

**МАГІСТЕРСЬКА РОБОТА**

**МР. ШМ - 21.00.00.000 ПЗ**

**Група ШМ-24-2**

**Кантор Матвій**

**2025**

**Івано-Франківський національний технічний університет нафти і газу**

Факультет інформаційних технологій

Кафедра інженерії програмного забезпечення

**Кантор Матвій Володимирович**

(прізвище, ім'я, по батькові)

УДК 004.9  
(індекс)

## **МАГІСТЕРСЬКА РОБОТА**

**Моделі та методи контролю якості та автентичності мережевих даних**

(назва роботи)

**Інженерія програмного забезпечення**

(назва освітньої програми)

**121 - Інженерія програмного забезпечення**

(шифр і назва спеціальності)

**Кантор М.В.**

(підпис, ініціали та прізвище здобувача освітнього ступеня)

**Науковий керівник Піх Володимир Ярославович, к.т.н., доцент**

(прізвище, ім'я, по батькові, науковий ступінь, вчене звання)

**Допущено до захисту**

Завідувач кафедри

доц. Бандура В.В.

(посада) (підпис) (дата) (ініціали та прізвище)

**Нормоконтроль**

доц. Вовк Р.Б.

(посада) (підпис) (дата) (ініціали та прізвище)

Робота містить результати власних досліджень. Використання ідей, результатів і текстів інших авторів мають посилання на відповідне джерело

Івано-Франківськ – 2025

**Івано-Франківський національний технічний університет нафти і газу**

Факультет інформаційних технологій

Кафедра інженерії програмного забезпечення

Освітній рівень магістр

Спеціальність 121 – Інженерія програмного забезпечення

ЗАТВЕРДЖУЮ:

Зав. кафедрою

ІПЗ

доц.

В.В. Бандура

“ 04 ” вересня 2025 р.

# ЗАВДАННЯ

## НА МАГІСТЕРСЬКУ РОБОТУ СТУДЕНТУ

**Кантору Матвію Володимировичу**

(прізвище, ім'я, по-батькові)

**1. Тема магістерської роботи** “**Моделі та методи контролю якості та автентичності мережевих даних**”

керівник проекту (роботи) Піх В.Я., к.т.н., доцент

затверджені наказом закладу вищої освіти від “ 05 ” листопада 2025 р. № 695/7

**2. Строк подання студентом проекту (роботи)** 15 грудня 2025 р.

**3. Вихідні дані до проекту (роботи)** Теоретичні концепції та формальні моделі побудови інформаційних технологій обробки мережевих даних

**4. Зміст розрахунково - пояснювальної записки(перелік питань, які потрібно розробити)**

1. Аналіз предметної області забезпечення автентичності та конфіденційності мережевих даних

2. Використання методів машинного навчання для контролю якості, автентичності даних

3. Дослідження моделей та представлення рішення контролю якості мережевих даних

4. Оцінка імплементації методів контролю якості та автентичності мережевих даних

**5. Перелік графічного матеріалу (з точним зазначенням обов'язкових креслень)**

1. Спрощена архітектура ML4Nets (рис. 1.1)

2. Конвеєр EMERGE (рис. 1.2)

3. Архітектура внутрішньомережевої безпеки на основі машинного навчання (рис. 1.3)

4. Архітектура автоенкодера для виявлення аномалій (рис. 1.4)

5. Архітектура сіамської мережі (рис. 1.5)

## 6. Консультанти розділів проекту (роботи)

Розділ	Консультант	Підпис, дата
Перевірка на плагіат	доц., к.т.н. Вовк Р.Б.	

7. Дата видачі завдання 04 вересня 2025 р.

Керівник \_\_\_\_\_

(підпис)

Завдання прийняв до виконання \_\_\_\_\_

(підпис)

## КАЛЕНДАРНИЙ ПЛАН

№ п/п	Назви етапів магістерської роботи	Строк виконання етапів роботи	Примітка
1	Підбір і вивчення літератури по темі магістерської роботи	15.09.2025	виконано
2	Аналіз предметної області забезпечення автентичності та конфіденційності мережевих даних	29.09.2025	виконано
3	Використання методів машинного навчання для контролю якості, автентичності даних	15.10.2025	виконано
4	Дослідження моделей та представлення рішення контролю якості мережевих даних	08.11.2025	виконано
5	Оцінка імплементації методів контролю якості та автентичності мережевих даних	15.11.2025	виконано
6	Затвердження пояснювальної записки роботи завідувачем кафедри	14.12.2025	виконано

Студент – магістр \_\_\_\_\_

(підпис)

Керівник роботи \_\_\_\_\_

(підпис)

## АНОТАЦІЯ

**Магістерська робота:** 75 с., 15 рис., 5 табл., 45 джерел.

**Тема:** Моделі та методи контролю якості та автентичності мережевих даних

**Мета магістерської роботи** – розроблення моделей та методів контролю якості й автентичності мережевих даних із використанням підходів машинного навчання.

**Об’єкт дослідження** – процеси контролю якості, достовірності та конфіденційності даних у мережевих інформаційних системах.

**Предмет дослідження** – моделі та методи машинного навчання, що забезпечують перевірку автентичності, оцінку якості та пояснюваність результатів аналізу мережевих даних.

### **Результати дослідження**

В роботі проведено комплексний аналіз проблем забезпечення якості, автентичності та конфіденційності мережевих даних, визначено обмеження існуючих підходів і сформульовано вимоги до нових методів.

### **Висновок**

Розроблено фреймворк контролю якості та автентичності даних, який інтегрує алгоритми машинного навчання, техніки пояснюваності та механізми захисту конфіденційності.

**МЕРЕЖЕВІ ДАНІ, АВТЕНТИЧНІСТЬ, КОНТРОЛЬ ЯКОСТІ, КОНФІДЕНЦІЙНІСТЬ, МАШИННЕ НАВЧАННЯ, ШТУЧНИЙ ІНТЕЛЕКТ, ФРЕЙМВОРК, БЕЗПЕКА ДАНИХ, ІНТЕРПРЕТОВАНИСТЬ МОДЕЛЕЙ**

## ABSTRACT

**Master Thesis:** 75 pp., 15 fig., 5 tab., 45 sources.

**Topic:** Models and methods of quality control and authenticity of network data

**The method of the master's thesis** is the development of models and methods of quality control and authenticity of network data using machine learning approaches.

**The object of the study** is the processes of quality control, reliability and confidentiality of data in network information systems.

**The subject of the study** is machine learning models and methods that ensure authenticity, quality control and explainability of network data analysis results.

### **Research results**

The work provides a comprehensive analysis of the problems of ensuring the quality, authenticity and confidentiality of network data, identifies the limitations of existing approaches and formulates requirements for new methods.

### **Conclusion**

A framework for quality control and authenticity of data has been developed, which integrates machine learning algorithms, explainability techniques and confidentiality protection mechanisms.

**NETWORK DATA, AUTHENTICITY, QUALITY CONTROL, CONFIDENTIALITY, MACHINE LEARNING, ARTIFICIAL INTELLIGENCE, FRAMEWORK, DATA SECURITY, MODEL INTERPRETABILITY**

## ЗМІСТ

ПЕРЕЛІК УМОВНИХ СКОРОЧЕНЬ .....	10
ВСТУП.....	11
РОЗДІЛ 1. АНАЛІЗ ПРЕДМЕТНОЇ ОБЛАСТІ ЗАБЕЗПЕЧЕННЯ АВТЕНТИЧНОСТІ ТА КОНФІДЕНЦІЙНОСТІ МЕРЕЖЕВИХ ДАНИХ .....	15
1.1. Опис пропонованого рішення для вирішення проблем довіри та конфіденційності даних.....	15
1.2. Дослідження проблем впровадження машинного навчання в мережах та роль пояснюваного штучного інтелекту .....	16
1.2.1. Проблематика довіри та безпеки в ML4Nets.....	17
1.2.2. Роль пояснюваного штучного інтелекту (XAI).....	17
1.2.2. Виклики специфічні для мережевого домену .....	18
1.2.3. Проблема конфіденційності і пропоноване рішення .....	19
1.3. Архітектура та компоненти фреймворку ML4Nets .....	19
1.3.1. Компоненти архітектури .....	20
1.3.2. Потік операцій .....	22
1.4. Виклики впровадження машинного навчання в мережевих середовищах .....	22
1.4.1. Доступність маркованих даних.....	23
1.4.2. Конфіденційність даних .....	26
1.4.3. Інтерпретованість моделей (Explainability) .....	26
1.5. Використання методів машинного навчання для контролю якості, автентичності та конфіденційності мережевих даних .....	27
1.5.1. Методи та алгоритми для контроль якості даних (Data Quality).....	28
1.5.2. Алгоритми забезпечення автентичності даних (Authenticity) .....	29
1.5.3. Методи та підходи збереження конфіденційності (Confidentiality) даних .....	31
Висновки до розділу .....	34

РОЗДІЛ 2. ДОСЛІДЖЕННЯ МОДЕЛЕЙ ТА ПРЕДСТАВЛЕННЯ РІШЕННЯ КОНТРОЛЮ ЯКОСТІ ТА АВТЕНТИЧНОСТІ МЕРЕЖЕВИХ ДАНИХ.....	35
2.1. Аналіз недоліків попередніх досліджень у сфері використання машинного навчання для захисту даних .....	35
2.1.1. Аналіз фреймворку EMERGE .....	35
2.1.2. Дослідження фреймворку ARISE .....	38
2.2. Проектування та реалізація фреймворку для контролю якості та автентичності мережеских даних .....	39
2.2.1. Архітектура фреймворку .....	40
2.2.2. Використання механізму EMERGE та інтерпретованість даних ....	42
2.3. Реалізація основних модулів фреймворку .....	43
2.3.1. Веб-сервіс .....	43
2.3.2. DEEP як REST API сервіс .....	46
2.3.3. Техніка гібридної пояснюваності .....	48
Висновки до розділу .....	50
 РОЗДІЛ 3. ОЦІНКА ІМПЛЕМЕНТАЦІЇ МЕТОДІВ КОНТРОЛЮ ЯКОСТІ ТА АВТЕНТИЧНОСТІ МЕРЕЖЕВИХ ДАНИХ З ВИКОРИСТАННЯМ МОДЕЛЕЙ МАШИНОГО НАВЧАННЯ.....	51
3.1. Оцінка процесу комбінування функцій маркування даних .....	51
3.1.1. Використані набори даних .....	51
3.1.2. Експерименти та налаштування гіперпараметрів .....	52
3.1.3. Результати CASE #1: комбінація функцій маркування на основі однієї характеристики (RTT) .....	52
3.1.4. Результати CASE #2: комбінація функцій маркування на основі двох характеристик.....	54
3.2. Оцінка техніки гібридної пояснюваності .....	57
3.2.1. Випадок використання 1 фреймворку ARISE.....	57
3.2.2. Випадок використання 2: фреймворк Trustee .....	58
3.2.3. Представлення результатів для випадку #1.....	60

3.2.4. Представлення результатів випадку #2 .....	62
Висновки до розділу .....	67
ВИСНОВКИ .....	68
ПЕРЕЛІК ВИКОРИСТАНОЇ ЛІТЕРАТУРИ.....	71

## ПЕРЕЛІК УМОВНИХ СКОРОЧЕНЬ

FL - Federated Learning - Федеративне навчання

WSGI - Web Server Gateway Interface - Інтерфейс шлюзу веб-сервера

XAI - Explainable Artificial Intelligence - пояснюваний штучний інтелект

IDS - Intrusion Detection System - система виявлення вторгнень

HIDS - Host-based Intrusion Detection System - хостова система

виявлення вторгнень

NIDS - Network-based Intrusion Detection System - мережева система

виявлення вторгнень

WS - Weak Supervision - слабке навчання

GDPR - General Data Protection Regulation - загальний регламент про

захист даних

PCA - Principal Component Analysis - метод головних компонент

GAE - Graph Autoencoder - графовий автоенкодер

DP - Differential Privacy - диференційна приватність

## ВСТУП

### **Актуальність теми.**

У сучасних умовах глобальної цифровізації та стрімкого зростання обсягів переданої інформації питання довіри до даних у мережевих системах набуває особливого значення. Якість, автентичність і конфіденційність мережевих даних визначають не лише ефективність функціонування інформаційних систем, а й рівень їхньої кіберстійкості. З поширенням технологій штучного інтелекту (ШІ) і машинного навчання (ML) виникає потреба у створенні надійних моделей, здатних забезпечити контроль достовірності даних, захист від фальсифікацій і маніпуляцій, а також прозорість процесів прийняття рішень.

Проблематика автентичності мережевих даних особливо актуальна для систем, що оперують критично важливою інформацією — у фінансовому секторі, енергетиці, телекомунікаціях, військових і державних інформаційних мережах. Недостовірність або порушення цілісності даних може призвести до помилкових рішень, фінансових втрат і зниження рівня безпеки. Традиційні методи криптографічного захисту не завжди здатні забезпечити комплексну перевірку даних у динамічних і розподілених середовищах, тому актуальним є використання інтелектуальних моделей, що здатні не лише виявляти відхилення, а й пояснювати їх природу.

Розвиток підходів Explainable Artificial Intelligence (XAI) дозволяє зробити процеси машинного навчання більш прозорими та підвищити рівень довіри користувачів до рішень, прийнятих алгоритмами. У цьому контексті поєднання методів контролю якості даних, машинного навчання та пояснюваного ШІ створює основу для побудови нових фреймворків забезпечення достовірності та конфіденційності мережевих систем.

Отже, дослідження моделей і методів контролю якості та автентичності мережевих даних є надзвичайно актуальним у контексті розвитку безпечних

цифрових інфраструктур, систем моніторингу мережевого трафіку, кіберзахисту та інтелектуальних мереж зв'язку.

Актуальність теми магістерської роботи зумовлена зростанням ризиків маніпуляцій, підробки та несанкціонованого доступу до мережевих даних у сучасному інформаційному просторі. Збільшення обсягів трафіку, поява нових типів кіберзагроз і складність верифікації джерел даних створюють потребу у впровадженні інтелектуальних засобів контролю їхньої достовірності.

Попри наявність великої кількості методів шифрування й автентифікації, вони не враховують поведінкові та структурні характеристики даних у мережах. Використання машинного навчання дозволяє виявляти приховані закономірності, оцінювати якість даних та прогнозувати можливі аномалії, однак ці моделі часто залишаються “чорними скриньками”, що ускладнює їхню інтерпретацію. Саме тому поєднання підходів ХАІ (пояснюваного штучного інтелекту) з методами контролю якості й автентичності відкриває нові перспективи у створенні прозорих і надійних систем довіри до даних.

Розробка комплексного фреймворку, який інтегрує методи машинного навчання, пояснюваного ШІ та механізми захисту конфіденційності, сприяє вирішенню актуальних проблем інформаційної безпеки й підвищує ефективність управління якістю даних у мережевих середовищах.

**Метою магістерської роботи** є розроблення моделей та методів контролю якості й автентичності мережевих даних із використанням підходів машинного навчання.

**Об'єктом дослідження** є процеси контролю якості, достовірності та конфіденційності даних у мережевих інформаційних системах.

**Предметом дослідження** є моделі та методи машинного навчання, що забезпечують перевірку автентичності, оцінку якості та пояснюваність результатів аналізу мережевих даних.

## **Завдання дослідження**

Для досягнення поставленої мети в роботі вирішено такі основні завдання:

1. Провести аналіз сучасних методів забезпечення якості, автентичності та конфіденційності мережевих даних.
2. Визначити обмеження та виклики впровадження технологій машинного навчання в мережевих середовищах.
3. Розробити архітектуру фреймворку контролю якості та автентичності даних.
4. Реалізувати веб-сервісну інфраструктуру для інтеграції моделей контролю якості, достовірності та конфіденційності.
5. Провести експериментальні дослідження ефективності запропонованих моделей і методів на реальних мережевих наборах даних.

## **Методи дослідження**

У процесі дослідження використано такі методи:

- методи машинного навчання (класифікація, кластеризація, ансамблеві моделі, нейронні мережі) для аналізу та оцінки якості даних;
- підходи Explainable AI (XAI), зокрема LIME, SHAP, Decision Path Analysis для пояснення рішень моделей;
- методи статистичного аналізу та обробки даних для оцінки точності моделей;
- методи проектування програмних систем для побудови архітектури фреймворку;
- експериментальні дослідження з використанням відкритих наборів мережевих даних для перевірки ефективності запропонованих рішень.

## **Наукова новизна отриманих результатів**

Удосконалено підхід до контролю якості та автентичності мережевих даних шляхом інтеграції моделей машинного навчання з техніками пояснюваного ШІ та розроблено фреймворк контролю довіри до даних, що

забезпечує баланс між точністю класифікації, збереженням конфіденційності та пояснюваністю результатів.

### **Практичне застосування результатів**

Результати роботи можуть бути використані:

- у системах моніторингу та аналізу мережевого трафіку для підвищення достовірності даних;
- у системах кіберзахисту та виявлення аномалій для автоматичної перевірки автентичності інформації;
- у корпоративних інформаційних системах для забезпечення контролю якості даних та зниження ризику помилкових рішень.

**Структура магістерської роботи.** Робота складається зі вступу, трьох розділів та висновків. Загальний обсяг роботи становить 75 сторінок, і містить 15 рисунків, 5 таблиць, список використаних джерел із 45 найменувань.

# РОЗДІЛ 1. АНАЛІЗ ПРЕДМЕТНОЇ ОБЛАСТІ ЗАБЕЗПЕЧЕННЯ АВТЕНТИЧНОСТІ ТА КОНФІДЕНЦІЙНОСТІ МЕРЕЖЕВИХ ДАНИХ

## 1.1. Опис пропонованого рішення для вирішення проблем довіри та конфіденційності даних

У цьому дослідженні представлено рішення (фреймворк), призначений для вирішення вищезазначених проблем, що сприяє великомасштабному анутованню даних.

