

Міністерство освіти і науки України  
Харківський національний університет імені В. Н. Каразіна  
Навчально-науковий інститут комп'ютерних наук та штучного інтелекту  
Кафедра комп'ютерних систем та робототехніки

«Затверджую»  
в.о. завідуючого кафедри  
комп'ютерних систем та робототехніки  
к. ф.-м. н., доцент Максим ХРУСЛОВ  
\_\_\_\_\_ к. ф.-м. н., доцент Максим ХРУСЛОВ  
«\_\_\_» червня 2025 р.

## Пояснювальна записка

до кваліфікаційної роботи  
бакалавра

на тему: «МОДЕЛЬ РОЗПІЗНАВАННЯ ОБ'ЄКТІВ РІЗНОЇ  
ПРИРОДИ У СИСТЕМАХ СПОСТЕРЕЖЕННЯ»

Спеціальність 123 – Комп'ютерна інженерія  
Галузь знань: 12 – Інформаційні технології.  
Освітня програма «Комп'ютерна інженерія».

**Захищено на засіданні**

**Екзаменаційної комісії № 44**  
протокол № \_\_ від \_\_.06.2025 р.

Оцінка \_\_\_\_\_ / \_\_\_\_\_

**Голова Екзаменаційної комісії**

\_\_\_\_\_ **ЧУГАЙ А. М.**

**Виконав:**

Студент групи КІ– 41  
**МОТИКА Володимир Назарович**



**Керівник:**

к.т.н., доцент кафедри комп'ютерних  
систем та робототехніки

**СТРІЛЕЦЬ Вікторія Євгенівна**



**Рецензент:** к.т.н., доцент, доцент  
кафедри математичного моделювання та  
аналізу даних

**КОРОБЧИНСЬКИЙ Кирил Петрович**



Харків – 2025

## АНОТАЦІЯ

Пояснювальна записка до кваліфікаційної роботи містить вступ, три розділи, висновки, список джерел і додатки. Загальний обсяг – 69 сторінок, із яких 45 сторінок основної частини, що включає 16 рисунків, 1 таблиці та 16 джерел.

**Мета роботи** – підвищення ефективності розпізнавання об'єктів різної природи на зображеннях за допомогою одноетапної (YOLOv8) та дворівневої (Faster R-CNN) моделей розпізнавання об'єктів різної природи, а також інтегрованого рішення, що поєднує їхні переваги.

**Об'єкт дослідження** – процес автоматизованого розпізнавання об'єктів на зображеннях і відео в системах спостереження.

**Предмет дослідження** – методи та алгоритми розпізнавання об'єктів, архітектурні рішення та практичні підходи до проєктування, тренування й інтеграції двох моделей детекції об'єктів – YOLOv8 і Faster R-CNN – у єдиній системі спостереження.

**Проблема**, яка вирішується в кваліфікаційній роботі – зменшення помилок і покращення продуктивності у задачах розпізнавання об'єктів.

**Область застосування** – системи відеоаналітики, безпеки та автоматизованого контролю.

У практичній частині дослідження створено програму, яка завантажує зображення, застосовує до них обидві моделі та порівнює результати за якісними метриками.

**Ключові слова:** *глибинне навчання, детекція об'єктів, YOLOv8, Faster R-CNN, комп'ютерний зір, точність, швидкодія.*

## ABSTRACT

The explanatory note for the qualification work includes an introduction, three chapters, conclusions, a list of references, and appendices. The total length is 69 pages, of which 45 pages comprise the main body, containing 16 figures, 1 table, and 16 references.

**Objective.** The purpose of this work is to improve the efficiency of detecting objects of various types in images by leveraging both a one-stage detector (YOLOv8) and a two-stage detector (Faster R-CNN), and an integrated solution that combines their respective advantages.

**Object of study.** The object of this research is the process of automated object recognition in images and video within surveillance systems.

**Subject of study.** The subject of this research is the methods and algorithms for object recognition, architectural designs, and practical approaches to designing, training, and integrating two object-detection models – YOLOv8 and Faster R-CNN – into a unified surveillance system.

**Problem statement.** The central problem addressed in this qualification work is the reduction of detection errors and the improvement of overall performance in object recognition tasks.

**Application domain.** This work is applicable to video analytics systems, security platforms, and automated monitoring and control systems.

**Practical component.** In the practical part of this study, a software application was developed that loads input images, applies both detection models, and compares their outputs using quantitative quality metrics.

**Keywords:** *deep learning, object detection, YOLOv8, Faster R-CNN, computer vision, accuracy, real-time performance.*

## ЗМІСТ

ПЕРЕЛІК УМОВНИХ ПОЗНАЧЕНЬ.....	5
ВСТУП .....	7
РОЗДІЛ 1. АНАЛІЗ ЗАДАЧІ РОЗПІЗНАВАННЯ ОБ’ЄКТІВ РІЗНОЇ ПРИРОДИ .....	10
1.1 Підходи до розпізнавання образів. Огляд та класифікація .....	10
1.2 Порівняння методів розпізнавання .....	15
1.3 Тенденції та напрямки розвитку систем розпізнавання об’єктів.....	18
Висновок до розділу 1 .....	20
РОЗДІЛ 2. ПРОЄКТУВАННЯ МОДЕЛЕЙ ДЛЯ РОЗПІЗНАВАННЯ ОБ’ЄКТІВ.....	22
2.1 Архітектура та принцип роботи моделі YOLOv8.....	22
2.2 Архітектура та принцип роботи моделі Faster R-CNN .....	24
2.3 Проєктування інтегрованої системи розпізнавання об’єктів .....	26
2.4 Навчальний набір даних (Pascal VOC 2012) .....	30
Висновок до розділу 2 .....	33
РОЗДІЛ 3. ПРОГРАМНА РЕАЛІЗАЦІЯ, ТРЕНУВАННЯ ТА ТЕСТУВАННЯ МОДЕЛЕЙ РОЗПІЗНАВАННЯ .....	35
3.1 Реалізація та навчання моделі YOLOv8 .....	35
3.2 Реалізація та навчання моделі Faster R-CNN .....	38
3.3 Порівняльний аналіз результатів розпізнавання і оцінка ефективності	42
3.4 Аналіз процесу навчання моделей .....	45
Висновки до розділу 3 .....	50
ВИСНОВКИ.....	51
СПИСОК ВИКОРИСТАНИХ ДЖЕРЕЛ.....	53
ДОДАТКИ.....	56
Додаток А.....	56
Додаток Б .....	58
Додаток В.....	61
Додаток Г .....	68

