate ce îi permite modificarea configurației parametrilor trimisi către orchestrator.
- **Evasion attack** Când adversarul are acces la datele de testare și le poate modifica în momentul inferenței.

În funcție de gradul de acces la sistemele gazdă, modelul poate fi inspectat în diferite moduri:

- **Black Box** Adversarul nu are abilitatea să inspecteze parametrii modelului înainte sau în timpul atacului.



- **Stale Box** Adversarul poate inspecta doar o versiune incipienta a modelului. Aceasta capabilitate apare si in federated learning cand adversarul are acces la rundele de antrenare ale clientului.
- **White Box** Adversarul are abilitatea de a inspecta direct parametrii modelului. Acest scenariu se bazeaza pe un grad de acces superior al adversarului asupra sistemului.

Pe baza acestor scenarii, atacatorul poate aplica o serie de tehnici pentru modificarea starii modelului. In contextul lucrarii de fata, se va discuta scenariul stale box, adversarul avand posibilitatea doar de a introduce un atac Data Poisoning in rundele de antrenare, datele corupte fiind introduse o singura data in procesul de reglare a modelului (fine-tuning).

### 2.3.3 Atacul Data poisoning

Acest tip de atac presupune coruperea datelor de antrenare sau de testare, prin diferite tehnici specifice, la nivelul dispozitivului clientului. In aceasta paradigma, atacul poate fi considerat la fel de bine targeted sau untargeted intrucat depinde de intentia adversarului si de potentialul risc in divulgarea punctului de exploatare de care dispune.

În practica atacurilor de tip Data Poisoning, adversarul poate manipula atât conținutul datelor, cât și etichetele acestora, efectele fiind de obicei greu de observat la nivel local. În mediul federated learning această dificultate este amplificată, deoarece orchestratorul nu are acces direct la datele brute ale clienților. Astfel, orice modificare realizată pe un dispozitiv compromis intră automat în procesul de antrenare, fiind tratată ca o contribuție legitimă. În mod particular, chiar și un număr redus de exemple otrăvite poate introduce un comportament persistent în modelul global, mai ales dacă atacul este repetat pe durata mai multor runde

In contextul solutiei propuse, se va discuta despre impactul atacului asupra unui set de imagini si tipurile de metode pentru a le altera, asa cum se poate vedea in tabelul 2.3.3.

Tip atac	Descriere
<b>Gaussian noise</b>	Introducerea unui zgomot aleator în imagini sau în vectorii de caracteristici, cu scopul degradării calității datelor și scăderii performanței modelului.
<b>Label flip</b>	Modificarea intenționată a etichetelor din setul de antrenare, astfel încât exemple corecte sunt asociate cu clase greșite, afectând procesul de învățare.
<b>Backdoor injection</b>	Inserarea unui artefact vizual sau a unui tipar specific într-un subset mic de date, astfel încât modelul să învețe un comportament anormal activat doar de acel trigger.

Tabel 2.2: Tipuri de atacuri Data Poisoning

Prin aceste tipuri de atacuri adversariale, impactul asupra modelului are loc pe o perioada determinata de timp, de obicei mai lunga, si produce variatii in predictia finala.

### 2.3.4 Impactul poisoning in Federated Learning

Mediul de invatare federata a introdus o serie de amenintari cibernetice preponderent la nivelul dispozitivelor clientilor, acestea fiind cele mai vulnerabile din punctul de vedere al aplicarii unui atac de otravire a datelor. Clientii sunt producatorii unui model calitativ care sa ofere predictii legitime in diferite scenarii.

Cand se discuta despre alterarea etichetelor (Label Flip Attack) atunci la o prima vedere, utilizatorul nu si-ar da seama decat in urma unei inspectii amanuntite. Pentru acest tip de atac, exista tehnici de verificare si filtrare ce pot determina daca un tip de informatie este catalogata corect inainte de antrenare.

Efectele unui data poisoning pot persista chiar și după eliminarea datelor corupte, deoarece modelul învață un tipar greșit care nu dispare imediat fără o reantrenare completă. Acest lucru este relevant mai ales pentru atacurile backdoor, care rămân inactive până la apariția unui trigger vizual, fără a afecta acuratețea generală. Din acest considerent, atacurile tintite (targeted attacks) sunt deosebit de periculoase in contextul unui mediu de invatare federata pentru ca datele corupte se ascund in interiorul configuratiilor particulare ale clientilor. Aceste configuratii sunt transmise mai departe la agregator care aplicand functia sa de agregare, amplifica negativ starea modelului.