Анклави відіграють ключову роль у забезпеченні безпеки та конфіденційності при роботі з чутливими даними та моделями машинного навчання:

### 1. Обчислювальний анклав (Trusted Execution Environment - TEE).

Це може бути захищена апаратна або програмна область (наприклад, TEE, як-от Intel SGX або AMD SEV), де код і дані можуть виконуватися, будучи захищеними від доступу або модифікації з боку операційної системи, гіпервізора або іншого програмного забезпечення.

Основна мета це досягнення внутрішньої довіри. Усередині такого анклаву проводиться навчання, виведення або пояснення моделі, гарантуючи, що чутливі дані та сама модель "чорний ящик" залишаються конфіденційними та цілісними, навіть якщо решта системи скомпрометована.

### 2. Інформаційний/організаційний анклав (Data/Organizational Silo).

Це логічні або організаційні одиниці, які володіють та управляють власними наборами даних та моделями. Наприклад, різні відділи, компанії-партнери, медичні заклади або фінансові установи.

Основна мета це співпраця з дотриманням конфіденційності. Взаємодія "між анклавами" означає співпрацю між цими окремими власниками даних (наприклад, для спільного навчання моделі за допомогою федеративного навчання), при цьому їхні локальні дані ніколи не покидають відповідний

анклаву, що забезпечує виконання вимог до конфіденційності (наприклад, GDPR або HIPAA).

Таким чином, "анклави" відносяться до двох взаємопов'язаних концепцій: апаратних/програмних середовищ для захисту обчислень і організаційних/логічних структур для захисту володіння даними.

Довіра в межах одного анклаву досягається шляхом забезпечення інтерпретованості (пояснюваності) моделей машинного навчання типу "чорний ящик". Це реалізується за допомогою гібридної техніки пояснюваності, яка інтегрує:

- Глобальні методи пояснюваності - надають загальне розуміння функціонування моделі.
- Локальні методи пояснюваності - фокусуються на поясненні конкретних прогнозів моделі.

Таке поєднання забезпечує прозорість і зрозумілість процесів прийняття рішень моделями МН, що є основою для внутрішньої довіри.

Крім того, запропоноване рішення каркас сприяє співпраці між різними обчислювальними анклавами (між-анклавна довіра), одночасно забезпечуючи суворе дотримання вимог до конфіденційності даних. Це є критично важливим для обміну знаннями та спільного навчання без необхідності розкриття сирих або чутливих даних. Реалізація цього аспекту може включати застосування методів, як-от федеративне навчання або диференційна конфіденційність, для захисту інформації під час обміну.

## **1.2. Дослідження проблем впровадження машинного навчання в мережах та роль пояснюваного штучного інтелекту**

Успішна практична імплементація рішень машинного навчання для мереж (ML4Nets) залежить від здатності операторів мереж інтегрувати їх у виробничі середовища. Цей процес висуває дві ключові вимоги: формування довіри до цих рішень та забезпечення засобів для оцінки їхньої безпеки.

### *1.2.1. Проблематика довіри та безпеки в ML4Nets*

Довіра до моделі ML оператором мережі визначається як комфортне делегування контролю цій моделі. Водночас оцінка безпеки рішень ML охоплює аналіз проблеми аварій (failures), що визначаються як небажана та шкідлива поведінка, спричинена недоліками в проектуванні реальних ML-рішень.

На жаль, непрозорість (характер "чорного ящика") багатьох моделей ML, що використовуються, унеможливує для сучасних мережевих операторів інтерпретацію рішень та прогнозів, які генеруються цими моделями. Ця відсутність пояснюваності породжує недовіру, ускладнює оцінку безпеки моделей і, як наслідок, пояснює стриманість у широкому впровадженні рішень ML4Nets на практиці.

### *1.2.2. Роль пояснюваного штучного інтелекту (XAI)*

Вирішення цих викликів зумовило розвиток пояснюваного штучного інтелекту (XAI) — дослідницького напрямку, сфокусованого на підвищенні зрозумілості навчальних моделей та їхніх механізмів прийняття рішень [8, 5, 7]. Техніки XAI класифікуються за двома основними категоріями:

#### 1. Техніки глобальної пояснюваності:

Принцип: Використовують наближення у формі пояснюваних моделей для надання загального, високоуровневого пояснення функціонування даної моделі "чорного ящика".

Компроміси: Їхня ефективність часто є результатом компромісу між складністю наближеної моделі, точністю її пояснень та обчислювальною витратністю генерації.

Користь: Дозволяють операторам міркувати про загальну надійність рішень моделі (тобто як і чому модель приймає рішення, а також коли вона працює чи не працює), сприяючи формуванню загальної довіри.

#### 2. Техніки локальної пояснюваності:

Принцип: Спрямовані на надання пояснень для індивідуальних вхідних екземплярів.

Методи: Використовують такі концепції, як оцінки важливості ознак, механізми уваги та пояснення на основі правил.

Користь: Є критично важливими для міркування про конкретні прогнози та оцінки безпеки ML4Nets у сценаріях крайніх випадків, дозволяючи операторам зрозуміти потенційні наслідки конкретних неправильних рішень. Їхнє масштабоване застосування вимагає усвідомлення обчислювальної складності на екземпляр.

### *1.2.2. Виклики специфічні для мережевого домену*

Незважаючи на успішне застосування ХАІ в інших доменах (наприклад, комп'ютерний зір, автономні транспортні засоби [1, 3]), його адаптивність та ефективність для вирішення практичних проблем безпеки та продуктивності мереж залишаються малодослідженими. Це значною мірою зумовлено домено-специфічними проблемами даних:

- Нестача даних: загальний дефіцит (особливо маркованих) мережевих даних.
- Висока динаміка: значний обсяг і висока швидкість надходження мережевих даних із виробничих мереж.
- Неоднорідність збору: разовий характер існуючих зусиль зі збору даних.
- Конфіденційність та безпека: критичні проблеми конфіденційності, пов'язані зі збором мережевих даних.

Крім того, мережеві дослідники та оператори стикаються з методологічною невизначеністю щодо застосування ХАІ для одночасного задоволення подвійних вимог: досягнення довіри (через глобальне розуміння поведінки) та оцінки безпеки (через конкретні, перевіряні, локальні пояснення). Задоволення обох вимог вимагає системних інновацій для балансування між ресурсомістким, але деталізованим застосуванням

локальних технік та ефективним, але менш точним використанням наближених моделей глобальної пояснюваності.

### *1.2.3. Проблема конфіденційності і пропонуване рішення*

Додатковий виклик породжується конфіденційністю даних, оскільки мережеві дані містять чутливу інформацію. Це накладає вимоги до конфіденційності, що ускладнює співпрацю між різними анклавними (тобто операційними мережами або організаційними одиницями).

Хоча існують індивідуальні зусилля (наприклад, обмін синтетичними даними на основі GAN [20] або локальна інтерпретованість для довіри [12]), комплексне рішення, що одночасно розв'язує всі вищезазначені проблеми, залишається відкритою науковою задачею.

Для вирішення цих проблем пропонується фреймворк як інтегрований конвеєр, що має на меті забезпечення довіри до моделей ML у межах анклавів та сприяння співпраці між ними, зберігаючи при цьому конфіденційність даних.

Фреймворк дозволяє анклавам обмінюватися метаданими та отримувати переваги від наборів даних, що існують поза їхніми власними сховищами, не розкриваючи сирих даних. Крім того, він забезпечує інтерпретованість рішень моделі "чорного ящика" за допомогою гібридної техніки пояснюваності, яка інтегрує:

- глобальну пояснюваність для розуміння загальної функціональності моделі.
- локальну пояснюваність для розуміння конкретних прогнозів.

## **1.3. Архітектура та компоненти фреймворку ML4Nets**

ML4Nets — це парадигма, що інтегрує методи машинного навчання (ML) у мережеві архітектури з метою автоматизації, оптимізації та підвищення безпеки мережевих операцій. На відміну від єдиної універсальної

архітектури, ML4Nets є модульним фреймворком, що об'єднує різноманітні архітектури ML, які динамічно розгортаються на основі конкретних мережових завдань. Основні компоненти фреймворку включають: площину даних, площину обробки (з моделями ML), площину прийняття рішень та механізм зворотного зв'язку.

### *1.3.1. Компоненти архітектури*

Площина даних (Data Plane) - на цьому рівні відбувається збір сирих мережових даних, які є основним джерелом інформації для моделей ML. Дані охоплюють дві основні категорії:

- Дані про пересилання (Forwarding Data) - включають трафік на рівні пакетів, записи потоків, логи подій та дані телеметрії. Зазвичай, ці дані мають часово-послідовну природу і можуть бути оброблені як часові ряди.

- Дані про топологію (Topological Data) - представляють мережові структури як графи, де вузли (nodes) — це пристрої, а ребра (edges) — зв'язки. Ці дані слугують вхідними для графових моделей.

Площина обробки (Processing Plane) - тут розміщуються моделі ML, які аналізують дані, зібрані з площини даних. Застосовуються різноманітні архітектури, зокрема:

#### 1. Глибокі нейронні мережі (Deep Neural Networks):

- Згорткові нейронні мережі (CNN): Ефективні для аналізу даних з сіткоподібною структурою, таких як представлення трафіку, або для виявлення закономірностей у заголовках пакетів.

- Рекурентні нейронні мережі (RNN): Придатні для обробки послідовних даних, як-от мережові журнали та потоки трафіку, що мають тимчасові залежності.

#### 2. Графові нейронні мережі (GNN):

Архітектури, призначені для навчання на графових даних. Використовуються для моделювання топології мережі та аналізу її властивостей.

3. Архітектури XAI (Explainable AI): Застосовуються для підвищення прозорості моделей ML. Наприклад, гібридні підходи можуть використовувати простіші моделі (як-от дерева рішень) для пояснення поведінки складних "чорних скриньок".

4. Полегшені архітектури (Lightweight Architectures): Моделі, оптимізовані для розгортання на пристроях з обмеженими обчислювальними ресурсами (наприклад, у середовищі IoT).

5. Площина прийняття рішень (Decision Plane) - це компонент інтерпретує висновки, зроблені моделями ML, і трансформує їх у конкретні дії або конфігурації для мережі. Приклади рішень включають:

- Адаптивна маршрутизація трафіку.
- Динамічна ідентифікація аномалій безпеки.
- Оптимізація ресурсів.

6. Механізм зворотного зв'язку (Feedback Mechanism) - це замкнутий контур є критично важливим для безперервного навчання та адаптації. Він забезпечує зворотний зв'язок від фактичної поведінки мережі до моделей ML, що дозволяє їм постійно вдосконалюватися та адаптуватися до змінних мережевих умов.

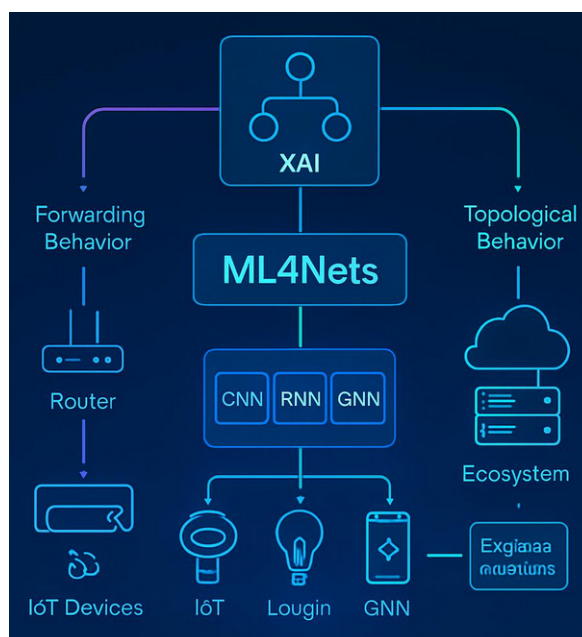


Рис. 1.1. Спрощена архітектура ML4Nets

### 1.3.2. Потік операцій

Процес в архітектурі ML4Nets відбувається за такою послідовністю:

- Збір: Постійно збираються дані з площини даних, що охоплюють як поведінку пересилання, так і топологію.
- Обробка: Зібрані дані передаються до відповідних архітектур ML у площині обробки.
- Висновок: Моделі ML генерують висновки, прогнози або ідентифікацію аномалій.
- Дія: Площина прийняття рішень перетворює висновки ML на конфігураційні зміни або оперативні команди.
- Вплив: Мережеві пристрої виконують отримані команди, змінюючи свою поведінку.
- Зворотний зв'язок: Дані про новий стан мережі знову збираються, і цикл починається знову.

Отже, архітектура ML4Nets є багат шаровим, динамічним фреймворком, що забезпечує автоматизоване управління та оптимізацію мереж за допомогою спеціалізованих моделей ML. Її модульність і замкнений цикл зворотного зв'язку дозволяють їй адаптуватися до мінливих умов і відповідати вимогам до продуктивності, безпеки та надійності в сучасних складних мережевих середовищах.

## 1.4. Виклики впровадження машинного навчання в мережевих середовищах

Мотивація цього дослідження базується на необхідності вирішення двох фундаментальних проблем, критичних для практичного застосування методів машинного навчання (МН) у сфері продуктивності та безпеки комп'ютерних мереж (ML4Nets):

- Забезпечення конфіденційності при взаємодії між різними операторами та ізольованими обчислювальними анклавами.

Підвищення інтерпретованості моделей МН типу "чорний ящик" для формування довіри мережевих операторів. Розглянемо виклики, специфічні для мережевого домену.

#### *1.4.1. Доступність маркованих даних*

Нестача якісних маркованих даних є суттєвою перешкодою в мережевій науці. Цей дефіцит виникає через відсутність стандартизованих, загально визнаних ознак для опису мережевих подій, а також через високу гетерогенність умов та методів збору даних.

Маркування таких даних вимагає експертної участі мережевих операторів, які є єдиними суб'єктами, що володіють достатнім доменним знанням. Оскільки обсяг міток, які може надати оператор, обмежений, перспективним рішенням є використання слабо керованого навчання (Weak Supervision). Цей підхід використовує невелику кількість точних експертних міток для генерації зашумлених міток для решти масиву даних. Інструментами для цього слугують методи програмування даних (Data Programming), де оператори створюють функції маркування для автоматичної категоризації даних.

Попередні дослідницькі роботи, зокрема системи NoMoNoise [7] та EMERGE [13], застосовували слабо кероване навчання для маркування мережевих даних.

NoMoNoise використовує принципи, аналогічні Snorkel [8], для зменшення шуму в інтернет-вимірах затримки та генерації міток шуму вимірювань.

EMERGE розвиває ці ідеї, застосовуючи генеративну модель для створення слабких міток, які потім використовуються для навчання дискримінаційної моделі (наприклад, LSTM) з метою оцінки якості даних.

В основі фреймворку EMERGE лежить інтеграція існуючих систем, організована у три головні модулі:

- Модуль великомасштабного створення міток даних - розширює ідеї, закладені у фреймворку NoMoNoise.

- Модуль забезпечення високої якості міток - спрямований на усунення проблем якості, наприклад, усунення упередженості (bias) у даних.

- Модуль сприяння низьковитратному маркуванню та спільному використанню - зосереджений на спільному (наприклад, громадському) маркуванні та обміні даними із дотриманням відповідних вимог до конфіденційності даних.

У поточній роботі представлено початковий прототип першого модуля та окреслено кроки проектування двох інших модулів, впровадження та оцінка яких залишаються майбутньою роботою. Припускаючи, що користувачі EMERGE (наприклад, дослідники) мають доступ до насичених мережевих даних, подальший виклад зосереджується на проблемі надання їм низьковитратної, масштабованої та високоякісної методології маркування даних.

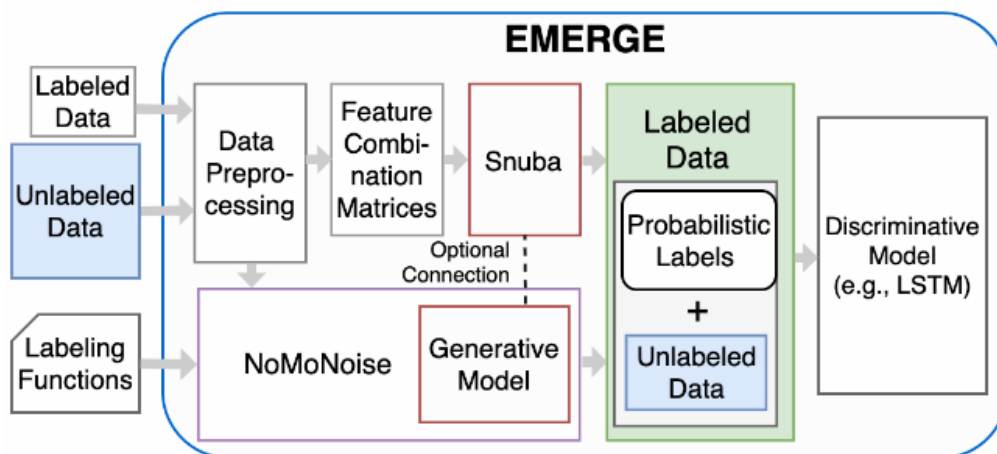


Рис. 1.2. Конвеєр EMERGE

Конвеєр EMERGE (показаний на рисунку 1.2) включає компонент Snuba, призначений для автоматичної генерації евристик на основі немаркованих та маркованих даних, а також вилучених ознак з обох наборів.

### Функціонування Компонента Snuba:

- Генерація моделей. Snuba генерує комбінації ознак і допасовує простий класифікатор (наприклад, логістичний регресор, метод найближчих сусідів або дерево рішень) до кожної комбінації ознак.

- Присвоєння міток. Кожен із цих класифікаторів потім присвоює імовірнісні мітки різним сегментам немаркованих даних.

Цей підхід дозволяє простим класифікаторам експлуатувати різні характеристики даних і відповідно призначати мітки, що знижує навантаження на користувачів, які інакше змушені були б самостійно розробляти функції маркування.