## ПЕРЕЛІК УМОВНИХ ПОЗНАЧЕНЬ

CNN	–	Convolutional Neural Network – згорткова нейронна мережа
RPN	–	Region Proposal Network – мережа регіональних пропозицій
IoU	–	Intersection over Union – коефіцієнт перетину прогнозованої та істинної області
mAP	–	mean Average Precision – середня точність (усереднена по класах)
mAP@0.5	–	mAP при $\text{IoU} \geq 0.5$
mAP@[0.5:0.95]	–	mAP усереднена по порогах IoU від 0.5 до 0.95 з кроком 0.05
P (Precision)	–	Precision – точність (правильні передбачення / усі передбачені положення)
R (Recall)	–	Recall – повнота (правильні передбачення / усі істинні об'єкти)
FPN	–	Feature Pyramid Network – піраміда ознак для багатомасштабного виявлення об'єктів
PAN	–	Path Aggregation Network – мережа для агрегування ознак різних рівнів
NMS	–	Non-Maximum Suppression – алгоритм відсікання зайвих перекриваючихся боксів
FPS	–	Frames Per Second – кількість кадрів за секунду (метрика швидкодії)
GPU	–	Graphics Processing Unit – графічний процесор
CPU	–	Central Processing Unit – центральний процесор
DL	–	Deep Learning – глибоке навчання
ML	–	Machine Learning – машинне навчання

SIFT	–	Scale-Invariant Feature Transform – інваріантний до масштабу дескриптор
HOG	–	Histogram of Oriented Gradients – гістограма напрямків градієнтів
SSD	–	Single Shot Multibox Detector – одноетапний детектор з фіксованими «дефолтними» вікнами
FCN	–	Fully Convolutional Network – повністю згортова мережа для пікельної сегментації
U-Net	–	U-architecture Network – згортова мережа з U-подібними пропусками і рескіп-лінками
DETR	–	Detection Transformer – детектор на основі архітектури трансформера
ViT	–	Vision Transformer – архітектура трансформера для обробки зображень

## ВСТУП

**Актуальність дослідження.** Сучасні системи відеоспостереження та моніторингу вимагають високої точності та швидкодії при виявленні та класифікації об'єктів у потоках зображень. З ростом обсягів відеоданих та поширенням «розумних» камер виникає потреба в автоматизованих рішеннях, здатних в реальному часі обробляти кадри й приймати рішення про наявність та тип об'єктів (людей, транспортних засобів, тварин тощо). Традиційні методи на основі ручного виділення ознак дедалі більше поступаються глибинним підходам, які забезпечують вищу точність та адаптивність. Наразі серед широко застосованих і перспективних архітектур детекції об'єктів особливо виділяються одноетапні (YOLOv8) та дворівневі (Faster R-CNN) моделі. Вибір оптимального рішення чи їх інтеграція є важливим для різноцільових застосунків – від контролю доступу й безпеки до інтелектуального аналізу трафіку.

**Об'єктом дослідження** є процес автоматизованого розпізнавання об'єктів на зображеннях і відео в системах спостереження.

**Предметом дослідження** є методи та алгоритми розпізнавання об'єктів, архітектурні рішення та практичні підходи до проектування, тренування й інтеграції двох моделей детекції об'єктів – YOLOv8 і Faster R-CNN – у єдиній системі спостереження.

**Метою** кваліфікаційної роботи є підвищення ефективності розпізнавання об'єктів різної природи на зображеннях за допомогою одноетапної (YOLOv8) та дворівневої (Faster R-CNN) моделей розпізнавання об'єктів різної природи, а також інтегрованого рішення, що поєднує їхні переваги.

### **Завдання дослідження:**

1. Провести огляд та класифікацію сучасних підходів до розпізнавання образів, виділити переваги та обмеження класичних і глибинних методів.

2. Розробити та описати архітектуру моделі YOLOv8 з урахуванням її основних компонентів: backbone, neck, head.
3. Розробити та описати архітектуру моделі Faster R-CNN, включно з RPN, RoI Pooling і класифікаційним модулем.
4. Побудувати інтегровану систему детекції, що комбінує YOLOv8 та Faster R-CNN, із застосуванням механізму злиття прогнозів (Fusion).
5. Реалізувати моделі в середовищі Google Colab на базі PyTorch із використанням датасету Pascal VOC.
6. Провести навчання та тестування моделей, зібрати основні показники ефективності (precision, recall, mAP@0.5, mAP@[0.5:0.95]).
7. Виконати порівняльний аналіз отриманих результатів, зробити висновки про доцільність використання кожного підходу в залежності від вимог системи.

#### **Методи дослідження:**

- аналіз літературних джерел – ознайомлення з класичними методами (SIFT, HOG, template matching) та глибинними підходами (CNN, R-CNN, YOLO, SSD, DETR);
- моделювання та проектування – побудова архітектур моделей та інтеграційної схеми;
- експериментальний метод – реалізація, навчання та валідація моделей на практичному датасеті Pascal VOC;
- кількісна оцінка – використання метрик precision, recall та mAP для зіставлення якості детекції;
- порівняльний аналіз – зіставлення параметрів швидкодії (FPS) і точності обох підходів.

**Практичне значення роботи.** Результати дослідження можуть бути використані для:

- розробки інтелектуальних систем відеоспостереження з високою пропускнуою здатністю у реальному часі;

- впровадження інтегрованих рішень на основі YOLOv8 і Faster R-CNN у завданнях охорони, контролю доступу та аналітики відеопотоку;
- оптимізації моделей для Edge-пристроїв із обмеженими ресурсами шляхом поєднання швидких і точних детекторів;
- подальших досліджень у галузі гібридних архітектур, self-supervised learning та трансформерів у комп'ютерному зорі.

Таким чином, виконана в межах кваліфікаційної роботи розробка та аналіз моделей розпізнавання об'єктів можуть бути основою для створення практично застосовних, високоефективних і гнучких систем відеоспостереження.

## РОЗДІЛ 1

### АНАЛІЗ ЗАДАЧІ РОЗПІЗНАВАННЯ ОБ'ЄКТІВ РІЗНОЇ ПРИРОДИ

#### 1.1 Підходи до розпізнавання образів. Огляд та класифікація

Сучасні системи розпізнавання образів поєднують багатоступеневі алгоритми, що включають підготовку даних, виділення ознак, навчання моделі та класифікацію чи детекцію об'єктів. Основні задачі комп'ютерного зору поділяють на *класифікацію образів* (визначення, до якого класу належить ціле зображення), *детекцію об'єктів* (визначення положення й класу кількох об'єктів у зображенні) та *сегментацію* (построкова класифікація пікселів з розрізненням різних об'єктів або лише класів; у т.ч. семантична або сегментація екземплярів (інстанс)). Задача класифікації зображень є базовою: система видає один чи кілька тегів для всього зображення, а детектори виявляють багатокласові об'єкти й креслять обмежувальні рамки або маски. Семантична сегментація позначає кожен піксель класом фону або об'єкта, а інстанс-сегментація (як у масці Mask R-CNN) додатково відрізняє окремі екземпляри одного класу.