Având în vedere aceste particularități, devine esențială analiza metodelor de apărare și a mecanismelor prin care pot fi detectate contribuțiile malițioase. Provocarile in acest domeniu conduc la o serie tot mai mare de utilitare sau platforme ce permit detectarea facila a acestor atacuri si imbunatatirea evenimentelor cu un potential risc in organizatii.

Mediul federated learning introduce un risc suplimentar: un adversar care controlează un număr mic de clienți poate influența disproporționat modelul global dacă este integrat într-un moment critic al antrenării. În absența unor mecanisme robuste de apărare, actualizările malițioase sunt tratate ca fiind legitime, iar agregatorul nu are nicio modalitate directă de a le verifica.

Un alt efect important al acestui tip de atac este degradarea treptată a performanței modelului. În scenariile untargeted, scăderea acurateții globale poate trece neobservată în primele runde de antrenare, dar devine evidentă odată ce modelul converge către o reprezentare eronată a datelor. În scenariile targeted, atacul poate compromite decizii critice doar într-un subset de cazuri, ceea ce face detectarea mult mai dificilă și impactul mult mai nociv, mai ales în aplicații sensibile cum ar fi securitatea, domeniul medical sau sistemele autonome.

## 2.4 Alte Notiuni

În vederea elaborării soluției propuse în această lucrare, se vor aminti celelalte concepte care stau la baza implementării propriu-zise. În acest capitol se vor detalia succint mecanismele ce stau la baza platformei propuse, modul de utilizare și scopul alegerii lor.

### 2.4.1 Docker

Docker reprezintă un set de servicii software de tip platformă ce utilizează virtualizarea la nivel de sistem de operare pentru a crea entități independente, numite containere. Aceste containere sunt create specific pentru a întreprinde anumite acțiuni și oferă un mediu izolat de execuție.

Un container este o instanță software ce vine împachetată cu programul aplicației și toate bibliotecile necesare dezvoltării ei. O imagine este vizualizată drept un șablon de instrucțiuni pentru crearea unui container cu un anumit tip de bibliotecă necesară dezvoltării unei aplicații. De obicei, imaginile sunt construite pe baza unui fișier denumit Dockerfile care propune o serie de comenzi pentru crearea unui mediu personalizat.

Toate containerle Docker rulează prin intermediul Docker Engine, un serviciu ce rulează la nivelul sistemului de operare și oferă suport cross-platform (Linux, Windows, macOS).

Containerizarea diferă de virtualizarea tradițională prin faptul că containerele partajează același kernel al sistemului de operare gazdă, în timp ce mașinile virtuale (VM) necesită fiecare un sistem de operare complet. Această abordare face containerele Docker mult mai ușoare (de ordinul MB vs GB) și mai rapide la pornire (secunde vs minute) comparativ cu VM-urile.

Arhitectura platformei Docker este asemănătoare modelului client-server, fiind compusă din următoarele componente:

- **Dockerd** un proces daemon, identificabil drept server, ce gestionează tot fluxul de servicii, de la imagini și containere până la volume și rețele.
- **Application Binary Interfaces (API)** o suită de interfețe de comunicare și control al serverului
- **Comanda docker** o interfață în linie de comandă, docker

În contextul platformei dezvoltate, Docker oferă un mediu de dezvoltare automat în care se regăsește gazduita platforma publică de simulare. Toate informațiile despre fiecare simulare sunt stocate la nivelul unei baze de date PostgreSQL, în timp ce aplicația web este publică printr-un container frontend, iar în spate regăsim un container backend pentru transmiterea de comenzi serverului.

Comunicatia dintre containere are loc în aceeași rețea locală Docker, iar informațiile despre fiecare simulare a clientului persistă în același volum partajat.

Pe lângă Dockerfile, serviciul Docker Engine oferă și posibilitatea creării unei configurații prestabilite pentru definirea de topologii de rețea de containere prin Docker Compose. Pe baza fișierelor Dockerfile în care sunt definite șabloanele imaginilor și a unui fișier de configurare YAML a topologiei relațiilor dintre containere, serviciul Docker Compose permite implementarea unei infrastructuri întregi prin rularea și ștergerea sa dintr-o serie de comenzi.