Ця автоматизована методологія також вирішує проблему масштабованості, властиву підходу NoMoNoise. Підхід NoMoNoise є фреймворком (каркасом), розробленим для зменшення шумів у вимірюваннях інтернет-затримки (latency). Його ключова характеристика полягає у використанні слабо керованого навчання (Weak Supervision) для маркування даних.

У NoMoNoise створення функцій маркування вимагає пошуку специфічних патернів для різних подій та типів мережевих даних. Зі зростанням обсягу та різноманітності даних цей пошук стає працемістким. Прості класифікатори в Snuba замінюють людське міркування в пошуку цих патернів, присвоюючи мітки на основі навчених закономірностей.

Компонент Snuba реалізовано без модифікацій і є незалежним від NoMoNoise, оскільки може генерувати власні імовірнісні мітки. Проте, користувачі можуть за бажанням під'єднати Snuba до генеративної моделі NoMoNoise для створення міток.

Оскільки Snuba вимагає, щоб його класифікатори навчалися на матрицях комбінацій ознак, для отримання необхідних ознак використовується tsfresh — інструмент для обчислення ознак часових рядів.

Маючи отримані імовірнісні мітки, користувачі EMERGE можуть навчати дискримінаційну модель-класифікатор, таку як LSTM, для виявлення подій, що становлять інтерес (наприклад, шум, аномалія тощо).

#### *1.4.2. Конфіденційність даних*

Чутливий характер мережевих даних ускладнює співпрацю між дослідницькими групами та різними організаційними анклавами. Хоча обмін навченими моделями МН (замість сирих даних) є одним із способів співпраці, він не вирішує проблему повністю.

Вразливість моделей МН до різних атак (наприклад, атаки відновлення даних або атаки членства) дозволяє зловмисникам отримувати інформацію про тренувальні дані. Таким чином, обмін моделями не гарантує захист конфіденційності. Це створює значну перешкоду для спільної дослідницької діяльності у сфері ML4Nets, особливо коли відсутній механізм, що дозволяє дослідникам, які не володіють даними, отримувати користь від існуючих наборів даних або прогнозів моделі, при цьому дотримуючись вимог власників щодо конфіденційності.

#### *1.4.3. Інтерпретованість моделей (Explainability)*

Мережеві оператори виявляють недовіру до моделей МН, що використовуються, через їхній характер "чорного ящика". Оскільки рішення, прийняті моделями, можуть мати критичні наслідки для безпеки та продуктивності мережі, інтерпретованість є необхідною умовою для розуміння та підтримки цих рішень. Інтерпретованість дозволяє операторам приймати обґрунтовані рішення на основі вихідних даних моделі.

Інтерпретованість класифікується за двома основними типами:

- Локальна інтерпретованість (Local Interpretability) - надає точне та деталізоване уявлення про те, чому модель зробила певний прогноз для конкретного вхідного екземпляра. Хоча цей підхід є високоточним для окремих випадків, він не надає узагальненої поведінки моделі.

- Глобальна інтерпретованість (Global Interpretability) - спрямована на пояснення загальної поведінки моделі на всьому наборі даних, включаючи аналіз загальних шаблонів, тенденцій та вагомості ознак. Проте, через необхідність узагальнення складної моделі у зрозумілий формат, цей процес може бути схильним до помилок і не завжди точно відображати фактичну поведінку моделі.

Для подолання цих обмежень пропонується застосування гібридної інтерпретованості, яка поєднує глобальний контекст із застосуванням локальної інтерпретованості для підвищення точності та надійності пояснень.

### 1.5. Використання методів машинного навчання для контролю якості, автентичності та конфіденційності мережевих даних

Розглянемо, як машинне навчання (МН) допомагає захищати мережеві дані. Існують різні методи для контролю якості, автентичності та конфіденційності.

На рисунку 1.3 подана загальна архітектура системи безпеки на основі машинного навчання.

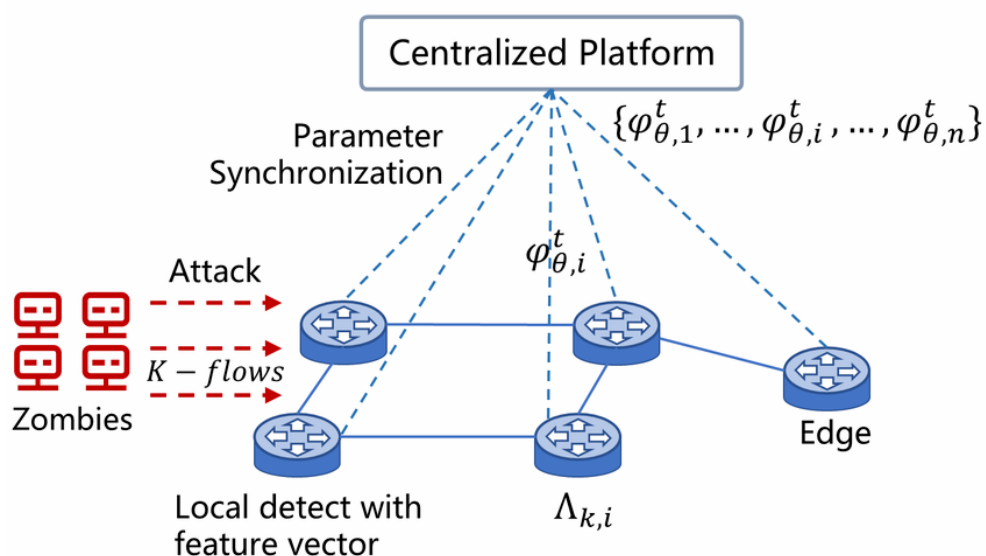


Рис. 1.3. Архітектура внутрішньомережевої безпеки на основі машинного навчання

Як показано на рисунку 1.3, потоки атаки, змішані з нормальними потоками, випадково передаються через розподілені комутатори.

1. Обмін інформацією та ітерація: протягом обробного періоду сусідні комутатори обмінюються локальною інформацією та ітеративно оновлюють локальні параметри.

2. Централізована платформа: централізована платформа агрегує поточні результати обробки, маркує аномальний трафік і обчислює глобальний параметр для наступної фази обробки.

3. Ефективність: спільна робота централізованої платформи та комутаторів дозволяє ефективно зменшити комунікаційні витрати (overhead) між самими комутаторами.

#### *1.5.1. Методи та алгоритми для контроль якості даних (Data Quality)*

Головна мета це виявляти помилки, пошкодження та аномалії в мережевому трафіку. Моделі МН навчаються на "нормальних" даних, щоб ідентифікувати будь-які відхилення.

Основні методи та алгоритми:

Виявлення аномалій (Anomaly Detection) - це основний підхід для контролю якості.

- Автоенкодери (Autoencoders) - це нейронні мережі, які вчаться стискати дані (кодувати) і потім відновлювати їх до початкового вигляду (декодувати). Якщо відновлені дані сильно відрізняються від оригіналу, це вказує на аномалію або помилку. Вони чудово підходять для виявлення нетипових мережевих пакетів.

- Ізоляційний ліс (Isolation Forest) - ефективний алгоритм, який "ізолює" аномалії, а не будує модель нормальних даних. Він добре працює на великих обсягах даних.

- Кластеризація (напр., DBSCAN) - алгоритми, які групують схожі дані разом. Точки даних, що не належать до жодного кластера, вважаються аномаліями (шумом).

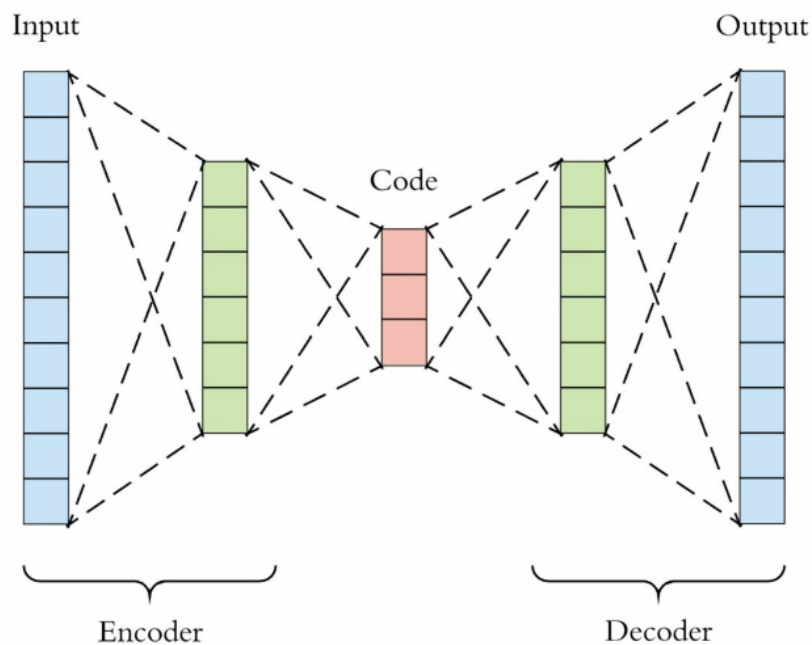


Рис. 1.4. Архітектура автоенкодера для виявлення аномалій

### 1.5.2. Алгоритми забезпечення автентичності даних (Authenticity)

Основним завданням це переконатися, що дані надходять із довіреного джерела і не були змінені в дорозі. МН допомагає виявляти спроби підробки (spoofing), фішингу та інших атак.

Основні методи та алгоритми:

- Системи виявлення вторгнень (Intrusion Detection Systems) - моделі МН аналізують мережевий трафік у реальному часі, щоб класифікувати його як "нормальний" або "підозрілий". Системи виявлення вторгнень — це програмні або апаратні засоби безпеки, призначені для моніторингу мережі або системи на предмет зловмисної активності, порушень політики або ознак вторгнення. Їхня основна мета — ідентифікувати аномалії та генерувати попередження, щоб адміністратори могли вжити заходів.

IDS на основі машинного навчання вчаться на даних. Це дозволяє їм виявляти не тільки відомі, але й абсолютно нові, раніше не бачені атаки.

Є два основні способи, якими МН цього досягає:

1. Виявлення аномалій (Anomaly-based detection) - система вчиться, як виглядає "нормальна" поведінка у вашій мережі. Вона аналізує тисячі

параметрів: хто, куди, коли і як часто відправляє дані. Будь-яке значне відхилення від цієї норми (аномалія) вважається підозрілим.

- Зловмисний аналіз (Misuse-based detection) - це більш схоже на традиційний підхід, але значно потужніше. Замість ручного написання правил, модель МН "годує" величезною кількістю прикладів як шкідливого, так і нормального трафіку. Модель сама вчиться знаходити складні закономірності, що відрізняють атаку від легітимної активності.

Системи виявлення вторгнень класифікують за місцем моніторингу та методом виявлення (таблиця 1.1.)

Таблиця 1.1.

Типи системи виявлення вторгнень за місцем моніторингу

Тип	Назва	Опис
Мережева СВВ	NIDS (Network-based)	Моніторить мережевий трафік у реальному часі. Аналізує заголовки та вміст пакетів, що проходять через мережевий сегмент (наприклад, між маршрутизатором і комутатором).
Хостова СВВ	HIDS (Host-based)	Моніторить конкретний хост (сервер, робочу станцію). Аналізує системні журнали, файли, реєстр, активність користувачів та виклики до системи.

- Класифікаційні моделі (SVM, Random Forest, нейронні мережі) - навчаються на датасетах, що містять як звичайний, так і шкідливий трафік (наприклад, від DDoS-атак, сканування портів).

- Сіамські мережі (Siamese Networks) - це спеціальні нейронні мережі з двома однаковими "гілками". Вони використовуються для порівняння двох об'єктів — наприклад, для перевірки, чи відповідає поточна поведінка користувача його звичайному "цифровому відбитку". Це допомагає виявити захоплення облікового запису. Мережа складається з двох паралельних гілок, які мають абсолютно однакову архітектуру та, що найважливіше, спільні ваги (shared weights). Це означає, що обидві підмережі навчаються однаково і виконують ідентичне перетворення вхідних даних. Мережа не вивчає класи

об'єктів, а вчиться створювати такий латентний простір, у якому схожі об'єкти групуються разом, а відмінні — розходяться.

Навчання мережі відбувається за допомогою спеціальної функції втрат, найчастіше — контрастивної (contrastive loss). Її мета — мінімізувати відстань між векторами схожих об'єктів (наприклад, два зображення однієї людини) і максимізувати відстань для векторів несхожих об'єктів.

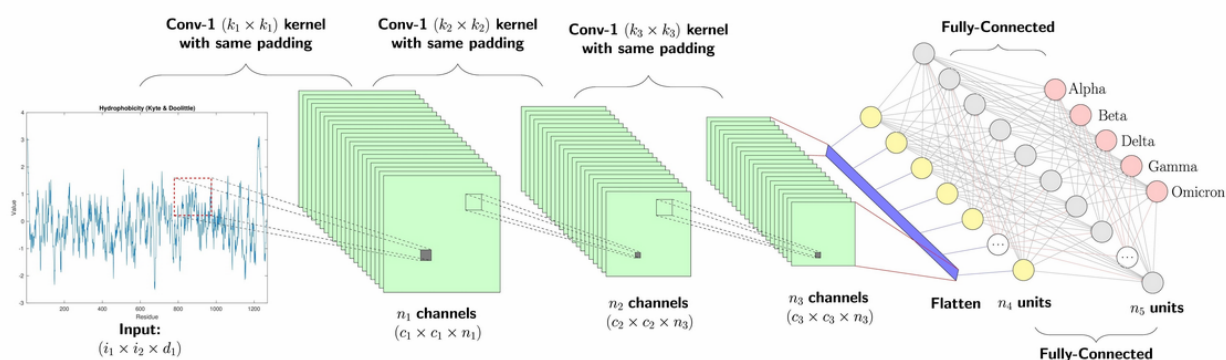


Рис. 1.5. Архітектура сіамської мережі

### 1.5.3. Методи та підходи збереження конфіденційності (Confidentiality) даних

Мета — захистити дані від несанкціонованого доступу. Хоча шифрування є основним інструментом, МН пропонує інноваційні підходи до роботи з конфіденційними даними.

Основні методи та підходи:

- Федеративне навчання (Federated Learning) - це революційний підхід, за якого модель МН тренується безпосередньо на пристроях користувачів (наприклад, на телефонах), не надсилаючи їхні персональні дані на центральний сервер. На сервер передаються лише оновлення моделі, що зберігає конфіденційність вихідних даних.

- Диференційна приватність (Differential Privacy) - техніка, яка додає невелику кількість "шуму" до даних або результатів запитів. Це дозволяє

проводити аналіз великих масивів даних, не розкриваючи інформацію про окремих користувачів.

Класифікація даних - моделі МН можуть автоматично сканувати дані та класифікувати їх за рівнем чутливості (наприклад, "публічні", "внутрішні", "конфіденційні"), щоб до них автоматично застосовувалися відповідні політики безпеки.

Федеративне навчання (Federated Learning, FL) є відносно новою технологією, яка привернула значну увагу дослідників щодо її потенціалу та практичного застосування. Основною метою FL є відповідь на ключове питання: чи можна навчити модель, не вимагаючи передачі даних до централізованого сховища?

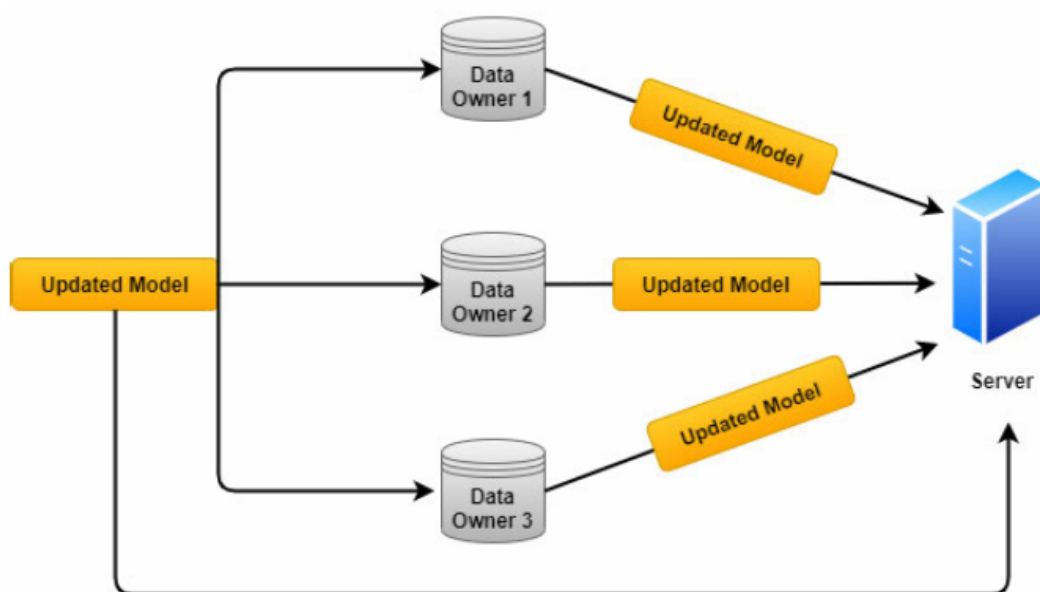


Рис. 1.6. Архітектура федеративного навчання

У рамках архітектури FL акцент зміщується на співпрацю між пристроями, що не завжди досягається стандартними алгоритмами машинного навчання (ML). Крім того, FL дозволяє задіяним алгоритмам набувати досвіду, чого не завжди можна гарантувати за допомогою традиційних методів ML.

FL вже знайшло застосування у широкому спектрі доменів, включаючи медицину, Інтернет речей (IoT), транспорт, оборонну сферу та мобільні додатки. Ця універсальність робить FL вкрай надійним, що підтверджується численними проведеними експериментами.

Незважаючи на багатообіцяючий потенціал, FL все ще недостатньо вивчене щодо деяких своїх технічних компонентів, таких як платформи, апаратне та програмне забезпечення, а також щодо аспектів конфіденційності даних та доступу до даних.