Підходи до розпізнавання образів традиційно поділяють на класичні (традиційні) методи, що базуються на ручному виділенні ознак і статистичній обробці, та сучасні підходи на основі глибинного навчання (головним чином з використанням згорткових нейронних мереж, CNN). З класичних методів широко застосовуються *шаблонні алгоритми* та *методи на основі ознак* (feature-based), де для кожного об'єкта вручну обчислюють дескриптори. Наприклад, алгоритм *Scale-Invariant Feature Transform* (SIFT) надає інваріантні до масштабу та повороту ключові точки з описом на основі градієнтів інтенсивності [12]. SIFT-ознаки забезпечують надійне відновлення локальних ознак в умовах змін освітлення і перспективи. Також відомі дескриптори *Histogram of Oriented Gradients* (HOG), запропоновані для виявлення людей. HOG-ознаки статистично підраховують гістограми напрямків градієнтів у регіоні зображення; вони дали досконалу точність на

контрольному наборі для виявлення пішоходів. Інші ручні ознаки (SURF, ORB, LBP та інші) спочатку широко застосовувалися в розпізнаванні образів, але поступово пішли на другий план.

**Шаблонне зіставлення (*template matching*)** – класичний метод, що полягає у пошуку прямого співпадіння зразка (шаблону) в зображенні [6]. Цей підхід оцінює ступінь схожості між об'єктом і заданим шаблоном (наприклад, за сумою квадратів різниць або нормалізованою крос-кореляцією). Шаблонний метод простий у реалізації: кожен фрагмент вхідного зображення по черзі підставляється під вікно шаблону, і обчислюється функція невідповідності. Проте шаблонне зіставлення чутливе до змін масштабу, повороту, освітлення і шуму. З часом на його основі з'явилися вдосконалення (деформовані шаблони, багатовимірні шаблони), але фундаментально ці методи мають обмежену гнучкість. Зазвичай їх поєднують з іншими алгоритмами або застосовують у завданнях з контрольованим середовищем.

Класичний статистичний підхід використовує ручні ознаки та прості моделі. Наприклад, після виділення ознак (SIFT/HOG) часто застосовують методи класифікації: **лінійні алгоритми** (лінійний дискримінантний аналіз, порогова класифікація), **методи найближчих сусідів** (k-NN), **машини опорних векторів** (SVM) тощо. Такі методи можуть виконувати класифікацію чи оцінювати ймовірності класів на основі статистичних моделей (Гаусівські суміші, байєсівські класифікатори). Також використовувалися й **багатошарові нейронні мережі з одним–двома прихованими шарами** – наприклад, перцептрони та мережі зі звужувальною/розширювальною архітектурою (на кшталт LeNet для розпізнавання цифр) – але їх потужність обмежувалася малою глибиною і залежністю від довільно відібраних ознак.

Однак останні десятиліття революція у розпізнаванні образів пов'язана із застосуванням глибинного навчання, зокрема **згорткових нейронних мереж** (CNN). Глибинні мережі здатні автоматично виявляти репрезентації ознак в даних, навчатися кінцево-кінцево від пікселів до класів і досягати безпрецедентно високих результатів. Наприклад, відомі архітектури для

класифікації зображень – AlexNet (2012), VGG, GoogleNet/Inception, ResNet і DenseNet – продовжують покращувати точність на стандартних наборах (ImageNet, CIFAR) завдяки збільшенню глибини та складності моделей. Зокрема, ResNet з 152 шарами згортання здобув 1-е місце на ILSVRC2015, забезпечивши помилку лише 3.57% на ImageNet [7]. Глибокі CNN одночасно виявляють низькорівневі ієрархічні ознаки та інтегрують їх на наступних рівнях, що дало революційні покращення класифікації (наприклад, менше 5% помилок на ImageNet у 2016 році).

Підходи на базі глибоких мереж застосовують також до задачі виявлення (детекції) об'єктів. Одним із ключових проривів стало впровадження регіональних мереж: R-CNN (2014), Fast R-CNN (2015) і Faster R-CNN (2015). У Faster R-CNN автори представили *мережу пропозицій регіонів* (RPN), яка інтегрує етап генерації кандидатних областей прямо всередині згорткової мережі [15]. Швидка мережа R-CNN отримала на PASCAL VOC 2007 mAP  $\approx 73.2\%$  з частотою всього 5 кадрів/с (включно з етапом RPN). Проте, незважаючи на високу точність, такі двоетапні методи відносно повільні через складність етапів пропозицій.

Альтернативний клас моделей – *одноетапні детектори*. Найвідоміші – класів YOLO («You Only Look Once») та SSD («Single Shot Multibox Detector»). У YOLO (Redmon et al., CVPR 2016) виявлення об'єктів формулюється як *єдине регресійне завдання* [14]: нейронна мережа одночасно видає координати вікон і класи об'єктів у цих вікнах. Базова модель YOLO обробляє зображення розміром 448×448 за 45 кадрів/с; спрощена версія (Fast YOLO) – до 155 кадрів/с. Таким чином, YOLO забезпечує дуже високу швидкість, придатну для реального часу, поступаючись дещо точністю: у роботі автори зауважували, що YOLO більше помиляється в локалізації, ніж двоетапні методи, але рідше видає хибно-позитивні детекції фону. SSD (Liu et al., ECCV 2016) також усуває етап регіональних пропозицій, використовуючи набір «дефолтних вікон» різних масштабів і співвідношень сторін на кожному рівні признаков [10]. SSD-прототип досягає приблизно 74% mAP на VOC2007 при

59 кадрів/с на NVIDIA Titan X, суттєво випереджаючи Faster R-CNN за швидкістю і зберігаючи порівнянну точність. SSD демонструє перевагу над іншими одноетапними детекторами у точності при малому розмірі вхідного зображення.

Для задачі сегментації (піксельної класифікації) глибинні мережі також дали важливі результати. Ранні роботи на цьому напрямі – *Fully Convolutional Networks* (FCN, 2015), які перетворюють звичайні CNN на “повністю згорткові”, що дозволяють видавати карти класів довільного розміру [11]. FCN надав змогу навчати мережі напряму «піксель-в-піксель» і досягати досконалих результатів на наборі PASCAL VOC 2012 (наприклад, близько 62% середньої точності за індексом IoU). На основі FCN згодом з’явилися й глибші архітектури: *U-Net* (2015), спочатку для біомедичних даних, та *Mask R-CNN* (2017), які поєднують детекцію і сегментацію інстансів. U-Net відрізняється бістатичною U-подібною структурою, яка поєднує контекстні характеристики з високодеталізованими картами, що дає дуже точне поділення об’єктів [16]. Він простий у навчанні і обробляє одне 512×512 зображення менше, ніж за секунду. Mask R-CNN розширює Faster R-CNN, додаючи на паралельній гілці згорткову маску для кожного Розуміння Регіону (RoI), і при цьому забезпечує точне сегментування окремих екземплярів [8]. Маскова гілка мало ускладнює Faster R-CNN, і система залишається відносно швидкою (близько 5 кадрів/с на сучасному GPU).