### 2.4.2 Python

### 2.4.3 Rest API

## Capitolul 3:

# Proiectare, Implementare si Testare

### 3.1 Cerintele Software

#### 3.1.1 Cerintele functionale

#### 3.1.2 Cerintele nefunctionale

### 3.2 Arhitectura platformei

#### 3.2.1 Containere

#### 3.2.2 Server

### 3.3 Testare

## Capitolul 4:

# Rezultate si Metrice Simulari

### 4.1 Evaluare Performante

#### 4.1.1 Scalabilitatea Simularilor

#### 4.1.2 Scalabilitatea platformei

### 4.2 Evaluare Rezultate

#### 4.2.1 Performante Gaussian Noise

#### 4.2.2 Performante Label-Flip

#### 4.2.3 Performante Backdoor

## Capitolul 5:

# Concluzii si dezvoltare ulterioara

### 5.1 Starea Curenta

### 5.2 Dezvoltare Ulterioara

### 5.3 Tabele

Tabelele sunt aranjări a informației într-o structură formată din linii și coloane, care permite o mai bună observare a acesteia.

Mai jos apar două exemple. Primul tabel este de dimensiune mică. Al doilea, din cauza dimensiunii mai mari, are o orientare inversată și este plasat singur pe o pagină.

Nume Complet	Funcție Ocupată
Joshua Roob	Manager de Proiect
Asa Hauck	Artist Grafic
Harley Hagenes	Programator

Tabel 5.1: Colaboratori la Realizarea Studiului

Stat	Oras	Latitudine	Longitudine
South Carolina	Corwinberg	86.609523	42.408007
Rhode Island	East Isaacmouth	63.17309	-13.786023
Mississippi	North Noblestad	-31.316834	5.280483
Illinois	Grahamland	-39.853659	-77.713676
Rhode Island	West Richardfurt	67.583131	31.858455
Florida	Port Roberta	25.276026	83.715344

Tabel 5.2: Locații de Conducere a Studiului

## 5.4 Imagini

Imaginile sunt utilizate în cadrul lucrării pentru exemplificarea unor idei în manieră vizuală.

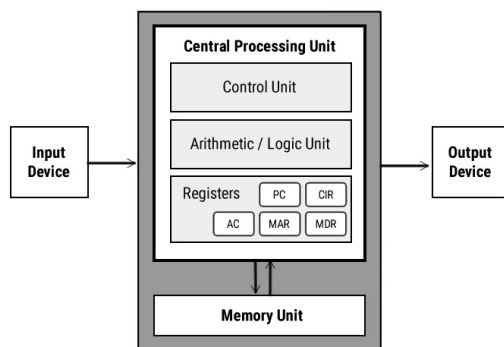


Figura 5.1: Arhitectura unui calculator<sup>1</sup>

## 5.5 Liste

Listele sunt simple serii de informații.

- Un item
- Unul dintre itemi
- Încă un item

Acestea pot conține itemi identificați prin numere dacă indexarea sau sortarea sunt necesare.

1. Primul item
2. Al doilea item
3. Al treilea item

## 5.6 Formule Matematice

$\text{\LaTeX}$  oferă un mod programatic de a construi formule matematice, după cum este cea de mai jos.

<sup>1</sup>Arhitectura ilustrată este de fapt cea von Neumann.

$$\sum \mathbf{F} = 0 \Leftrightarrow \frac{d\mathbf{v}}{dt} = 0$$

## 5.7 Note de Subsol. Citări

Notele de subsol pot fi utile în cazul explicațiilor suplimentare (cum a fost cea referitoare la imaginea inclusă, la care sintaxa este puțin diferită din cauza plasării notei în cadrul legendei) sau a citărilor<sup>2</sup> care nu se pretează a fi trecute în bibliografie din cauza utilizării lor punctuale.

Pe de altă parte, sursele bibliografice citate intens [1] sunt marcate corespunzător și notate în bibliografie.

## 5.8 Etichete. Referințe

În cadrul surselor  $\text{\LaTeX}$  a acestui document, apar *tag*-uri `\label` care creează o etichetă utilă referințelor interne. Acestea din urmă indică elemente din cadrul documentului curent (de exemplu, către tabelul 5.1).

Mai pot apărea referințe externe, către resurse din Internet (de exemplu, către *website*-ul Wikipedia).

---

<sup>2</sup>Cristian Lupascu, Cezar Plesca și Mihai Togan, “Privacy Preserving Morphological Operations for Digital Images”, în (iun. 2020), pp. 183–188, doi: 10.1109/COMM48946.2020.9141997



# Bibliografie

## Cărți

- [1] Mihai Togan, *Cryptographic Technologies for Data Protection in Cloud*, Editura Matrix Rom, 2017.

## Articole Științifice

- [2] Ionuț Dumitru și Mihai Togan, “Client Module with Multifactor Authentication for Remote Electronic Signature Generation Using Cryptography API: Next Generation”, în *Journal of Military Technology* 3 (iun. 2020), pp. 5–10, DOI: 10.32754/JMT.2020.1.01.