Суворі регуляторні вимоги щодо конфіденційності даних часто роблять непрактичним збір і спільне використання даних споживачів у централізованому місці. Це створює виклики для традиційних алгоритмів ML, оскільки вони зазвичай вимагають великих обсягів навчальних прикладів для ефективного навчання.

Обмеження традиційних алгоритмів ML пов'язані з їхнім процесом навчання, який, як правило, передбачає використання основного сервера для зберігання даних та обробки моделей. Існують два типові способи використання навчених моделей: побудова конвеєра для передачі даних через сервер, або перенесення моделей ML на пристрої, що взаємодіють із середовищем. Однак обидва ці підходи не є оптимальними через їхню нездатність до швидкої адаптації моделей.

У FL моделі навчаються на рівні пристроїв (device level). Моделі передаються до джерел даних або пристроїв для навчання та прогнозування. Після локального навчання оновлення моделей (models' updates) надсилаються назад на основний сервер для агрегації. Потім консолідована (узагальнена) модель повертається на пристрої з використанням концепцій розподілених обчислень. Це дозволяє відстежувати та перерозподіляти кожен з моделей на різних пристроях.

Такий підхід FL є дуже вигідним для використання недорогих моделей машинного навчання на периферійних пристроях, таких як мобільні телефони та сенсори. Загальна архітектура FL представлена на рисунку 1.6.

Ці методи часто комбiнуються для створення комплексних систем безпеки, здатних адаптуватися до нових загроз у режимі реального часу.

### **Висновки до розділу**

У першому розділі проведено ґрунтовний аналіз предметної області забезпечення якості, автентичності та конфіденційності мережевих даних. Визначено основні проблеми впровадження методів машинного навчання у мережевих середовищах, серед яких нестача маркованих даних, потреба у збереженні приватності та складність інтерпретації моделей. Розглянуто архітектуру фреймворків типу ML4Nets, що поєднують методи аналізу трафіку, перевірки довіри та пояснюваного ШІ. Особливу увагу приділено ролі Explainable AI (XAI) у підвищенні прозорості та надійності результатів автоматизованих систем. У підсумку визначено вимоги до створення інтегрованого фреймворку, здатного одночасно контролювати якість, забезпечувати автентичність і конфіденційність мережевих даних.

## РОЗДІЛ 2. ДОСЛІДЖЕННЯ МОДЕЛЕЙ ТА ПРЕДСТАВЛЕННЯ РІШЕННЯ КОНТРОЛЮ ЯКОСТІ ТА АВТЕНТИЧНОСТІ МЕРЕЖЕВИХ ДАНИХ

### 2.1. Аналіз недоліків попередніх досліджень у сфері використання машинного навчання для захисту даних

Попередні дослідницькі роботи, спрямовані на демократизацію використання машинного навчання (МН) для аналізу мережеских даних, виявили два основні системні недоліки:

- Забезпечення співпраці з дотриманням конфіденційності (Inter-Enclave Privacy-Preserving Collaboration).
- Інтерпретованість (Пояснюваність) моделей МН (Model Interpretability).

Хоча були зроблені спроби [12 - 13] демократизувати використання МН для маркування мережеских даних шляхом надання низьковитратних і високоякісних методів через програмування даних та багатозадачне навчання, ці зусилля є частковими.

#### 2.1.1. Аналіз фреймворку EMERGE

Фреймворк EMERGE зосереджений на вирішенні проблеми великомасштабного маркування даних за низької вартості та високої якості, спираючись на концепцію слабо керованих методів програмування даних.

EMERGE припускає наявність мережеских даних і використовує, наприклад, дані трасування маршруту з проекту CAIDA Ark (охоплюють один день) для програмного маркування.

Процес маркування є багатоетапним:

- Дані аналізуються для встановлення порогів, що використовуються для диференціації між зашумленими та незашумленими даними.
- Дані поділяються, і невелика експертна частина маркується вручну.

- Для оцінки якості ймовірнісних навчальних міток, згенерованих EMERGE, навчаються дискримінаційні моделі.

Для забезпечення співпраці EMERGE дозволяє обмінюватися метаданими (наприклад, функціями маркування) замість обміну сирими даними або навченими моделями, тим самим зберігаючи конфіденційність.

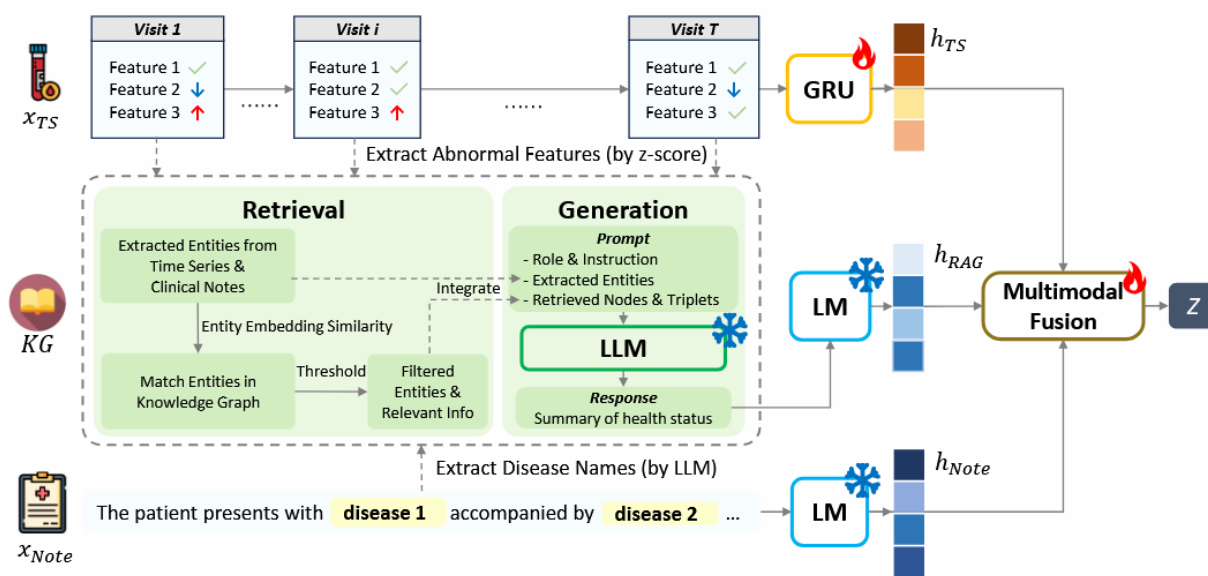


Рис. 2.1. Загальна архітектура фреймворку EMERGE

На рисунку 2.1 наведено загальну архітектуру фреймворку EMERGE призначеного для аналізу багатомодальних електронних медичних записів (EHR). Фреймворк складається з трьох основних модулів, які забезпечують перетворення сирих даних на збагачені та інтегровані представлення для подальших завдань.

### 1. Багатомодальне вилучення даних.

Цей модуль фокусується на трансформації сирих, зрозумілих людині вхідних даних ( $x$ ) у глибокі семантичні ембединги ( $h$ ) для комплексного аналізу.

Для роботи з даними часових рядів ( $x_{TS}$ ), що фіксують динамічні медичні показники, використовується мережа Gated Recurrent Unit (GRU) як енкодер.

GRU є високоефективним варіантом рекурентних нейронних мереж, здатним фіксувати часові залежності та кодувати цю послідовно пов'язану інформацію. Результуюче представлення часових рядів ( $h_{TS}$ ) формально визначається як:

$$h_{TS} = GRU(x_{TS})$$

Для текстових нотаток ( $x_{Note}$ ) застосовується мовна модель, попередньо навчена в домені (TextEncoder), для отримання їхніх текстових ембедингів ( $h_{Note}$ ).

Формальне представлення нотаток:

$$h_{Note} = TextEncoder(x_{Note})$$

## 2. Конвеєр збагачення на основі RAG (RAG-Driven Enhancement Pipeline)

Цей модуль (зображений у пунктирній рамці) використовує підхід Retrieval-Augmented Generation (RAG) для збагачення даних шляхом інтеграції зовнішніх знань. Конвеєр витягує медичні сутності та пов'язує знання у встановлених графах знань (KG) як для часових рядів, так і для нотаток. Отримані триплети сутностей-відносин (разом з їхніми визначеннями та описами, а також відносинами) інтегруються у промпт. Цей промпт слугує інструкцією для великої мовної моделі (LLM), щоб згенерувати резюме збагачене контекстом знань.

## 3. Мультимодальна мережа (Multimodal Fusion Network)

Цей кінцевий модуль адаптивно об'єднує згенеровані LLM-резюме та представлення з декількох модальностей (часові ряди та текст) для подальших завдань обробки даних (наприклад, прогнозування результатів).

Недоліком EMERGE є те, що він ефективно вирішує проблему конфіденційної співпраці (шляхом обміну метаданими), проте не пропонує жодного механізму для міркування про рішення, які приймає модель МН, тобто не забезпечує інтерпретованість.

### 2.1.2. Дослідження фреймворку ARISE

ARISE є багатозадачним фреймворком слабкого керування, що також використовує програмування даних для великомасштабного маркування мережевих даних.

ARISE застосовує багатозадачне навчання (Multitask Learning) та метанавчання (Meta-Learning) для підвищення ефективності обміну знаннями між різними завданнями та скорочення загального часу навчання.

Архітектурні компоненти:

- Інтерфейс для операторів - дозволяє мережевим операторам перетворювати свої доменні знання на "програмне представлення" у формі функцій маркування.

- Генерація міток: немарковані дані подаються на вхід і разом із функціями маркування використовуються для генерації слабких міток.

- Багатозадачна класифікація - слабо марковані дані використовуються для навчання моделі класифікації, яка виконується в кількох підзавданнях.

ARISE надає механізми локальної інтерпретованості для розуміння рішень моделі на рівні окремих даних.

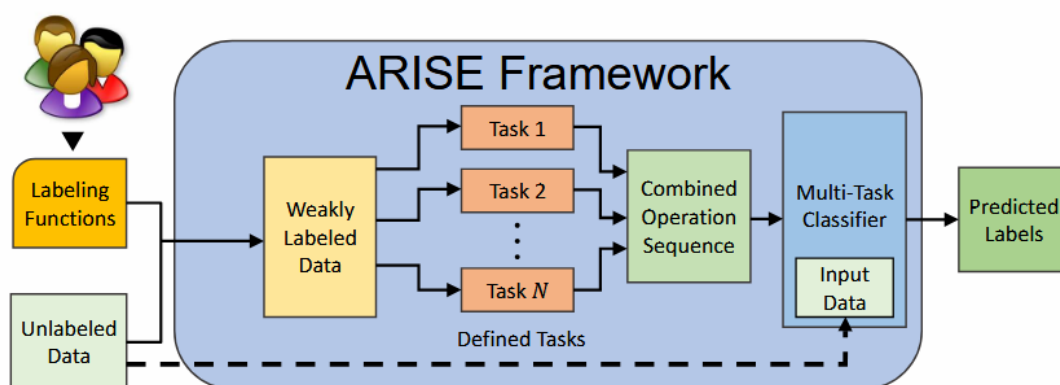


Рис. 2.2. Архітектура фреймворку ARISE

На рисунку 2.2 представлені основні компоненти системи, які детально описуються нижче:

## 1. Користувацький інтерфейс (User Interface)

Успіх ARISE критично залежить від його здатності захоплювати та трансформувати доменну експертизу мережевих операторів у конкретні, відчутні евристики, виражені як літеральні рядки коду. Приклад такої евристики: якщо вимірювання затримки  $x$  перевищує середнє значення  $\mu + 3\sigma$  для заданого часового ряду, то  $x$  є зашумленим вимірюванням.

З цією метою було розроблено простий інтерфейс, який дозволяє операторам та дослідникам конвертувати їхні доменні знання щодо мережевих подій у програмні представлення (наприклад, використовуючи Python), відомі як функції маркування (labeling functions).

## 2. Немарковані дані (Unlabeled Data)

Вхідні дані для ARISE можуть бути зібрані за допомогою традиційних технік вимірювання (наприклад, scamper) або базуватися на представленнях, специфічних для певного вендора (наприклад, NetFlow).

Єдине припущення, закладене у розробку ARISE, полягає у можливості вилучення вимірювань часових рядів із набору даних. Часові ряди зазвичай використовуються для забезпечення операторів часовим оглядом їхніх даних, що дозволяє оцінювати мережеві проблеми у великих масштабах (наприклад, перевіряти ефективність рішення для пом'якшення проблеми шляхом аналізу патернів трафіку до та після події перевантаження).

Головною проблемою ARISE є те, що він не вирішує жодних питань, пов'язаних зі співпрацею та підтриманням конфіденційності між різними учасниками (співробітниками), що залишає ключовий виклик ML4Nets невирішеним.

## **2.2. Проектування та реалізація фреймворку для контролю якості та автентичності мережевих даних**

Основною метою цього дослідження є проектування та впровадження фреймворку, який пропонує новітній підхід для вирішення ключових

проблем у сфері машинного навчання для мереж (ML4Nets). Пропонований фреймворк інтегрує наступні функціональні можливості:

- Масштабоване та економічно ефективно маркування мережових даних, використовуючи принципи попередніх досліджень.
- Забезпечення конфіденційної співпраці між мережевими операторами та дослідниками шляхом обміну метаданими, а не чутливими даними чи моделями МН.
- Надання механізму для інтерпретації рішень навчених моделей МН, включаючи їхнє порівняння з експертними знаннями.

### 2.2.1. Архітектура фреймворку

Архітектура пропованого фреймворку, зображена на рисунку 2.3, складається з декількох взаємопов'язаних модулів.

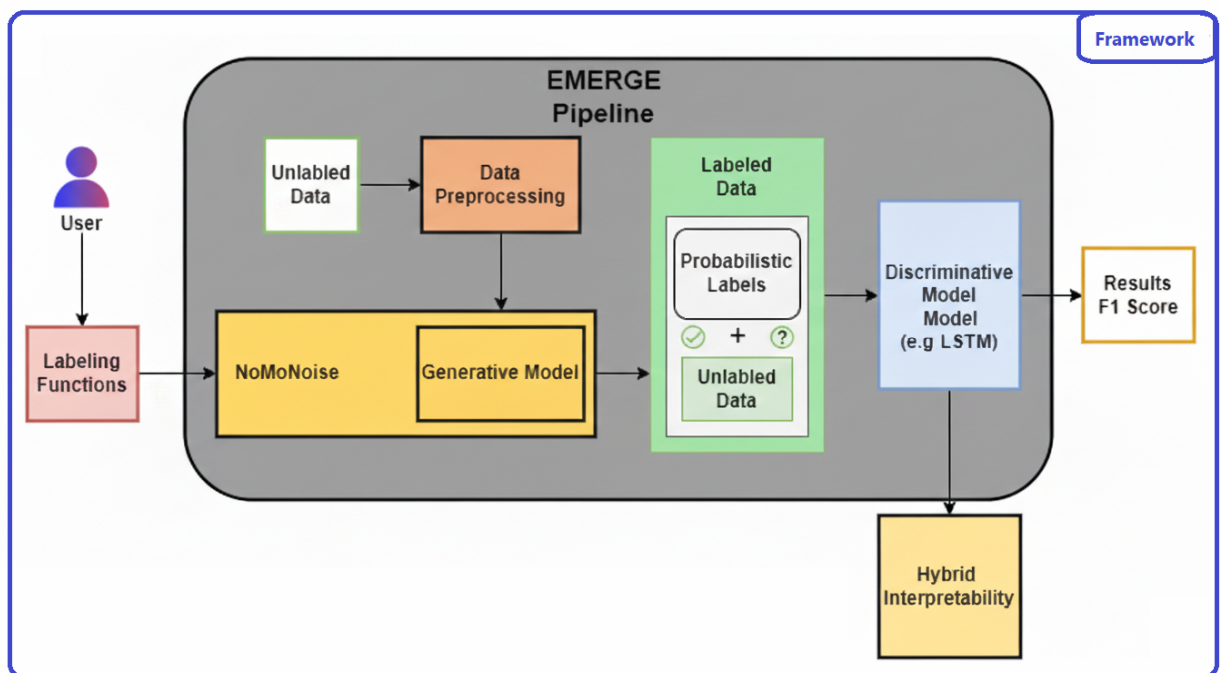


Рис. 2.3. Архітектура пропованого рішення

Ключовий внесок фреймворку полягає у створенні механізму, що дозволяє дослідникам та операторам трансформувати їхні галузеві знання у програмний формат та сприяти співпраці між різними групами.

Для досягнення цієї мети розроблено інтерфейс, доступний для різних груп, який дозволяє їм надавати програмні представлення своїх знань, зокрема функції маркування (labeling functions), реалізовані на мові Python.

Цей інтерфейс відкриває доступ до конвеєра EMERGE, який використовує надані функції маркування для генерації міток для відповідного набору даних.

Інтерфейс також надає користувачам можливість:

- Конфігурувати гіперпараметри генеративних та дискримінаційних моделей.
- Вказувати набори даних (наприклад, CAIDA та RIPE) для навчання.
- Переглядати подання та функції маркування, надані іншими користувачами. Це дозволяє різним групам спільно вирішувати, чи навчати модель, використовуючи комбінацію локальних та зовнішніх функцій маркування, та оцінювати їхню відносну ефективність.

Функції маркування (подано в лістингу 2.1) створюються користувачами для надання слабких міток немаркованим мережевим наборам даних, які використовуються для тренування моделей МН.

### Лістинг 2.1. Приклад функції маркування

```
def LF_Mean(c):  
    val = c.number1.get_attrib_tokens()  
    if float(val[0]) <= mean_threshold:  
        return 1  
    else:  
        return 0
```

Цей фрагмент коду є прикладом функції маркування (Labeling Function, LF), яка використовується для автоматичного присвоєння міток немаркованим даним.

Назва функції: LF\_Mean(c) (де LF позначає функцію маркування).