Новітнім напрямком стали *трансформери для зору*. Удосконалено концепцію детекторів на основі «запитів» (queries) – прикладом є *DETR* (Detection Transformer, ECCV 2020). [2] DETR розглядає детекцію як задачу множинного передбачення (set prediction) та використовує архітектуру енкодер–декодер з механізмом самоуваги. Головні плюси DETR – відсутність необхідності в ручних компонентах (засори, NMS, якорі). Автори показали, що DETR дає порівнянну з Faster R-CNN якість на COCO при трішки повільнішому виконанні, а також гарно масштабується до паноптичної сегментації. Крім того, *Vision Transformer* (ViT, 2020) застосовує чисту

архітектуру трансформера до зображень: воно розбивається на патчі й обробляється як послідовність «слів».[5] Показано, що при достатньо великій передтренуванні ViT досягає чудових результатів на класифікації, навіть кращих, ніж традиційні CNN.

Таким чином, сучасні підходи до розпізнавання образів охоплюють широкий спектр методів: від класичних шаблонних і статистичних методів з ручно визначеними ознаками (SIFT, HOG, бібліотеки дескрипторів, прості нейронні мережі, SVM тощо) до різноманіття архітектур глибоких мереж. Сучасні алгоритми згорткових нейронних мереж дозволяють будувати єдині системи, які навчаються «від пікселів до результатів», мінімізуючи необхідність ручної настройки ознак. Типова структура таких систем включає етапи: подання зображення на вхід CNN, виділення ієрархічних ознак, генерування передбачень класів (для класифікації або детекції) або прогнозу піксельних масок (для сегментації). На рис. 1.1 умовно показано загальну схему процесу розпізнавання образів із зазначенням основних компонентів.

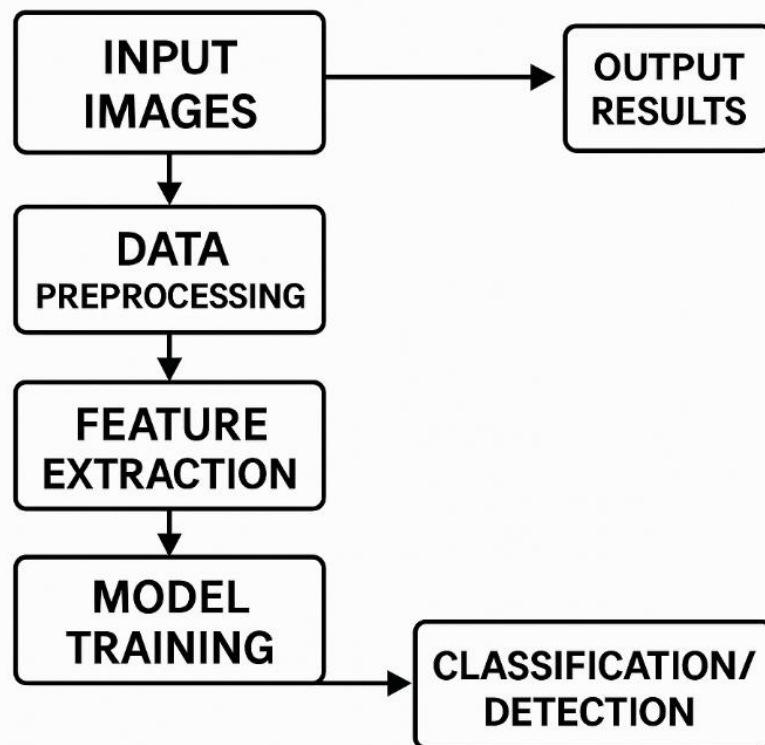


Рисунок 1.1 – Основні етапи процесу розпізнавання образів

## 1.2 Порівняння методів розпізнавання

Різні методи розпізнавання образів оцінюють за сукупністю критеріїв: *точність* (accuracy), *швидкодія* (продуктивність), *ресурсомісткість* (пам'ять, обчислювальні потужності), *інтерпретованість* результатів тощо. Класичні методи та сучасні глибокі моделі мають протилежні властивості за цими вимірами.

**Точність і надійність.** Загалом глибокі моделі (CNN) переважають класичні у складних задачах з нерегулярними візуальними об'єктами. Ручні ознаки обмежені областю застосування: наприклад, HOG чудово працює для статичних об'єктів людського тіла в контролюваному середовищі, але вимагає «налаштувань» при різких змінах ракурсу чи освітлення.

Сучасні CNN здатні навчитися інваріантним ознакам із набору даних, тому демонструють набагато вищу точність на загальних завданнях. Зокрема, з появою глибоких архітектур помилкові класифікації стали рідкістю навіть на багатокласових задачах (ImageNet і т.д.), де похибка була знижена з  $\approx 15\%$  (до CNN) до 3–4% на передових мережах, таких як ResNet. У детекції, хоча Faster R-CNN і SSD дають порівнянну точність, двоетапні методи (R-CNN сімейства) традиційно показують дещо кращий mAP завдяки ретельнішому обробленню регіональних пропозицій. Наприклад, Faster R-CNN 2015 дав mAP 73.2% на VOC2007, тоді як перша версія YOLO – лише  $\approx 63\%$  при набагато вищій швидкодії. SSD ж показує mAP  $\sim 74\%$  при 59 FPS. Маски сегментації, що формуються Mask R-CNN, забезпечують високу точність розділення окремих об'єктів, але лише ціною суттєвого ускладнення моделі відносно Faster R-CNN. U-Net та FCN демонструють високу точність семантичної сегментації за відносно невеликих обсягів даних (завдяки агресивній аугментації). Іншими словами, за критерієм якості розпізнавання перевага на боці глибоких моделей. Водночас інтерпретація рішень CNN суттєво важча: класична модель “пояснюється” вручну підібраними фільтрами або граничними значеннями,

тоді як глибока CNN – «чорний ящик», де ознаки автоматично формуються в латентному просторі.

**Швидкодія (FPS) та реальний час.** Одноетапні детектори (YOLO, SSD) спроектовані для високої пропускної здатності. Наприклад, базовий YOLO працює ~45 кадр/с, Fast YOLO – понад 155 кадр/с. SSD на Titan X обробляє ~59 кадр/с при збереженні майже такої ж точності, як Faster R-CNN. Для порівняння, двоетапний Faster R-CNN (VGG-16) демонстрував лише ~5 кадр/с. Це означає, що в системах, де потрібна обробка відео в реальному часі або мінімальна затримка (наприклад, відеоспостереження чи роботи), одноетапні мережі часто переважні. У класифікаційних задачах пікову швидкість дають оптимізовані мережі (наприклад, ResNet-50 може передавати кілька сотень зображень за секунду на сучасному GPU). У сегментації U-Net на GPU виконує один зразок за час  $<1$  с, що також задовольняє вимоги більшості реальних застосувань у медицині та робототехніці. Водночас, класичні методи зазвичай легші: шаблонне зіставлення чи обчислення HOG-ознак на CPU може бути доволі швидким на невеликих зображеннях, проте при збільшенні розміру даних їхня складність також зростає.