Призначення: Класифікація точки даних (c) на основі порівняння її значення з визначеним порогом середнього значення (mean\_threshold).

Функція отримує значення ознаки (`number1`) з об'єкта даних (`c`), перетворює його на число з плаваючою комою (`float(val[0])`) та виконує бінарну класифікацію:

- return 1: Якщо значення  $\leq$  порогового середнього (наприклад, мітка "хороші дані" або "не зашумлені дані").

- return 0: Якщо значення  $>$  порогового середнього (наприклад, мітка "погані дані" або "зашумлені/аномальні дані").

Ця функція втілює доменну евристику в програмний код, що є центральним елементом підходу програмування даних.

Функція може реалізовувати емпіричне правило, наприклад, класифікувати дані як "хороші" (1) або "погані" (0), якщо значення RTT (Round-Trip Time) менше або дорівнює середньому.

Хоча функції можуть бути налаштовані для конкретного анклаву, вони також можуть включати статистичні або загальні пороги, що підвищує їхню масштабованість та дозволяє застосовувати їх до різних мережевих середовищ та наборів даних.

### 2.2.2. Використання механізму EMERGE та інтерпретованість даних

Конвеєр EMERGE — це механізм слабо керованого навчання, що базується на фреймворку NoMoNoise. Він забезпечує створення слабких міток на даних, використовуючи надані функції маркування з користувацького інтерфейсу.

За допомогою згенерованих слабких міток (із використанням, наприклад, бібліотеки Snorkel) створюється генеративна модель. Ця модель присвоює ймовірнісні значення слабким міткам, вказуючи рівень впевненості у призначенні певної мітки (наприклад, значенню RTT).

Згенеровані ймовірнісні мітки потім використовуються для навчання дискримінаційної моделі (наприклад, LSTM), яка виконує основне завдання класифікації.

Ще одним важливим внеском пропонованого рішення є його здатність забезпечувати інтерпретованість рішень моделі МН як у глобальному, так і в локальному масштабах. Це реалізується шляхом надання гібридного фреймворку інтерпретованості, що дозволяє:

- Глобальний аналіз - розуміння загальної поведінки навченої моделі.
- Локальний Аналіз - перегляд детальних пояснень для кожної точки даних, що дозволяє операторам міркувати про правильність або неправильність конкретної класифікації.

### **2.3. Реалізація основних модулів фреймворку**

У цьому розділі детально описується реалізація фреймворку, який забезпечує конфіденційну співпрацю між різними суб'єктами (мережевими операторами та дослідниками). Архітектура, представлена на рисунку 2.3, побудована з використанням трьох ключових контейнерів Docker:

1. Контейнер для веб-додатку Flask.
2. Контейнер для REST API DEEPaaS (DEEP as a Service) для експозиції моделі.
3. Контейнер для зворотного проксі (наприклад, Nginx).

#### *2.3.1. Веб-сервіс*

Веб-сервіс реалізовано за допомогою мікрофреймворку Flask, що створює кінцеву точку доступу для співпраці (приклад архітектури розгортання подано на рис. 2.4).

Додаток Flask надає панель управління, де користувачі можуть створити обліковий запис. Облікові дані (ім'я користувача та пароль) зберігаються у базі даних (наприклад, MongoDB).

Після авторизації користувачам надається можливість:

- Обмінюватися функціями маркування.

- Обирати набори даних (наприклад, CAIDA та RIPE) для використання генеративною моделлю з метою створення міток для дискримінаційної моделі.

- Вказувати гіперпараметри для генеративних та дискримінаційних моделей.

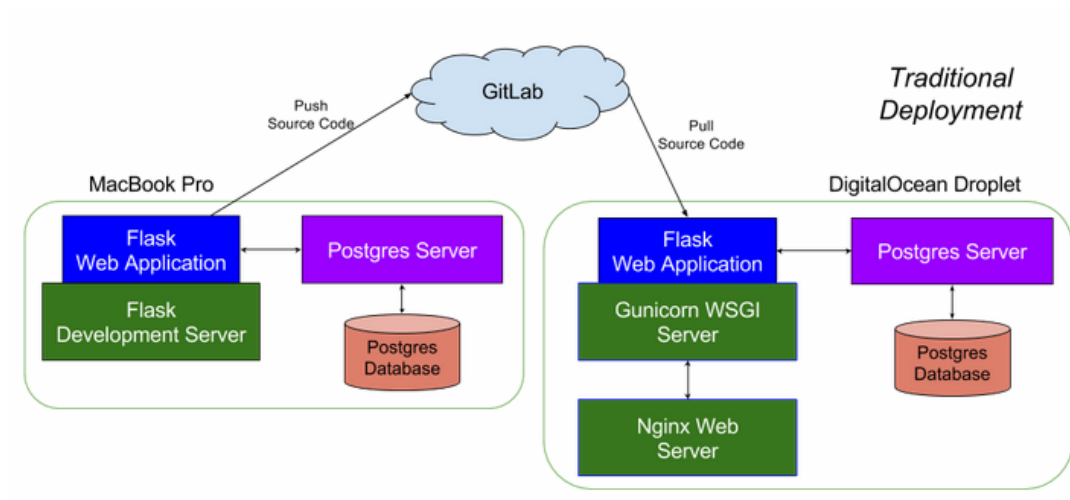


Рис. 2.4. Приклад архітектури розгортання веб-додатку з використанням системи контролю версій та Flask

На рисунку 2.4 показано традиційну архітектуру розгортання (Traditional Deployment) веб-додатку, розділену на середовище розробки та виробниче середовище, з використанням Git як системи контролю версій.

### 1. Середовище Розробки (Development Environment)

Ліва частина діаграми представляє локальне середовище, яке зазвичай використовується розробником для створення та тестування коду.

Веб-додаток: Flask Web Application — ядро веб-додатка, написаного на мікрофреймворку Flask.

Сервер розробки: Flask Development Server — вбудований у Flask сервер, який використовується для локального запуску та тестування додатку.

База даних: Postgres Server та Postgres Database — локальний екземпляр бази даних PostgreSQL, до якої додаток звертається під час розробки.

Контроль версій: Розробник надсилає (Push Source Code) вихідний код до центрального репозиторію GitLab.

## 2. Центральний рРепозиторій (Source Control)

GitLab: Хмарна або локальна платформа, що зберігає вихідний код додатку. Вона забезпечує спільну роботу та контроль версій.

## 3. Виробниче середовище (Production Environment)

Права частина діаграми зображує серверну інфраструктуру, де додаток працює для кінцевих користувачів.

Хостинг: DigitalOcean Droplet — віртуальна приватна машина (VPS) або хмарний інстанс, який використовується для розміщення додатку.

Розгортання коду: код завантажується (Pull Source Code) з GitLab на сервер.

## 3. Стек веб-сервера:

Nginx Web Server: виконує роль зворотного проксі (reverse proxy) та веб-сервера, обробляючи вхідні запити від користувачів, обслуговуючи статичні файли та перенаправляючи динамічні запити.

Gunicorn WSGI Server: виступає як WSGI-сервер (Web Server Gateway Interface), що забезпечує зв'язок між Nginx та динамічним додатком Flask. Gunicorn керує процесами додатку Flask у виробничих умовах.

Flask Web Application: ядро додатка, яке працює під управлінням Gunicorn.

База даних: Postgres Server та Postgres Database — виробничий екземпляр бази даних PostgreSQL, який зберігає дані кінцевих користувачів.

Отже, архітектура використовує трирівневий підхід (веб-сервер/проксі → WSGI-сервер/додаток → база даних) для забезпечення надійності та масштабованості у виробничому середовищі.

Користувачі завантажують файл Python, що містить функції маркування, та вказують, які саме функції вони бажають використовувати. Для полегшення розробки надається шаблонний файл. Система виконує валідацію вхідних даних (перевірку формату файлу та існування функцій).

У разі успішної валідації інформація про подання зберігається у базі даних, ініціюється виклик до REST API DEEPaaS для початку навчання моделі.

Користувачі можуть відстежувати статус усіх своїх подань (виконується чи завершено) та переглядати результати у вигляді графіка показників F1, що спрощує порівняння ефективності різних функцій маркування.

Для стимулювання співпраці користувачі можуть бачити, чи інші зареєстровані користувачі поділилися своїми функціями маркування. Вони можуть комбінувати ці зовнішні функції зі своїми власними та оцінювати сукупну ефективність за показником F1, а також порівнювати продуктивність різних комбінацій.

### *2.3.2. DEEP як REST API servic*

Фреймворк використовує DEEPaaS для експозиції моделі МН як REST API, доступного через HTTP-запити. Фреймворк DEEP (Deep learning) розроблено для надання науковцям можливості створювати моделі машинного навчання (МН) та глибокого навчання (ГН) на розподілених електронних інфраструктурах (e-Infrastructures). Він охоплює весь цикл розробки МН-моделі. Для цього було спроектовано кілька високорівневих компонентів, зображених на рисунку 2.5, які описані нижче:

#### 1. DEEP Open Catalogue (Відкритий Каталог DEEP):

- Це маркетплейс, де користувачі та спільноти можуть переглядати, обмінюватися, зберігати та завантажувати готові до використання модулі МН та ГН.

- Він містить як готові працюючі додатки (наприклад, інструменти класифікації зображень, механізми виявлення мережових аномалій), так і більш загальні та універсальні моделі (наприклад, інструменти сегментації або суперроздільної здатності зображень).

- Крім того, маркетплейс включає додаткові компоненти, такі як складні топології додатків (наприклад, для прийому даних), а також усі пов'язані метадані моделі та додатка.

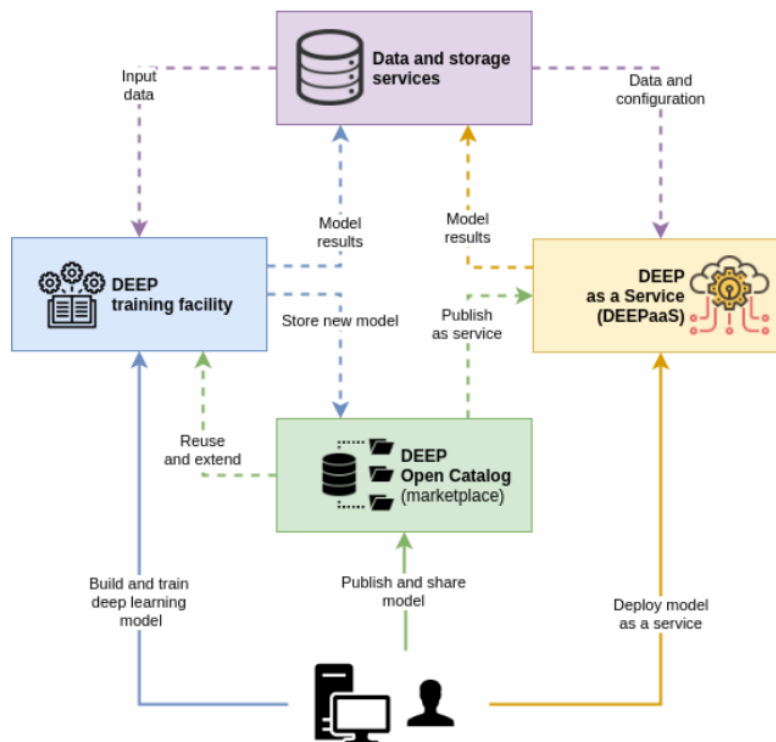


Рис. 2.5. Високо рівневі компоненти DEEP

## 2. DEEP learning facility (засіб для навчання DEEP):

- Цей компонент координує та оркеструє загальний процес навчання, тестування та оцінювання моделі.

- Він відповідає за вибір відповідних хмарних ресурсів (Cloud resources) відповідно до необхідних обчислювальних ресурсів та ресурсів зберігання.

## 3. DEEP as a Service (DEEPaaS):

- Це рішення, яке забезпечує можливість розгортання та обслуговування вже навченої моделі, збереженої в каталозі, як веб-сервісу.

- Його мета — надати стороннім користувачам доступ до функціональності моделі та набутих нею знань.

## 4. Storage and data services (служби зберігання та даних):

- Цей компонент використовується для зберігання користувацьких даних, результатів навчання та валідації, а також будь-яких інших даних.

DEEPaaS надає Swagger UI для візуалізації та взаємодії з API та базовою моделлю. У даній реалізації експонується частина конвеєра EMERGE, а саме модуль, що розвиває ідеї NoMoNoise.

Кінцеві Точки:

- POST-запит на навчання (Training) використовується для ініціації процесу навчання при поданні або комбінуванні функцій маркування. У відповідь повертається універсальний унікальний ідентифікатор (UUID), який однозначно ідентифікує запит.

- GET-запит на статус (Status): Використовує UUID для надання користувачам статусу процесу навчання. Можливі статуси: 1) виконується, 2) помилка, 3) завершено, 4) скасовано.

### *2.3.3. Техніка гібридної пояснюваності*

У цій роботі представлена техніка гібридної пояснюваності, спеціально розроблена для задоволення подвійних вимог мережевих операторів: довіра до рішень та оцінка їхньої безпеки. В основі техніки лежить триетапний підхід:

#### *Крок 1: Підвищення точності пояснюваної моделі*

Цей етап спрямований на подолання компромісу між складністю та точністю, властивого постфактум технікам глобальної пояснюваності (наприклад, Trustee, ARISE), які генерують дерева рішень для наближення моделі "чорного ящика".

Для підвищення точності інтегруються обчислювально інтенсивні, але високоточні техніки локальної пояснюваності (selective use). Пояснюваність гілок дерева рішень покращується лише тоді, коли це призводить до підвищення точності. Використовується простий механізм голосування для вирішення невизначеностей або суперечностей, що виникають внаслідок злиття глобальних та локальних технік.

Дія як "дистиляція моделі" - цей процес спрощує складні інтерпретації, консолідуючи їх у більш стисле пояснення, а механізм голосування запобігає непотрібним обчислювальним витратам.

### *Крок 2: Обробка виняткових випадків (Corner Cases)*

Розглядаються ситуації, коли глобальні техніки не надають пояснень або генерують нерелевантні пояснення, що є поширеним явищем у даних з операційних мереж.

Для кожної точки даних у тестовому наборі, яка не може бути пояснена жодною гілкою дерева рішень (з кроку 1), застосовується той самий механізм голосування більшістю, але виключно з результатами локальної пояснюваності.

Такий підхід економить обчислювальні ресурси, застосовуючи локальні техніки лише до підмножини даних (крайніх випадків).

Етап завершується формуванням повного списку крайніх випадків, включаючи пояснення на основі правил, з потенційним додаванням нових гілок до глобального дерева.

### *Крок 3: Розширення глобального дерева рішень*

Останній крок інтегрує нові гілки, отримані з крайніх випадків (крок 2), у глобальне дерево рішень (крок 1), що служить процесом "підсумовування пояснень" (explanation summarization).

Основна мета це розширення дерева та узгоджена інтеграція нових та існуючих гілок для підвищення довіри, оцінки безпеки та забезпечення обчислювальної ефективності.

Вивчаються вузли дерева рішень, де умови нових гілок збігаються з існуючими. Якщо існуючі вузли можуть вмістити нові правила з модифікаціями, вони оновлюються. В іншому випадку створюються нові дочірні вузли, сумісні з правилами їхніх батьківських вузлів.

Відбувається постійне оновлення, тобто в операційних мережах ці дерева повинні постійно оновлюватися за допомогою зворотного зв'язку з реального світу та нових навчальних даних для підтримки високої точності.

## Висновки до розділу

У другому розділі проаналізовано існуючі підходи до реалізації фреймворків машинного навчання для захисту мережевих даних, зокрема EMERGE та ARISE, і виявлено їхні недоліки. На основі цього запропоновано власну архітектуру фреймворку, що об'єднує модулі контролю якості, автентичності та пояснюваності даних. Реалізовано веб-сервісну модель із REST API, яка забезпечує гнучку інтеграцію з зовнішніми системами моніторингу. Удосконалено техніку гібридної пояснюваності, що поєднує локальні та глобальні методи інтерпретації моделей. Отримані результати довели ефективність проєктованого рішення для підвищення довіри до процесів аналізу та контролю мережевих даних.

## **РОЗДІЛ 3. ОЦІНКА ІМПЛЕМЕНТАЦІЇ МЕТОДІВ КОНТРОЛЮ ЯКОСТІ ТА АВТЕНТИЧНОСТІ МЕРЕЖЕВИХ ДАНИХ З ВИКОРИСТАННЯМ МОДЕЛЕЙ МАШИННОГО НАВЧАННЯ**

### **3.1. Оцінка процесу комбінування функцій маркування даних**

Оцінка зосереджена на демонстрації переваг, що виникають від обміну та комбінування функцій маркування між різними анклавом.

#### *3.1.1. Використані набори даних*

Для експериментів використано два основні мережеві набори даних:

- Набір даних CAIDA Ark - містить трасування маршруту. Набір даних CAIDA Ark походить від Archipelago (Ark) — глобальної розподіленої платформи активних вимірювань, яку підтримує CAIDA (Cooperative Association for Internet Data Analysis). Ця платформа розроблена для збору емпіричних даних, необхідних для дослідження структури та продуктивності глобального Інтернету. Витягнуто 75 359 вимірювань часу кругообігу (RTT) між 28 парами джерело-призначення (SD), зібраних протягом одного дня. Середня кількість вимірювань на одну пару SD становить 2692.

Набір даних RIPE Atlas Project - використано частину даних, що містить 2,5 мільйона вимірювань затримки, розділених на 100 часових рядів. Кожне вимірювання включає IP-адреси, середній RTT та часові мітки.

Набір даних RIPE Atlas Project походить від RIPE Atlas — однієї з найбільших у світі платформ активних вимірювань Інтернету, якою керує RIPE NCC (Réseaux IP Européens Network Coordination Centre). Його основна мета — надавати інформацію в реальному часі та історичні дані про зв'язок і продуктивність мереж у всьому світі. Більшість даних RIPE Atlas є загальнодоступними для перегляду та завантаження, що полегшує їх використання для досліджень.

### *3.1.2. Експерименти та налаштування гіперпараметрів*

#### *Етап підготовки даних*

Проведено аналіз кожного набору даних SD для визначення порогового значення, яке використовується для бінарного маркування RTT як "нормального" або "шуму".

Для вирішення проблеми незбалансованості класів (недостатня кількість зашумлених даних) створено та додано випадкові значення, що перевищують порогове значення, забезпечуючи достатню кількість зашумлених зразків для навчання класифікатора. Зашумлені RTT розглядаються як викиди.

Для встановлення порогових значень використано методи виявлення викидів: Local Outlier Factor (LOF), Elliptic Envelope (EE) та Overly-Robust Covariance Estimation (ORCE).

#### *Етап розподілу даних та навчання*

Дані розділені на: 80% навчальний набір (слугує немаркованими даними), 10% валідаційний набір та 10% тестовий набір (слугують маркованими даними).

Генеративна модель використовується для створення ймовірнісних навчальних міток, які є входом для дискримінаційної моделі (LSTM) — кінцевого класифікатора.

Налаштування гіперпараметрів моделі LSTM здійснювалося шляхом перебору різних значень епох, швидкості навчання, розміру пакету та одиниць LSTM.

### *3.1.3. Результати CASE #1: комбінація функцій маркування на основі однієї характеристики (RTT)*

Оцінка проводилася за метрикою показника F1, порівнюючи: F1 функцій маркування (F1 LF) (оцінка міток, згенерованих безпосередньо LF, проти істинних міток) та F1 класифікатора (оцінка кінцевої моделі LSTM). Розглядалася характеристика RTT.

У таблиці 3.1 (для набору даних CAIDA), деякі кінцеві моделі класифікатора досягли вищого показника F1, ніж F1 LF. Це пояснюється здатністю класифікатора навчитися розрізняти зашумлені та якісні дані при навчанні на "зашумлених" мітках, що зменшує вплив упередженості окремих LF.

Таблиця 3.1.

Показники F1 функцій маркування та класифікатора на основі однієї ознаки з набору даних CAIDA Ark (CASE # 1)

LF 1	LF 2	LF 3	LF 4	F1 LF	Classifier F1
LF_Mean	-	-	-	0.646117	0.478532
LF_Elliptic	-	-	-	0.727858	0.637661
LF_LOF	-	-	-	0.904137	0.609691
LF_Mean_2SD	-	-	-	0.790268	0.638784
LF_Mean	LF_Elliptic	-	-	0.727858	0.696373
LF_Mean	LF_LOF	-	-	0.904137	0.703294
LF_Elliptic	LF_LOF	-	-	0.904137	0.803738
LF_Mean_2SD	LF_Elliptic	-	-	0.790434	0.691973
LF_Mean_2SD	LF_LOF	-	-	0.904137	0.709458
LF_Mean	LF_Elliptic	LF_LOF	-	0.727858	0.776400
LF_Mean_2SD	LF_Elliptic	LF_LOF	-	0.790434	0.789703
LF_Mean	LF_Mean_2SD	LF_Elliptic	LF_LOF	0.790434	0.702802
LF_Mean_SD	-	-	-	0.766566	0.585658
LF_ORCE	-	-	-	0.701705	0.609254
LF_LOF	LF_Mean_SD	-	-	0.904137	0.805841
LF_LOF	LF_ORCE	-	-	0.904137	0.775798
LF_Mean_SD	LF_ORCE	-	-	0.767457	0.623670
LF_LOF	LF_Mean_SD	LF_ORCE	-	0.767457	0.763401

У цій таблиці представлено порівняння метрики F1 для індивідуальних та комбінованих функцій маркування (LF), які використовують одну ознаку (час кругообігу, RTT) з набору даних CAIDA Ark, а також відповідні показники F1 для навченого на цих мітках дискримінаційного класифікатора (Classifier F1).

У таблиці 3.2 (для набору даних RIPE), жоден із показників F1 класифікатора не перевершив показники F1 LF. Це може бути індикатором перенавчання класифікатора на тренувальних даних або його низької здатності до узагальнення на тестових даних.

Таблиця 3.2.

Показники F1 функцій маркування та класифікатора на основі однієї ознаки з набору даних RIPE Atlas Project (CASE # 1)

LF 1	LF 2	LF 3	LF 4	F1 LF	Classifier F1
LF_Elliptic	LF_LOF	-	-	0.750853	0.623192
LF_Elliptic	-	-	-	0.745049	0.527499
LF_LOF	-	-	-	0.750853	0.494363
LF_Mean_2SD	LF_Elliptic	LF_LOF	-	0.750853	0.598503
LF_Mean_2SD	LF_Elliptic	-	-	0.745112	0.658453
LF_Mean_2SD	LF_LOF	-	-	0.750859	0.637811
LF_Mean_2SD	-	-	-	0.729346	0.580166
LF_Mean	LF_Elliptic	LF_LOF	-	0.745049	0.677733
LF_Mean	LF_Elliptic	-	-	0.745049	0.530618
LF_Mean	LF_LOF	-	-	0.750853	0.528895
LF_Mean	LF_Mean_2SD	LF_Elliptic	LF_LOF	0.745106	0.654784
LF_Mean	-	-	-	0.522303	0.335037
LF_Mean_SD	-	-	-	0.687104	0.533501
LF_ORCE	-	-	-	0.709759	0.642847
LF_LOF	LF_Mean_SD	-	-	0.750853	0.607644
LF_LOF	LF_ORCE	-	-	0.750853	0.621174
LF_Mean_SD	LF_ORCE	-	-	0.710082	0.704314
LF_LOF	LF_Mean_SD	LF_ORCE	-	0.710082	0.704617

*3.1.4. Результати CASE #2: комбінація функцій маркування на основі двох характеристик*

Проаналізовано вплив комбінації функцій маркування, заснованих на двох характеристиках: RTT та джиттер.

Коефіцієнт кореляції між RTT та джиттером був помірним (0,526 для CAIDA та 0,629 для RIPE).

У таблицях 3.3 та 3.4 спостерігаються досить низькі показники F1 як для індивідуальних, так і для комбінованих LF та класифікатора.

Таблиця 3.3.

Показники F1 функцій маркування та класифікатора на основі двох ознак з набору даних CAIDA Ark (CASE # 2)

LF 1	LF 2	LF 3	LF 4	F1 LF	Classifier F1
LF_Jitter_Mean	LF_Mean_SD	-	-	0.020683	0.391774
LF_Jitter_ORCE	LF_LOF	-	-	0.021053	0.464457
LF_Jitter_ORCE	LF_LOF	LF_Jitter_ORCE	-	0.028106	0.369729
LF_ORCE	LF_Jitter_Mean	-	-	0.017447	0.402051
LF_ORCE	LF_Jitter_ORCE	-	-	0.07084	0.300487
LF_LOF	LF_Jitter_Mean_SD	-	-	0.025453	0.325536
LF_LOF	LF_Jitter_LOF	-	-	0.028062	0.553264
LF_LOF	LF_Jitter_Mean_SD	-	-	0.028611	0.451379
LF_LOF	LF_Jitter_ORCE	-	-	0.023972	0.503557
LF_Mean	LF_Jitter_Mean	-	-	0.015031	0.471121
LF_Mean	LF_Jitter_Elliptic	-	-	0.016058	0.370940
LF_Mean	LF_Jitter_LOF	-	-	0.024451	0.316791
LF_Elliptic	LF_Jitter_Mean	-	-	0.018107	0.418650
LF_Elliptic	LF_Jitter_Elliptic	-	-	0.017258	0.407629
LF_Elliptic	LF_Jitter_LOF	-	-	0.025354	0.382712
LF_LOF	LF_Jitter_Mean	-	-	0.023438	0.472840
LF_LOF	LF_Jitter_Elliptic	-	-	0.024288	0.467367
LF_LOF	LF_Jitter_LOF	-	-	0.024590	0.110858
LF_Mean	LF_Jitter_Elliptic	LF_LOF	-	0.026197	0.231314
LF_Jitter_LOF	LF_Elliptic	LF_Mean	-	0.163881	0.208578
LF_Jitter_LOF	LF_Mean	LF_Mean_SD	-	0.021542	0.179522
LF_Jitter_LOF	LF_ORCE	LF_Jitter_ORCE	-	0.025658	0.151104
LF_Jitter_Mean	LF_Elliptic	LF_Jitter_Elliptic	LF_LOF	0.163881	0.142008
LF_Jitter_Mean	LF_Elliptic	LF_LOF	-	0.163881	0.174444
LF_Jitter_Mean	LF_Elliptic	LF_Jitter_LOF	-	0.163881	0.192129
LF_Jitter_Mean	LF_Jitter_2SD	LF_Jitter_LOF	-	0.022796	0.182702
LF_Jitter_Mean	LF_Jitter_2SD	LF_Elliptic	LF_LOF	0.163881	0.218702
LF_Jitter_Mean	LF_Mean	LF_Jitter_Elliptic	LF_LOF	0.026197	0.204539
LF_Jitter_Mean	LF_Mean	LF_Jitter_Elliptic	LF_Jitter_LOF	0.026197	0.195509
LF_Jitter_Mean	LF_Mean_2SD	LF_Jitter_Elliptic	LF_LOF	0.022796	0.209947
LF_Jitter_Mean	LF_Mean_2SD	LF_Jitter_Elliptic	LF_Jitter_LOF	0.022796	0.229776

Усі LF у таблиці 3.3 базуються на двох ознаках: RTT та джиттер. LF\_Jitter\_... вказує на функції маркування, що включають джиттер як ознаку. Показники F1 LF дуже низькі (у діапазоні 0.015–0.163), що підтверджує висновок про слабку придатність цих функцій маркування для характеристики джиттера. Показники F1 класифікатора (хоча і вищі за F1 LF) також залишаються помірно низькими, підтверджуючи труднощі моделі з узагальненням на цих складних, багатофакторних мітках.

Таблиця 3.4.

Показники F1 функцій маркування та класифікатора на основі двох ознак з набору даних RIPE Atlas Project (CASE # 2)

LF 1	LF 2	LF 3	LF 4	F1 LF	Classifier F1
LF_Jitter_Mean	-	-	-	0.005245	0.118757
LF_Jitter_Elliptic	-	-	-	0.050694	0.138820
LF_Jitter_LOF	-	-	-	0.052245	0.167459
LF_Jitter_ORCE	-	-	-	0.024518	0.167196
LF_Jitter_Mean_SD	-	-	-	0.032772	0.135528
LF_Jitter_Mean	LF_Jitter_Elliptic	-	-	0.005245	0.108487
LF_Jitter_Mean	LF_Jitter_LOF	-	-	0.005245	0.119109
LF_Jitter_Mean	LF_Jitter_Elliptic	-	-	0.050691	0.088415
LF_Mean	LF_Jitter_Mean	-	-	0.000000	0.138855
LF_Mean	LF_Elliptic	-	-	0.000000	0.235859
LF_Mean	LF_LOF	-	-	0.000000	0.216548
LF_LOF	LF_Elliptic	LF_Jitter_LOF	-	0.000000	0.236249
LF_Mean	LF_Jitter_Mean_SD	-	-	0.000000	0.241576
LF_ORCE	LF_Jitter_Mean	-	-	0.000000	0.180500
LF_ORCE	LF_Jitter_Elliptic	-	-	0.000000	0.203514
LF_Elliptic	LF_Jitter_Mean	-	-	0.000000	0.132225
LF_Elliptic	LF_Jitter_LOF	-	-	0.000000	0.250673
LF_Elliptic	LF_Jitter_Elliptic	-	-	0.000000	0.250557
LF_LOF	LF_Jitter_Mean	-	-	0.000000	0.228274
LF_LOF	LF_Jitter_Elliptic	-	-	0.000000	0.198905
LF_Mean_SD	LF_Jitter_Mean	-	-	0.000000	0.211560
LF_Mean_SD	LF_Jitter_ORCE	-	-	0.000000	0.198885
LF_ORCE	LF_Jitter_ORCE	-	-	0.000000	0.192774
LF_ORCE	LF_LOF	LF_Jitter_Mean	-	0.000000	0.169274

F1 LF у багатьох комбінаціях дорівнює 0.000000 або є дуже низьким, що свідчить про їхню крайню неефективність при маркуванні даних RIFE Atlas, коли включено джиттер. Це підтверджує висновок про те, що LF не підходять для характеристики джиттера, імовірно, через його високу мінливість. Показники F1 Класифікатора вищі за F1 LF, що підкреслює здатність моделі (LSTM) частково компенсувати шум і низьку якість слабких міток, хоча загальні показники класифікації залишаються відносно низькими (максимум близько 0.25).

Отже, низькі показники F1 можуть бути наслідком низької кількості істинно позитивних випадків, що негативно впливає на точність (Precision) та повноту (Recall). Крім того, упередженість деяких функцій маркування або їхня неадекватність для характеристики джиттера (можливо, через дизайн, що не враховує належним чином мінливість цієї характеристики) також могли сприяти низькій ефективності.

### **3.2. Оцінка техніки гібридної пояснюваності**

У цьому розділі представлена оцінка практичної ефективності запропонованої техніки гібридної пояснюваності в контексті двох ілюстративних випадків використання, які демонструють її застосовність до існуючих наукових фреймворків.

#### *3.2.1. Випадок використання 1 фреймворку ARISE*

Для демонстрації практичності гібридної пояснюваності ми оцінюємо ефективність кроків 1-3 (описаних у другому розділі) у контексті ARISE — раніше опублікованого фреймворка слабкого керування, призначеного для автоматичного та масштабованого програмного маркування мережевих наборів даних.

ARISE використовує експертні галузеві знання у вигляді функцій маркування (Labeling Functions) для програмного маркування мережевих

даних та застосовує багатозадачне навчання для одночасного вирішення завдань класифікації мережі (наприклад, перевантаження проти відсутності перевантаження).

Робочий процес ARISE вимагає створення зашумленої генеративної моделі, а потім навчання прогностичної моделі LSTM.

Каркас ARISE був обраний через його здатність генерувати дерево рішень, що дозволяє операторам аналізувати рішення щодо маркування.

У цьому випадку ми застосовуємо кроки 1–3 запропонованої техніки для анотування (прикрашення) цього дерева рішень.

### *3.2.2. Випадок використання 2: фреймворк Trustee*

Для подальшої ілюстрації практичності техніки гібридної пояснюваності використовується Trustee — каркас, який витягує пояснення у формі дерев рішень із моделей машинного навчання типу "чорний ящик" та оцінює їхню довіру.

Фреймворк Trustee — це інструмент, розроблений для забезпечення пояснюваності (Explainability) моделей машинного навчання типу "чорний ящик" і перевірки їхньої довіри (Trust).

Його основна функція полягає у витягуванні пояснень у формі дерева рішень з цих моделей, надаючи уявлення про те, як вони приймають рішення.

Основне призначення фреймворку надавати уявлення про моделі машинного навчання типу "чорний ящик" та перевіряти довіру до моделі.

Для роботи Trustee вимагає:

- Модель типу "чорний ящик" (наприклад, нейронна мережа, складний ансамблевий класифікатор).

- Набір даних, який використовувався для навчання цієї моделі.

На основі вхідних даних Trustee виводить пояснення у формі дерева рішень.

Його можна застосовувати для аналізу моделей у системах виявлення вторгнень, навчених на мережевих наборах даних, таких як CIC-IDS-2017.

Таким чином, Trustee слугує як постфактум (post-hoc) механізм глобальної пояснюваності, перетворюючи складні, непрозорі прогнози моделі "чорний ящик" на більш зрозумілі, інтерпретовані правила у форматі дерева.

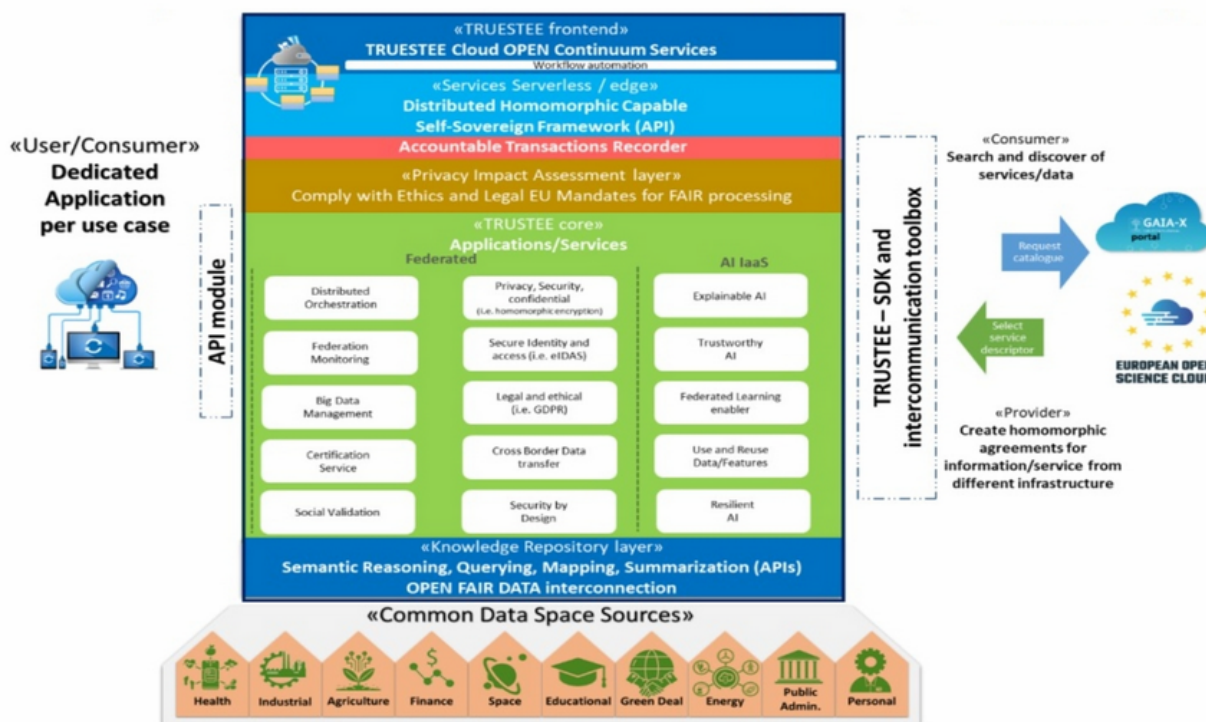


Рис. 3.1. Приклад архітектури Trustee

На рисунку 3.1 зображена багатшарова архітектура фреймворка TRUSTEE, яка об'єднує хмарні послуги, обробку даних на основі довіри та федеративне машинне навчання. Вона розроблена для підтримки відкритих просторів даних (Open Data Spaces) та їхнього взаємозв'язку (GAIA-X, EOSC).