**Ресурсномісткість і складність моделі.** Глибинні мережі зазвичай потребують значних обчислювальних ресурсів і пам'яті. Наприклад, ResNet-152 містить сотні мільйонів параметрів [13], а детектори на його базі (Faster R-CNN) вимагають великої відеопам'яті для роботи з пакетами. Додаткова обчислювальна складність необхідна для двоетапних детекторів: етап генерації пропозицій був значним «вузьким місцем», яке усунув Faster R-CNN, але плата за це – спільне навчання двох компонентів. У порівнянні, одноетапні детектори (YOLO, SSD) обчислюють всі етапи в єдиному проході через мережу, що знижує загальні затрати, але все одно вимагає великих прискорювачів для високої роздільності. Навпаки, класичні алгоритми найчастіше легші. Наприклад, побудова гіста HOG і виклик лінійного SVM працюють досить швидко навіть на CPU, і пам'ять потрібна тільки для збереження вагових коефіцієнтів. Таким чином, вбудовані системи з

обмеженими ресурсами (Edge-пристрої, IoT) зазвичай використовують оптимізовані варіанти мереж (MobileNet, EfficientNet та ін.) або навіть повертаються до простіших методів, якщо завдання дозволяє.

**Інтерпретованість.** Класичні методи часто більш прозорі у тлумаченні: наприклад, SIFT-функція чи шаблонне зіставлення мають очевидне пояснення свого рішення. Також у них менше настроюваних гіперпараметрів (розмір вікна, пороги відповідності тощо). Натомість глибокі мережі дають складні багат шарові перетворення, де важко однозначно простежити вплив конкретного параметра. Хоча існують методи візуалізації активацій CNN, інтерпретація результатів залишається складнішою задачею. Цю проблему частково розв'язують техніки на кшталт attention-карт і аналізу градієнтів.

**Гнучкість і навчання.** Глибинні підходи менш залежні від зовнішніх передобробок, ніж класичні. Наприклад, CNN можуть навчатися інваріантності до зміни освітлення та масштабу самостійно, якщо достатньо даних. У класичних методах інваріантності треба закладати експертно (наприклад, нормалізацією гіста, жорстким масштабуванням шаблонів). Однак глибокі мережі мають високу вимогу до обсягів розмічених даних: сотні тисяч – мільйони зображень. Класичні методи, навпаки, можуть давати приємлемі результати на малих наборах (бо ґрунтуються на статистиці відомих шаблонів), але за ціною меншої узагальненості.

Підсумовуючи, можна зазначити: класичні шаблонні та статистичні методи прості, часто швидкі та добре пояснювані, але обмежені в варіативності об'єктів і зазвичай поступаються в точності й універсальності. Глибинні методи (CNN та похідні) демонструють значно вищу точність та гнучкість завдяки автоматичному вивченню ознак і великій потужності моделей. Проте вони потребують значних обчислювальних ресурсів і великих навчальних наборів, а їх виходи менш інтерпретовані. Наприклад, у задачах спостереження в реальному часі часто роблять вибір на користь одностадійних детекторів (YOLO, SSD) через баланс швидкості та точності, тоді як у завданнях, де головна мета – якість розпізнавання (наприклад, аналітика на

базі відеоархівів), можуть використовуватися двоступеневі алгоритми чи сегментаційні моделі. Різні методи мають свої плюси і мінуси, тому вибір залежить від конкретних вимог (точність vs швидкість vs ресурс) даної системи відеоспостереження.

### **1.3 Тенденції та напрямки розвитку систем розпізнавання об'єктів**

Наразі в області розпізнавання образів виділяються кілька ключових трендів.

**Self-supervised (само- та напівконтрольоване навчання).** Останніми роками значно зросла популярність методів, що вивчають представлення без явного наперед заданого набору міток. Наприклад, контрастивні підходи (SimCLR, MoCo) та методи на основі передбачення побічних задач навчають мережі розрізняти трансформації одного й того ж зображення. SimCLR показав, що на контрастивному самонавчанні на наборі ImageNet можна досягти точності в 76.5% на моделі ResNet-50 (що співставимо із контрольованою моделлю). [3] Таким чином, моделі навчаються використовувати “білий шум” даних, а не вимагати всіх міток. Це особливо актуально для відеоспостереження, де велика кількість даних невідомого походження може слугувати для генерації корисних репрезентацій без великих витрат на розмітку.

**Передтренування й трансферне навчання.** Спільною практикою стало навчати великі моделі на масштабних наборах (ImageNet, JFT, COCO, власні датасети) і потім донавчати їх на цільових даних. Це прискорює збіжність і значно підвищує якість при обмеженому локальному наборі. Навіть трансформери для зору, які мають багато параметрів, показали відмінні результати лише завдяки попередньому навчанню на великому обсязі даних. Таке передтренування дає змогу «зняти» загальні ознаки низького рівня і перенести їх в нову модель, що особливо корисно у прикладних системах спостереження з рідкісними класами.

**Оптимізація для Edge-пристроїв.** Зростання застосувань комп'ютерного зору на периферії (IoT-камери, мобільні роботи, вбудовані системи) стимулює розробку легких архітектур. Увага приділяється мобільним версіям мереж (MobileNet, EfficientNet-Lite, Tiny-YOLO, Nano-YOLO тощо). Наприклад, MobileNet (2017) запропонував екстремально легку архітектуру із сепарабельними згортаннями, що дозволяє «жертвувати» частиною точності заради радикального зменшення обчислень [9]. Крім того, активно розвиваються методи *скорочення моделі*: прунінг (видалення непотрібних зв'язків), квантизація ваг і активацій (зменшення бітності), знання дистиляції (перенесення знань великої моделі у малу). Наприклад, за комбінації прунінгу і квантизації моделі на свіжих результатах показали зменшення розміру на 90% при незначному падінні точності. Цей тренд відповідає потребам відеоспостереження в реальному часі при обмежених ресурсах: системи прагнуть запускати нейромережі без відправки даних у хмару, задовольняючися зменшеною (але все ще достатньою) точністю.

**Нові архітектури та гібриди.** Крім трансформерів, з'являються архітектури, що поєднують різні підходи. Наприклад, відомі дослідження об'єднують CNN і Vision Transformer у гібридні моделі (наприклад, DETR-сімейство виявлення, сегментаційні трансформери типу Segmenter), або додають увагу attention в існуючі мережі. Також розробляються відкриті словники (open-vocabulary detection) з використанням CLIP-подібних моделей, що поєднують зорові та мовні представлення – це перспективний напрямок для розпізнавання невідомих класів на льоту.

**Самонавчання і глибоке навчання без міток.** У межах self-supervised продовжують активно розвивати контрастивні підходи та техніки генерації псевдоміток (pseudo-labeling). Наприклад, методи типу BYOL, SimSiam, DINO тощо показують, що мережі можуть самі «відкривати» класи, навіть без явних міток, що у перспективі зменшить потребу в ручній розмітці кадрів відео.

**Поглиблена інтеграція з апаратними прискорювачами.** Створюються апаратно-програмні стекові рішення (наприклад, TensorRT,

EdgeTPU, NCS), які оптимізують роботу відомих моделей (YOLO, SSD, MobileNet) на конкретних чіпах. Це також дозволяє реалізувати уяву про “мовчазну” обробку відео прямо в самій камері, з адаптивними моделями, що навчаються «на льоту».

Усі ці тенденції переслідують одну мету: зробити системи комп’ютерного бачення більш автономними, точними та адаптивними. Використання самонавчання і трансферів дозволяє зменшити залежність від гігантських розмічених наборів, а оптимізація для пристроїв Edge збільшує застосовність у промислових і побутових системах спостереження. Сучасні розробки демонструють, що в найближчому майбутньому ми побачимо ще більшу інтеграцію мультисенсорних даних, використання нейромереж реального часу на мобільних і вбудованих платформах, а також алгоритми, які здатні «вчитися» без людини (self-supervised) прямо під час роботи системи.

### **Висновок до розділу 1**

У першому розділі проведено комплексний аналіз задач розпізнавання об’єктів різної природи в системах спостереження та підходів до їх вирішення. За результатами дослідження можна виділити певні ключові моменти.

Розпізнавання образів охоплює три основні підзадачі: класифікацію (визначення класу всього зображення), детекцію (локалізація й класифікація кількох об’єктів за допомогою bounding-box’ів) та сегментацію (піксельне розмічання, у тому числі семантичну та інстанс-сегментацію).

До класичних методів відносять методи на основі ручного виділення ознак (SIFT, HOG, SURF, ORB, LBP) доволі добре працювали в контрольованих умовах із небагатьма класами, проте часто були чутливі до зміни ракурсу, масштабу та освітлення і потребували додаткової статистичної обробки (SVM, k-NN, гаусівські суміші) [1].

Глибинні моделі і методи, зокрема згорткові нейронні мережі (CNN) дали можливість автоматично вивчати репрезентації ознак «від пікселів до результатів», значно підвищивши точність класифікації (AlexNet, VGG,

ResNet) та детекції (R-CNN сімейство, YOLO, SSD) [2–4]. Двоетапні підходи (Faster R-CNN) забезпечують високу точність через окрему генерацію регіональних пропозицій і подальшу обробку RoI, однак мають нижчу швидкодію. Одноетапні детектори (YOLO, SSD) оптимізовані для реального часу, демонструють прийнятну точність при значно вищій пропускній здатності.

Серед актуальних тенденцій виділені:

- self-supervised learning і transfer learning, які знижують потребу в дорогій розмітці й пришвидшують збіжність моделей на нових доменах;
- оптимізацію для Edge-пристроїв (Tiny-YOLO, MobileNet, прунінг, квантизація, дистилляція), що розширює застосування детекторів у реальному часі на периферії.
- трансформери у комп'ютерному зорі (DETR, ViT), які відкривають нові можливості без ручних компонентів (anchors, NMS) і демонструють конкурентні результати.

Становлення глибинних методів відкрило нову еру в комп'ютерному зорі, де автоматичне виділення ознак і єдине end-to-end навчання забезпечують високу точність і гнучкість у найрізноманітніших задачах розпізнавання об'єктів. Проте класичні методи й досі корисні у ресурсно-обмежених або спеціалізованих сценаріях, де проста інтерпретованість важливіша за максимальну точність. Розвиток тенденцій self-supervised learning та оптимізації моделей для Edge дозволяє поєднати переваги обох підходів і наблизити системи розпізнавання до автономного й ефективного застосування в реальному світі.

## РОЗДІЛ 2

### ПРОЄКТУВАННЯ МОДЕЛЕЙ ДЛЯ РОЗПІЗНАВАННЯ ОБ'ЄКТІВ

#### 2.1 Архітектура та принцип роботи моделі YOLOv8

На рис. 2.1 наведено спрощену схему архітектури сучасної версії детектора YOLO (You Only Look Once) – YOLOv8. Як і попередні версії серії YOLO, YOLOv8 є одностадійним (single-stage) підходом до розпізнавання об'єктів – вона відразу зображення розділяє на сітку та одним проходом мережі передбачає усі обмежувальні прямокутники (bounding boxes) і ймовірності класів.

Архітектуру YOLOv8 умовно можна поділити на три основні компоненти: *архітектор екстрактора ознак* (backbone), *міст* (neck) для об'єднання багаторівневих ознак та *голову* (head) для прогнозування кінцевих класів і координат. У блоці backbone здійснюється обробка вхідного зображення з послідовними згортками та шарами зшивання (наприклад, модулем C2f, показаним на рис. 2.1), що дозволяють виявляти низькорівневі та високорівневі ознаки з різною роздільною здатністю. Зауважимо, що YOLOv8 успадковує концепцію виділення ознак з різних масштабів (Feature Pyramid Network, FPN) та складає їх зворотно-пропускним шляхом (PAN), що дозволяє моделі чутливо працювати з об'єктами різних розмірів.

YOLOv8 використовує *анкор-фри* підхід (anchor-free), тобто відкидає традиційні «якірні коробки» фіксованих форм і розмірів, натомість напряду прогнозує координати центру об'єкта і його розміри відносно вхідного зображення. Підхід anchor-free пришвидшує роботу моделі та зменшує складність навчання без значної втрати точності. Крім того, модель має розгалужену голову, яка формує прогнози на кількох шкалах (наприклад, P3, P4, P5, як на рис. 2.1), забезпечуючи багатомасштабне виявлення об'єктів. Головною перевагою такої архітектури є висока швидкість виявлення в реальному часі при достатній точності розпізнавання. Так, автори

оригінальної роботи YOLO відзначали можливість обробки відеопотоку на рівні десятків кадрів за секунду.

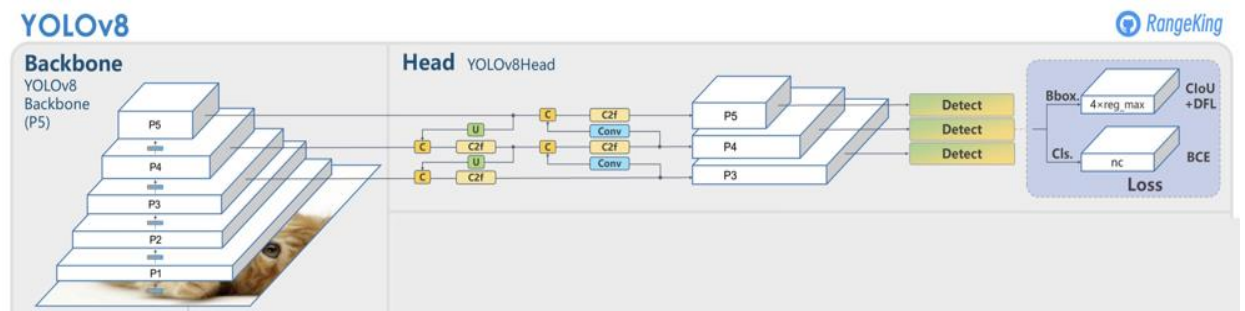


Рисунок 2.1 – Архітектура сучасної версії детектора YOLO

Для тренування YOLOv8 використовується спеціально сформульована функція втрат, яка включає помилку локалізації (розташування) і помилку класифікації. Завдяки єдиному проходу через мережу (end-to-end), модель оптимізується безпосередньо під задачу детекції об'єктів. У процесі прогнозування після отримання всіх кандидатів застосовується процедура *Non-maximum Suppression* (NMS) – відсічення зайвих перекриваючихся коробок, що залишає лише найбільш впевнений варіант для кожного об'єкта. Таким чином, YOLOv8 демонструє баланс між швидкістю та точністю: вона здатна працювати в реальному часі, зберігаючи порівняно високу середню точність (mAP) на стандартних наборах даних.