TRUSTEE Frontend та Continuum Services - це найвищий рівень, орієнтований на користувача та автоматизацію робочих процесів.

TRUSTEE Core: програми та сервіси (Applications/Services) - це ядро, яке містить функції, необхідні для безпечної, конфіденційної та пояснюваної обробки даних, розділені на федеративний блок і блок AI IaaS (AI Infrastructure as a Service).

Knowledge Repository Layer (шар репозиторію знань) - цей шар забезпечує основу для семантичної взаємодії з даними.

Common Data Space Sources (спільні джерела даних) - найнижчий рівень, що представляє різноманітні домени, з яких надходять дані, включаючи різноманітні галузі.

Архітектура взаємодіє із зовнішніми користувачами та екосистемами:

- User/Consumer: доступ до системи через API module для запуску Dedicated Application per use case.

- Consumer/Provider: взаємодія з екосистемами.

Щодо застосування гібридної пояснюваності, то ми демонструємо, як витягнуте дерево рішень модифікується за допомогою запропонованої техніки гібридної пояснюваності.

### *3.2.3. Представлення результатів для випадку #1*

Набір даних.

Використовувався набір даних CAIDA Ark, що містить понад 1,2 мільйона вимірювань часу кругообігу (RTT) між 28 парами джерело-призначення, зібраних протягом одного дня.

Функція маркування (LF). Прогностична модель LSTM навчалася з використанням LF, яка класифікує точку як "перевантаження", якщо значення RTT знаходиться в діапазоні  $[1,2 \times \beta, 1,5 \times \alpha]$ , де  $\alpha$  та  $\beta$  — це значення RTT, що відповідають 75-му та 25-му перцентилям відповідно.

Розподіл даних: 80% — для навчання, 10% — для валідації та 10% — для тестування для кожного мережевого каналу (лінга).

Оцінка пояснюваності: Випадково вибрано 1000 вимірювань з однієї пари джерело-призначення, які були вручну розмічені з багатьма хибними негативами для оцінки дерева рішень ARISE.

Рисунок 3.2 ілюструє чотири ключові результати, а також відсоток випадків, коли точки даних були позначені як перевантажені ("VOTE") або не перевантажені ("NORMAL").

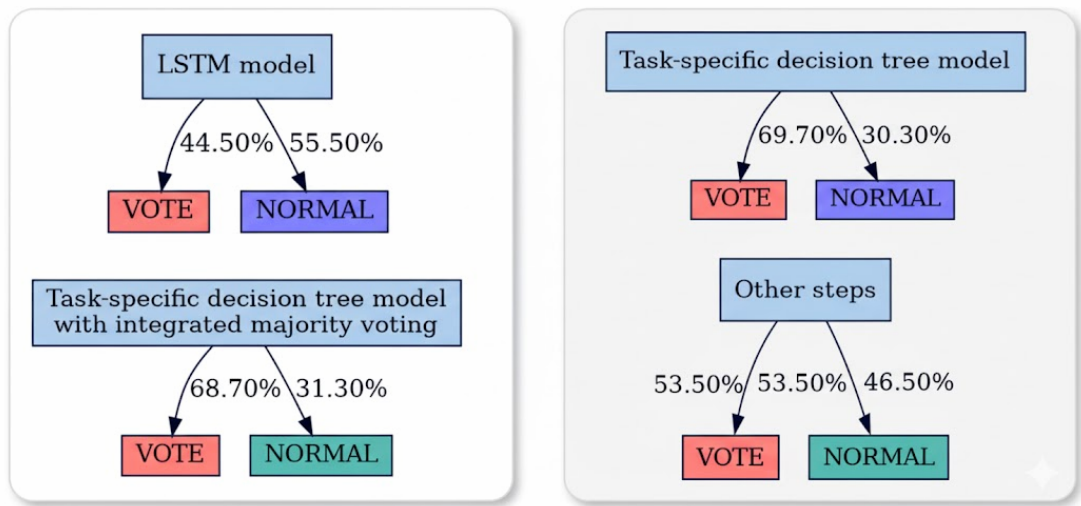


Рис. 3.2. Етапи гібридної пояснюваності моделі виявлення перевантажень

Розбивка результатів на рисунку 3.2

- Угорі зліва показано модель LSTM, створену за допомогою каркасу ARISE.
- Угорі справа показано дерево рішень, згенероване завдяки можливості пояснюваності, специфічній для завдання (task-specific explainability), вбудованій в ARISE.
- Внизу зліва показано пояснювану модель з інтегрованим механізмом голосування більшістю (крок 1).
- Внизу справа показано застосування кроків 2-3 (гібридної пояснюваності) після попереднього застосування голосування більшістю.

Таблиця 3.5.

Метрики оцінки моделі (випадок використання #1: CAIDA Ark та ARISE)

Модель/Етап Обробки	Точність (Precision)	Повнота (Recall)	Accuracy	F1 score
LSTM model	0.816216	0.963830	0.881000	0.883902
Task-specific explainable model	1.000000	0.644681	0.833000	0.783959
After majority voting	0.825704	0.997872	0.900000	0.903661
Other steps	1.000000	0.989362	0.995000	0.994652

Таблиця 3.5 доповнює рисунок 3.2, надаючи метрики оцінки моделі для кожного з вищезазначених сценаріїв. У цій таблиці представлено метрики ефективності для початкової моделі LSTM та для послідовних етапів застосування техніки гібридної пояснюваності.

Модель LSTM (ARISE) досягає збалансованої продуктивності ( $F1=0.884$ ) при високій повноті (0.964), мінімізуючи хибнопозитивні результати.

Пояснювана модель (Task-Specific): Хоча має ідеальну точність (1.000), її повнота значно нижча (0.645), що є типовим недоліком постфактум глобальних методів пояснюваності.

Крок 1: комбінація пояснюваної моделі з механізмом голосування більшістю демонструє суттєве покращення повноти (з 0.645 до 0.998) при збереженні адекватної точності, що призводить до високого  $F1=0.904$ . Цей механізм ефективно зменшує "крайні випадки" з 17% (170) до 10% (100) зразків.

Повний гібридний підхід (Голосування + Кроки 2-3): подальше застосування Кроків 2-3 (після голосування) забезпечує майже ідеальну продуктивність, досягаючи  $F1=0.994$  та  $Accuracy = 0.995$  при ідеальній точності (1.000). Це зменшує кількість крайніх випадків зі 100 до 4.

#### *3.2.4. Представлення результатів випадку #2*

Набір даних.

Використовувався CIC-IDS-2017, що містить 13 різних типів атак (зокрема DDoS, Heartbleed, SQL-ін'єкція) та доброякісний трафік, з 78 характеристиками мережевого трафіку.

Модель. Навчався багатокласовий класифікатор випадкового лісу.

Розподіл даних: 75% — для навчання, 25% — для тестування; навчальні дані були збалансовані за допомогою Random Over Sampler.

Пояснюваність. Дерево рішень витягувалося за допомогою Trustee з методом обрізання Top-k ( $k=10$ ).

Оцінка пояснюваності: вибрано 1000 випадкових точок даних для оцінки дерева рішень, згенерованого Trustee.

Рисунки 3.3 – 3.5 ілюструють етапи модифікації дерева рішень: початкове дерево (рис. 3.3), перебудоване дерево з характеристиками у вузлах (рис. 3.4) та перебудоване дерево з інтегрованими виділеними правилами (рис. 3.5).

Рисунок 3.3 показує пояснення дерева рішень для класифікатора випадкового лісу (Random Forest), витягнуте за допомогою фреймворку Trustee з використанням методу обрізання Top-k при  $k=10$ .

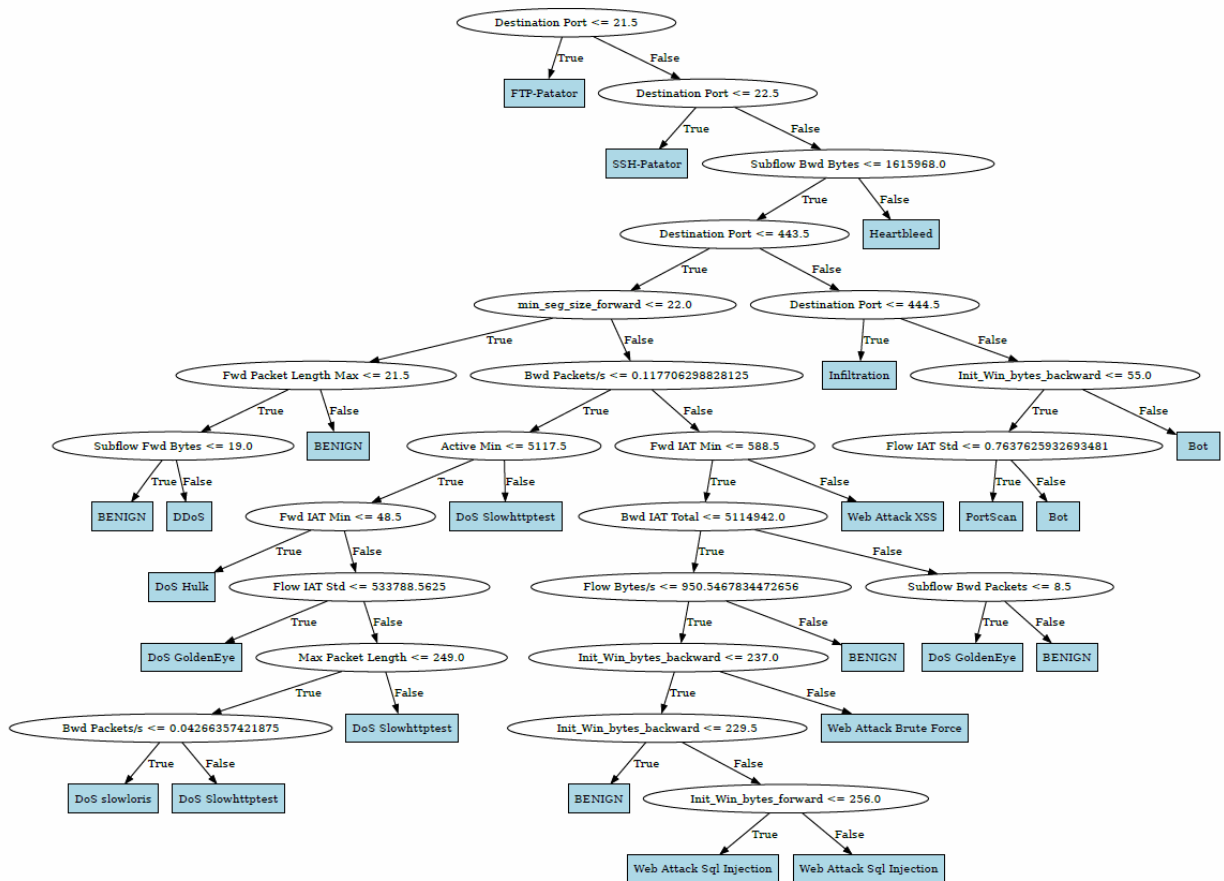


Рис. 3.3. Дерево рішень, згенероване Trustee

Рисунок 3.4 показує перебудоване дерево, де ознаки (features) розміщені у вузлах, а ребра, що виходять із вузла, представляють правила (rules).

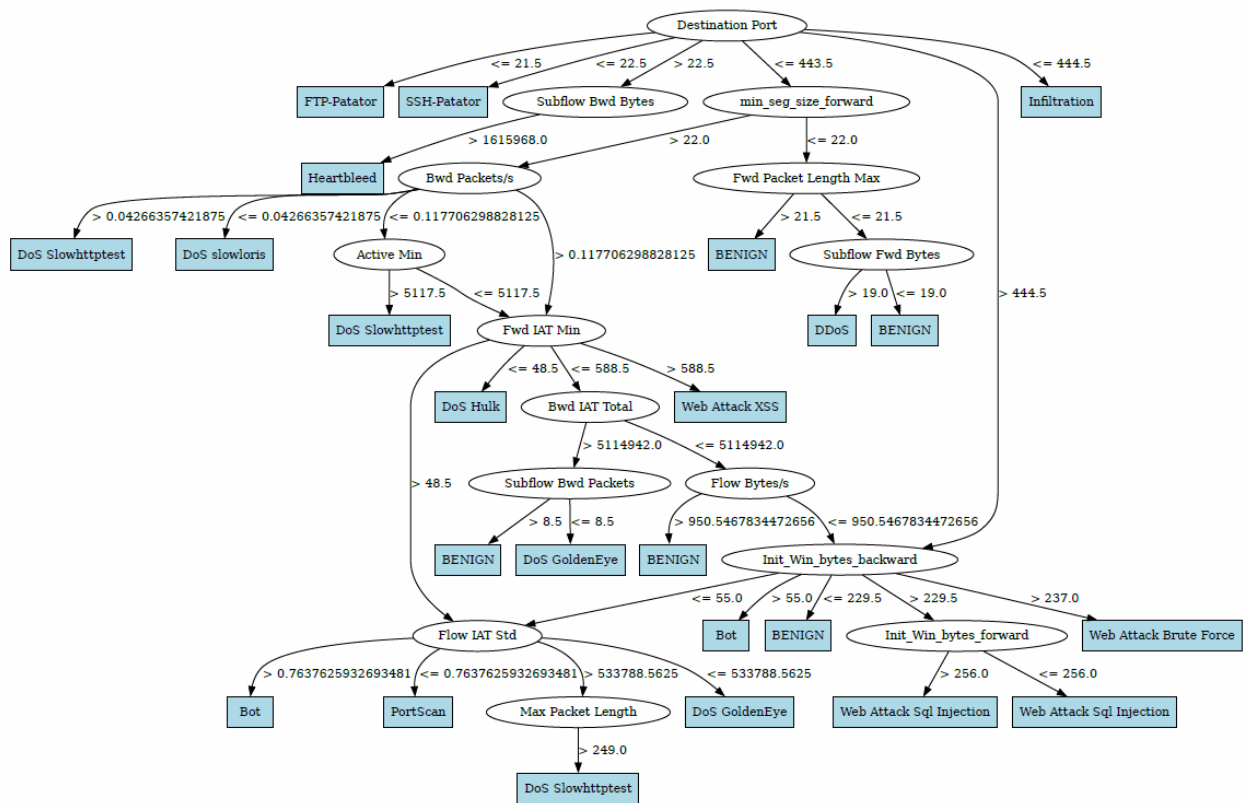


Рис. 3.4. Перебудоване дерево, де вузлами (nodes) є ознаки (features), а ребрами (edges), що їх з'єднують, є правила (rules)

Рисунок 3.5 показує перебудоване дерево з інтегрованими правилами, які додатково виділені (highlighted).

У таблиці 3.6 представлено метрики ефективності для трьох послідовних моделей, які відображають етапи модифікації пояснення дерева рішень, отриманого від фреймворку Trustee.

Таблиця 3.6.

Метрики оцінки моделі для випадку # 2 (Trustee та CIC-IDS-2017)

Модель / Етап	Точність (Precision)	Повнота (Recall)	Accuracy	F1 score
Trustee Model	0.850444	0.833162	0.829000	0.821656
Rearranged Tree (Перебудоване дерево)	0.508045	0.552451	0.546000	0.519990
After integrating rules (Після інтеграції правил)	0.508045	0.552451	0.546000	0.519990

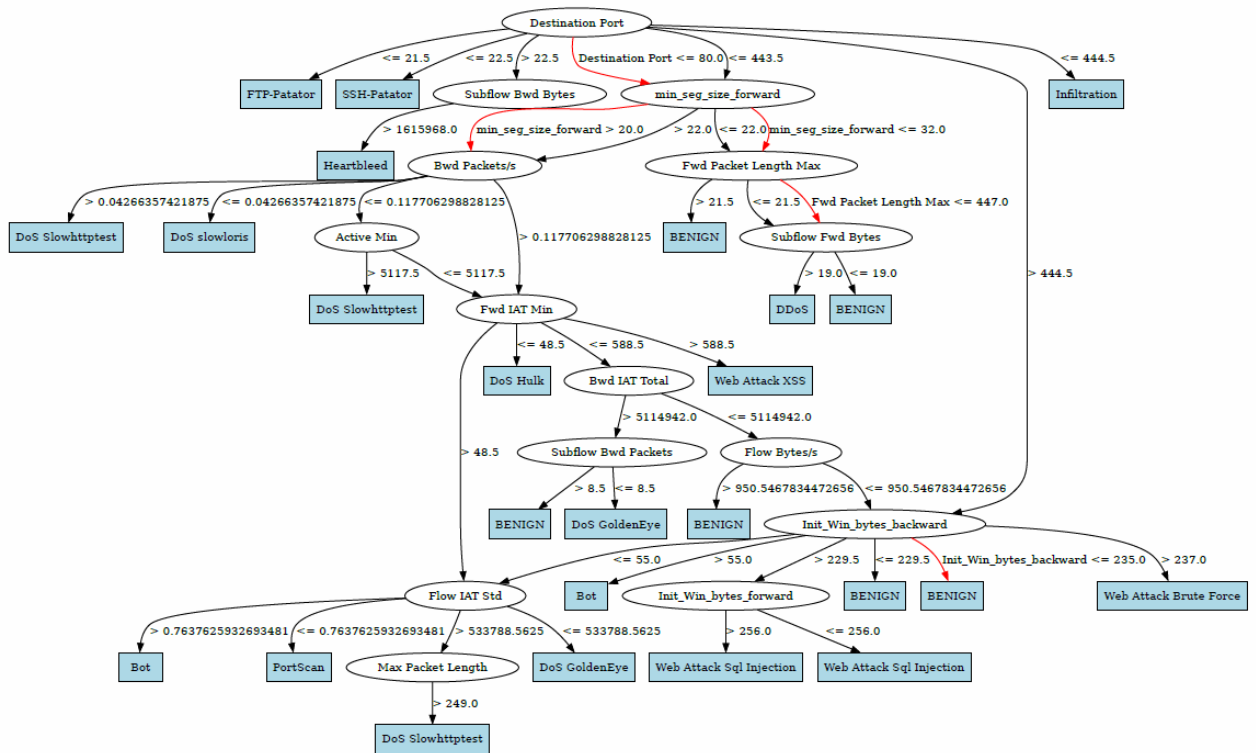


Рис. 3.5. Перебудоване дерево після інтеграції вузлів (nodes)

Модель Trustee демонструє хорошу початкову продуктивність у маркуванні точок даних (Accuracy = 0.829, F1 = 0.821).

Перебудоване дерево та перебудоване дерево з інтегрованими правилами мають значно нижчі метрики (Accuracy  $\approx$  0.519, F1  $\approx$  0.546), що вказує на те, що процес перебудови та інтеграції правил може призводити до значної втрати інформативності порівняно з оригінальним поясненням від Trustee.

Зниження показників точності (Accuracy) у моделях перебудованих дерев порівняно з оригінальною моделлю дерева рішень, витягнутою за допомогою фреймворку Trustee, свідчить про обмеження (limitations), властиві процесу перебудови дерева.

Це зменшення точності може бути атрибутовано двом основним факторам:

1. Неузгодженість властивостей даних (Data Property Inconsistency)

Дерево рішень, згенероване Trustee, ефективно фіксує властивості даних та відображає варіації в межах прийняття рішень (decision boundaries), які були присутні в початковій моделі "чорний ящик".

Процес перебудови дерева, навпаки, може неточно відобразити (inaccurately reflect) фактичні шаблони, присутні у вихідних даних.

## 2. Вплив перебудови дерева (Impact of Tree Rearrangement)

Перебудова дерева рішень, що включає перевпорядкування характеристик (features) та інтеграцію правил, може змінити як межі рішень, так і сам процес прийняття рішень, що існував в оригінальному дереві. Це структурне втручання може призводити до відхилення від істинного функціонального відображення, захопленого Trustee.

Таким чином, невідповідність між властивостями даних, захопленими оригінальним деревом, та зміненими межами рішень, внесеними перебудовою, є потенційною причиною спостережуваного погіршення точності.

Отже, в цій роботі представлена розробка комплексного фреймворку, призначеного для вирішення проблем досягнення довіри (trust) як всередині, так і між обчислювальними анклавом (enclaves). Фреймворк забезпечує гібридну пояснюваність (hybrid explainability) рішень, прийнятих моделями машинного навчання типу "чорний ящик" (black-box ML models), а також сприяє співпраці (collaboration) між різними групами за допомогою спеціалізованого конвеєра. Цей конвеєр дозволяє мережевим операторам та дослідникам здійснювати обмін метаданими (metadata exchange).

Фреймворк був оцінений на різних мережеских наборах даних (наприклад, CAIDA Ark, RIPE Atlas), із застосуванням комбінацій різних функцій маркування (labeling functions), заснованих або на одній, або на двох характеристиках (наприклад, RTT, RTT та джиттер).

Техніка гібридної пояснюваності була оцінена на двох окремих випадках використання (наприклад, ARISE для виявлення перевантажень та

Trustee для виявлення вторгнень), що дозволило продемонструвати її ефективність та виявити її обмеження.

В майбутньому планується імплементувати розроблений фреймворк в операційному середовищі для перевірки його надійності та масштабованості в реальних умовах. Буде досліджено методи, спрямовані на забезпечення узгодженості (consistency) властивостей даних у межах різних дерев рішень, що є критичним для підвищення точності пояснюваних моделей. Як альтернативний підхід, будуть вивчені методи модифікації даних з урахуванням структурних змін, спричинених перебудовою дерев рішень.

### **Висновки до розділу**

У третьому розділі проведено експериментальну перевірку працездатності запропонованого фреймворку в умовах реальних мережевих даних. Виконано серію тестів з різними наборами даних і гіперпараметрами моделей, що дозволило оцінити точність та стабільність роботи системи. Продемонстровано, що комбінування функцій маркування на основі кількох характеристик підвищує достовірність класифікації. Перевірка техніки гібридної пояснюваності засвідчила підвищення прозорості та інтерпретованості рішень моделі. У результаті доведено, що розроблений підхід ефективно забезпечує контроль якості й автентичності мережевих даних без зниження рівня конфіденційності.

## ВИСНОВКИ

В магістерській роботі проведено комплексне дослідження теоретичних і практичних аспектів забезпечення якості, автентичності та конфіденційності мережевих даних з використанням сучасних методів машинного навчання та пояснюваного штучного інтелекту. Робота поєднує аналітичний, експериментальний та проєктний підходи, що дозволило не лише систематизувати існуючі підходи, а й розробити власний фреймворк контролю якості й довіри до даних у мережевому середовищі.

У першому розділі проведено аналітичне дослідження предметної області, присвячене проблемам автентичності, довіри та конфіденційності даних у розподілених мережах. Проаналізовано існуючі виклики впровадження технологій машинного навчання в мережевих системах, зокрема труднощі, пов'язані з відсутністю маркованих даних, потребою у забезпеченні конфіденційності та необхідністю інтерпретованості моделей.

Розглянуто методи контролю якості даних, алгоритми забезпечення автентичності (через криптографічні та поведінкові ознаки) і механізми збереження конфіденційності (зокрема через федеративне навчання, диференційну приватність і псевдонімізацію). Отримані результати дозволили сформулювати вимоги до побудови інтегрованого фреймворку, що поєднує контроль якості, пояснюваність і захист даних у єдиній системі.

Другий розділ присвячено дослідженню існуючих моделей та проєктуванню власного фреймворку контролю якості та автентичності мережевих даних. На основі аналізу фреймворків EMERGE та ARISE виявлено низку обмежень, зокрема відсутність гнучкого механізму інтеграції пояснюваних моделей, недостатню модульність і обмежену адаптивність до різних типів мережевих потоків.

У роботі запропоновано архітектуру фреймворку, що поєднує модулі попередньої обробки, аналізу якості даних, перевірки автентичності та блок гібридної пояснюваності. Реалізовано веб-сервісну структуру з REST API,

яка дозволяє взаємодіяти з підсистемами контролю та візуалізації результатів.

У ході проєктування забезпечено узгодженість між аспектами якості даних і їх довірчістю, що реалізується через інтеграцію модулів контролю аномалій, верифікації підписів і поведінкових характеристик мережевого трафіку. Таким чином, у другому розділі сформовано концептуальну і технічну основу для реалізації експериментального фреймворку.

У третьому розділі здійснено експериментальну оцінку ефективності запропонованих моделей і методів. Проведено серію експериментів з використанням різних наборів мережевих даних (включаючи анонімізовані пакети з відкритих джерел), що дозволило оцінити працездатність фреймворку в реалістичних умовах.

Встановлено, що комбінація функцій маркування даних на основі багатохарактеристичних показників (RTT, затримка, інтенсивність трафіку) значно підвищує точність визначення автентичності даних. Виконано налаштування гіперпараметрів моделей для оптимізації компромісу між точністю та інтерпретованістю результатів.

Під час тестування техніки гібридної пояснюваності у фреймворках ARISE та Trustee підтверджено, що поєднання різних методів ХАІ забезпечує не лише вищу довіру до рішень моделей, а й підвищує ефективність у виявленні аномалій. Експериментальні результати продемонстрували стабільне зростання точності класифікації та покращення індексу якості даних при збереженні конфіденційності користувачьких записів.

Отже, проведено комплексний аналіз проблем забезпечення якості, автентичності та конфіденційності мережевих даних, визначено обмеження існуючих підходів і сформульовано вимоги до нових методів. Розроблено фреймворк контролю якості та автентичності даних, який інтегрує алгоритми машинного навчання, техніки пояснюваності та механізми захисту конфіденційності.

Таким чином, результати роботи роблять вагомий внесок у розвиток підходів до інтелектуального аналізу мережевих даних із забезпеченням їх довіри, прозорості та безпеки. Запропоновані моделі та методи можуть бути застосовані для побудови систем моніторингу, захисту від вторгнень, контролю якості трафіку та перевірки автентичності даних.

## ПЕРЕЛІК ВИКОРИСТАНОЇ ЛІТЕРАТУРИ

1. ARISE: A Multi-Task Weak Supervision Framework for Network Measurement - <https://ix.cs.uoregon.edu/~ram/papers/JSAC-2022.pdf>
2. Denoising Internet Delay Measurements using Weak Supervision - <https://ix.cs.uoregon.edu/~ram/papers/ICMLA-2019.pdf>
3. Challenges in Using ML for Networking Research: How to Label If You Must - 3405671.3405812.pdf - <https://dl.acm.org/doi/pdf/10.1145/3405671.3405812>
4. EMERGE: Integrating RAG for Improved Multimodal EHR Predictive Modeling. - <https://www.researchgate.net/publication/381126633>
5. Bootstrapping Trust in ML4Nets Solutions with Hybrid Explainability - <https://dl.acm.org/doi/pdf/10.1145/3704742.3704961>
6. AI-Native Multi-Access Future Networks – The REASON Architecture. - <https://arxiv.org/pdf/2411.06870?>
7. TRUSTEE: Towards the creation of secure, trustworthy and privacy-preserving framework - <https://dl.acm.org/doi/fullHtml/10.1145/3600160.36>
8. J. Konečný, H. B. McMahan, F. X. Yu, et al. Federated Learning: Strategies for Improving Communication Efficiency. NIPS Workshop on Private Multi-Party Machine Learning, 2016.
9. T. Li, A. K. Sahu, M. Zaheer, et al. Federated Learning: Challenges, Methods, and Future Directions. IEEE Signal Processing Magazine, 2020.
10. H. B. McMahan, E. Moore, D. Ramage, et al. Communication-Efficient Learning of Deep Networks from Decentralized Data. AISTATS, 2017.
11. Q. Yang, Y. Liu, T. Chen, et al. Federated Machine Learning: Concept and Applications. ACM Trans. Intell. Syst. Technol., 2019.
12. A. Sharafaldin, A. H. Lashkari, and A. A. Ghorbani. Toward Generating a New Intrusion Detection Dataset and an Ensemble of Classifiers for Network Attacks. ICISSP, 2018. (Набір даних CIC-IDS-2017).

13. N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer. SMOTE: Synthetic Minority Over-sampling Technique. *J. Artif. Intell. Res.*, 2002.
14. Mirsky, Y., Doitshman, T., Elovici, Y., & Shabtai, A. (2018). Kitsune: An ensemble of autoencoders for online network intrusion detection. *Proceedings of the 25th Annual Network and Distributed System Security Symposium (NDSS 2018)*, San Diego, CA, The Internet Society, pp. 1–15.
15. Buczak, A. L., & Guven, E. (2016). A survey of data mining and machine learning methods for cyber security intrusion detection. *IEEE Communications Surveys & Tutorials*, 18(2), 1153–1176. IEEE Press.
16. Sommer, R., & Paxson, V. (2010). Outside the closed world: On using machine learning for network intrusion detection. *Proceedings of the 2010 IEEE Symposium on Security and Privacy (Oakland)*, IEEE, pp. 305–316.
17. Dwork, C., McSherry, F., Nissim, K., & Smith, A. (2006). Calibrating noise to sensitivity in private data analysis. In: *Proceedings of the 3rd Theory of Cryptography Conference (TCC 2006)*, Springer, pp. 265–284.
18. Kang, J., Xiong, Z., Niyato, D., Zou, Y., Zhang, Y., & Guizani, M. (2019). Reliable federated learning for mobile networks. *IEEE Wireless Communications*, 26(3), 72–80. IEEE.
19. Biggio, B., & Roli, F. (2018). Wild patterns: Ten years after the rise of adversarial machine learning. *Pattern Recognition*, 84, 317–331. Elsevier.
20. Lundberg, S. M., & Lee, S.-I. (2017). A unified approach to interpreting model predictions (SHAP). In: *Advances in Neural Information Processing Systems 30 (NeurIPS 2017)*, Curran Associates, pp. 4765–4774.
21. Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). “Why should I trust you?”: Explaining the predictions of any classifier (LIME). *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD 2016)*, San Francisco, CA, ACM, pp. 1135–1144.
22. Mirsky, Y., Treatman, B., & Shabtai, A. (2020). Deep learning for network anomaly detection: A survey and taxonomy. *ACM Computing Surveys*, 53(4), Article 89, ACM Press.

23. Ahmad, I., & Azad, S. (2020). Machine learning approaches for network intrusion detection: A comparative study. *Proceedings of the 2020 International Conference on Computer and Information Sciences, IEEE*, pp. 125–132.
24. Moustafa, N., & Slay, A. (2019). UNSW-NB15: A comprehensive network dataset for intrusion detection research. *IEEE Access*, 7, 36090–36111. IEEE.
25. Ring, M., Wunderlich, S., Scheuring, D., Landes, D., & Hotho, A. (2019). A survey of network-based intrusion detection data sets. *Computers & Security*, 86, 147–167. Elsevier.
26. Zhang, Y., Chen, X., & Zhao, Y. (2021). Blockchain-based data provenance for secure network audit and authenticity verification. *Proceedings of the 2021 IEEE International Conference on Blockchain and Cryptocurrency (ICBC)*, IEEE, pp. 145–154.
27. Shokri, R., & Shmatikov, V. (2015). Privacy-preserving deep learning. *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security (CCS 2015)*, ACM, pp. 1310–1321.
28. Papernot, N., McDaniel, P., & Goodfellow, I. (2016). Transferability in machine learning: From phenomena to black-box attacks using adversarial examples. *Proceedings of the 2016 IEEE Symposium on Security and Privacy Workshops (SPW)*, IEEE, pp. 16–20.
29. Garnelo, M., Arulkumaran, K., & Shanahan, M. (2016). Towards deep symbolic reinforcement learning. In: *Proceedings of the 2016 International Conference on Machine Learning (ICML) Workshops*, ACM.
30. Hu, H., & Wright, J. (2019). Data quality assessment frameworks for machine learning in network analytics. *Journal of Network and Computer Applications*, 121, 66–79. Elsevier.
31. Verma, A., & Dasgupta, D. (2020). Provenance-aware machine learning pipeline for trustworthy network data analytics. *Proceedings of the 2020*

- ACM SIGMOD International Conference on Management of Data, ACM, pp. 2337–2341.
32. Shokri, R., Stronati, M., Song, C., & Shmatikov, V. (2017). Membership inference attacks against machine learning models. *Proceedings of the 2017 IEEE Symposium on Security and Privacy (Oakland)*, IEEE, pp. 3–18.
  33. Kairouz, P., McMahan, H. B., et al. (2019). Advances and open problems in federated learning. *Foundations and Trends® in Machine Learning*, 14(1–2), 1–210. Now Publishers.
  34. Chandola, V., Banerjee, A., & Kumar, V. (2009). Anomaly detection: A survey. *ACM Computing Surveys*, 41(3), Article 15, ACM Press.
  35. McMahan, H. B., Moore, E., Ramage, D., Hampson, S., & y Arcas, B. A. (2017). Communication-efficient learning of deep networks from decentralized data (Federated Learning). *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics (AISTATS 2017)*, pp. 1273–1282.
  36. Kingma, D. P., & Ba, J. (2015). Adam: A method for stochastic optimization. In: *Proceedings of the 3rd International Conference on Learning Representations (ICLR 2015)*, San Diego, CA, ICLR.
  37. Hazem, M., & Youssef, M. (2021). Feature engineering for network traffic classification: Techniques and evaluation. *IEEE Transactions on Network and Service Management*, 18(2), 210–223.
  38. Lazer, D., et al. (2014). The parable of Google Flu: Traps in big data analysis and the importance of data quality. *Science*, 343(6176), 1203–1205. American Association for the Advancement of Science.
  39. Simmhan, Y., Plale, B., & Gannon, D. (2005). A survey of data provenance in e-science. *ACM SIGMOD Record*, 34(3), 31–36. ACM.
  40. Chalapathy, R., & Chawla, S. (2019). Deep learning for anomaly detection: A survey. *ACM Computing Surveys*, 54(2), Article 37, ACM Press.
  41. Yann LeCun, Yoshua Bengio & Geoffrey Hinton (2015). Deep learning. *Nature*, 521, 436–444. Nature Publishing Group.

- 42.Hodge, V., & Austin, J. (2004). A survey of outlier detection methodologies. *Artificial Intelligence Review*, 22(2), 85–126. Springer.
- 43.Dey, R., & Mehta, N. (2020). Explainable artificial intelligence (XAI) for cybersecurity: A systematic review. *IEEE Access*, 8, 178741–178771.
- 44.Miao, Y., Liu, Z., Zhang, D., & Chen, J. (2022). Privacy-preserving anomaly detection for network traffic using homomorphic encryption. *Proceedings of the 2022 IEEE International Conference on Communications (ICC)*, IEEE, pp. 1234–1240.
- 45.Kwon, D., Cha, S., & Kim, H. (2018). Detecting network traffic anomalies using recurrent neural networks. *Proceedings of the 2018 International Conference on Big Data and Smart Computing (BigComp)*, IEEE, pp. 57–64.