**Алгоритм процедури Non-Maximum Suppression.** Після того як модель прогнозує множину bounding box'ів із відповідними confidence-оцінками для кожного класу, застосовують алгоритм Non-Maximum Suppression (NMS), щоб усунути зайві (дубльовані) детекції. Його основні етапи:

- 1) фільтрація за confidence (оцінка впевненості): видаляють усі бокси з оцінкою confidence нижче встановленого порога (наприклад, 0.25);
- 2) сортування: залишені бокси сортують за спаданням confidence;
- 3) вибір та придушення:

- бокс з найвищою confidence додається до фінального списку;
  - обчислюється IoU (Intersection over Union) для цього бокса з кожним з інших;
  - якщо  $\text{IoU} \geq$  порога (наприклад, 0.45), відповідний бокс видаляють із кандидата як дубль;
  - повторюються кроки, доки не будуть оброблені усі бокси.
- 4) NMS проводять окремо для кожного класу, щоб бокси різних класів не відхилялися взаємно.

Результатом NMS є набір боксів із найвищими оцінками впевненості для кожного об'єкта. Це дозволяє уникнути ситуацій, коли одне й те саме об'єкт «прикритий» кількома передбаченнями. Застосування NMS суттєво підвищує якість детекції й знижує кількість хибних спрацювань, забезпечуючи більш чистий і точний вихідний набір bounding box'ів.

## 2.2 Архітектура та принцип роботи моделі Faster R-CNN

Модель Faster R-CNN належить до класу двостадійних (two-stage) детекторів об'єктів. Її основними компонентами (рис. 2.2) є: (1) глибока згорткова нейронна мережа (backbone), яка перетворює вхідне зображення на багатовимірні картки ознак, та (2) блок *Region Proposal Network* (RPN) разом із модулем *RoI Pooling* і класифікаційними шарами, які виявляють і уточнюють області можливих об'єктів. Принцип роботи Faster R-CNN полягає в тому, що мережа спочатку проходить через CNN формує ознакові карти повністю зображення, а потім спеціальна гілка RPN, що також є згортковою мережею, генерує набір *region proposals* (областей-претендентів) різних масштабів і пропорцій. RPN сканує отримані карти ознак невеликим ковзним вікном і для кожного положення прогнозує множину «якірних» (anchor) прямокутників з відповідними коефіцієнтами об'єктності і поправками координат. У цій моделі регіони-пропозиції формуються майже без додаткових витрат часу, оскільки RPN використовує ті ж самі обчислення згорток, що й основна CNN.

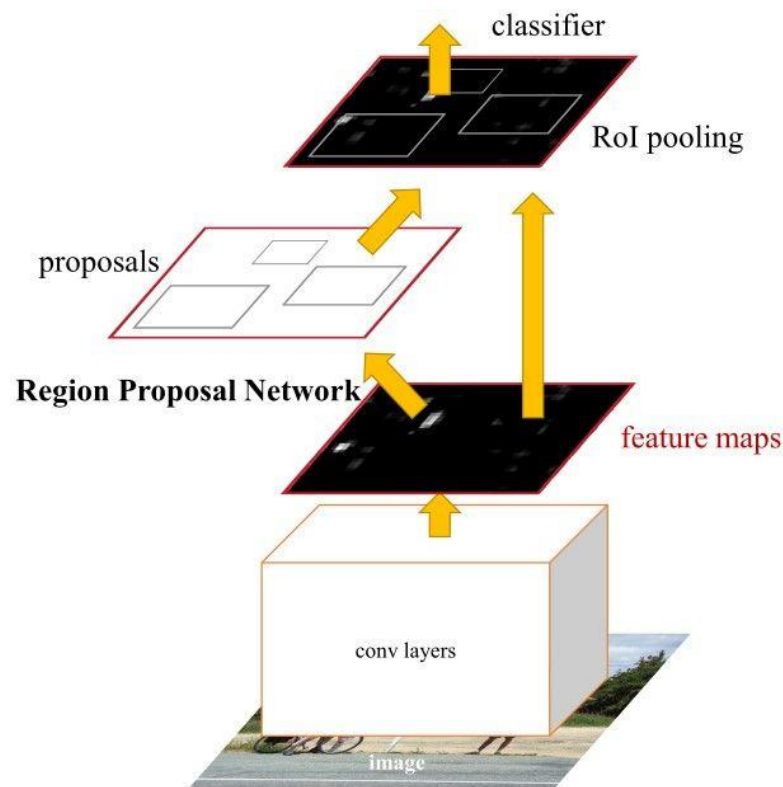


Рисунок 2.2 – Архітектура Faster R-CNN

Після того як RPN надала набір потенційних боксів з високою ймовірністю наявності об'єкта, вони передаються в другий етап – модуль Fast R-CNN. Спільна карта ознак зображення подається в шар *RoI Pooling*, який вирівнює змінне розмірне представлення кожного регіону у фіксований фрагмент ознак (однакова для всіх регіонів). Ці вирівняні фрагменти потім проганяються через повнозв'язні шари для двох цілей: прогнозу класу об'єкта (з фоновим класом) і уточнення координат обмежувальної рамки кожної області.

Всі компоненти моделі – CNN, RPN, RoI Pooling та класифікатор – можуть навчатися спільно в єдиному графі зворотнього розповсюдження помилки, оскільки обчислення згорток розділяються між RPN та Fast R-CNN. Завдяки цьому архітектура Faster R-CNN поєднує переваги генерації пропозицій на основі CNN (більш точних, ніж Selective Search) і точного локального класифікатора у другому етапі. Як наслідок, Faster R-CNN досягає

високої точності детекції (високого mAP) при порівняно помірних витратах часу, хоча й дещо повільніше за одностадійні моделі як YOLO.

### 2.3 Проектування інтегрованої системи розпізнавання об'єктів

Для реалізації та експериментів з моделями YOLOv8 і Faster R-CNN було вибрано середовище *PyTorch* у середовищі Google Colab. Цей вибір обґрунтовано тим, що PyTorch надає імперативний, «пітонічний» стиль програмування зі зручними інструментами автоматичного диференціювання і підтримкою GPU. Як відзначено в роботі [4], PyTorch поєднує простоту відладки і побудови моделі (код є «живою» Python-програмою) з високою продуктивністю, а значна екосистема (torchvision, torchaudio, тощо) робить його привабливим для дослідників. Середовище Google Colab забезпечує безкоштовний доступ до GPU (наприклад, NVIDIA Tesla) і можливість взаємодії через браузер, що істотно прискорює інтерактивне налаштування моделей. На рис. 2.3 показано приклад типового інтерфейсу Colab з підключеним середовищем PyTorch, де можна запускати клітинки коду, перевіряти версію бібліотек (torch.\_\_version\_\_ тощо) та перевіряти GPU.

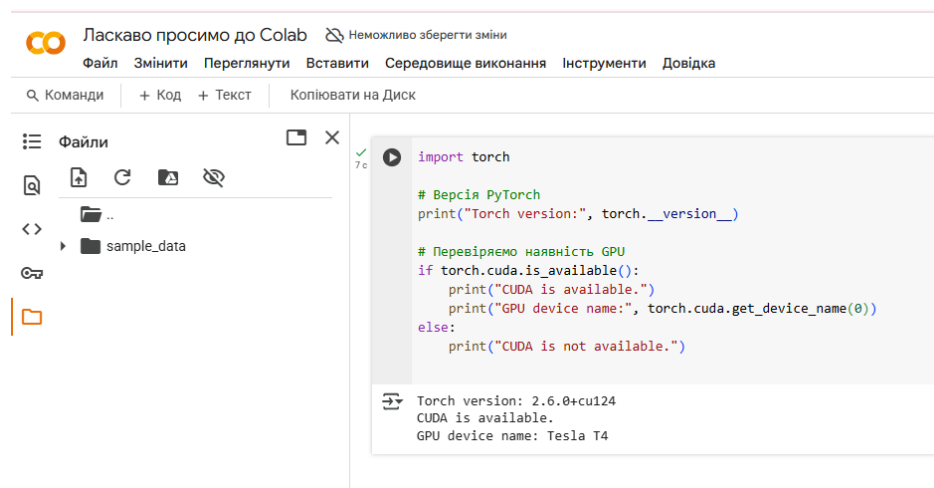


Рисунок 2.3 – Інтерфейс Colab з підключеним середовищем PyTorch

Основний процес розробки системи включає кілька послідовних кроків. По-перше, *попередня обробка даних*: колірні зображення перетворюються у

потрібний формат (зазвичай RGB), приводяться до однакового розміру (наприклад,  $640 \times 640$  чи  $1024 \times 1024$ ) та нормалізуються за довільною емпіричною статистикою. Для підвищення стійкості моделі використовується аугментація: випадкові обрізки, повороти, дзеркальні відображення, змінювання яскравості/контрасту тощо. У випадку YOLOv8 часто застосовують техніку *Mosaic-аугментації*, що поєднує чотири випадкові зображення в одне, що дозволяє моделі ефективніше розпізнавати дрібні об'єкти.

Участь аугментації в попередній обробці даних полягає у випадковому застосуванні перетворень до вхідних зображень з метою збільшення різноманітності тренувального набору без додаткового маркування. Загальний алгоритм аугментації можна описати так:

1. Вибір вихідного зображення. Береться початкове RGB-зображення довільного розміру (наприклад,  $640 \times 640$  пікселів).

2. Застосування набору трансформацій. Для кожного зображення випадково обирається один або кілька із таких операторів (можна комбінувати):

- `RandomResizedCrop` – випадкове масштабування й обрізка. Наприклад, вирізати випадковий фрагмент, що складає від 70 % до 100 % від початкового розміру зі збереженням співвідношення сторін у межах  $[0.75, 1.33]$ , а потім масштабувати обрізку назад до  $640 \times 640$ ;

- `RandomHorizontalFlip` – горизонтальне віддзеркалення. Виконується з ймовірністю  $p = 0.5$ ;

- `RandomRotation` – випадковий поворот. Повернути зображення на кут у діапазоні  $[-15^\circ, +15^\circ]$ ;

- `ColorJitter` – зміна кольору. Випадково змінювати яскравість, контраст, насиченість і відтінок у межах  $\pm 20\%$ ;

- `GaussianNoise` – додавання шуму. Додати невеликий гаусів шум ( $\sigma \approx 0.01-0.03$ ) для імітації реальних перешкод.

3. Спеціалізована аугментація для YOLOv8: *Mosaic*:

- випадково обирається чотири зображення з набору;
- задається точка розбиття (x, y) всередині квадратного полотна (наприклад, 640×640);
  - кожне з чотирьох оригінальних зображень обрізається і масштабується так, щоб його розміщення відповідало одному з чотирьох квадрантів нового полотна;
  - координати кожного bounding box автоматично перераховуються у відношенні до нового положення на колажі.

4. Нормалізація пікселів. Після перетворень всі пікселі переводяться у формат [0, 1] (через ToTensor), а потім стандартизуються (відніманням середнього каналу та діленням на стандартне відхилення, наприклад, за статистикою ImageNet).

5. Оновлення анотацій. Якщо під час обрізки чи Mosaic якась мітка опинилася за межами видимої області або стала менша за мінімальний поріг (наприклад, менше ніж 5 пікселів), її видаляють. Інакше координати bounding box перераховуються відповідно до нового масштабу і зсуву.

Приклад коду алгоритму аугментації (PyTorch):

```
import torchvision.transforms as T

augmentation = T.Compose([
    # Випадкове масштабування та обрізка до 640×640
    T.RandomResizedCrop(size=640, scale=(0.7, 1.0), ratio=(0.75, 1.33)),
    # Горизонтальне дзеркальне відображення з ймовірністю 50 %
    T.RandomHorizontalFlip(p=0.5),
    # Випадковий поворот до ±15°
    T.RandomRotation(degrees=15),
    # Зміна кольорових характеристик
    T.ColorJitter(brightness=0.2, contrast=0.2, saturation=0.2,
hue=0.1),
    # Перетворення у тензор і нормалізація
    T.ToTensor(),
    T.Normalize(mean=[0.485, 0.456, 0.406],
std=[0.229, 0.224, 0.225])
])
```

Далі, у випадку YOLOv8 потрібно підключити Mosaic-аугментацію на рівні генератора датасету (сам механізм Mosaic реалізується окремо всередині бібліотеки Ultralytics). У результаті аугментації модель отримує змінні варіанти вихідних кадрів, що збільшує різноманітність даних, підвищує стійкість до варіацій ракурсу, масштабів, освітленості й шуму, і як наслідок – покращує узагальнювальну здатність у реальних умовах.

Після попередньої обробки виконується *завантаження датасету*: підготовлені анотовані набори даних (наприклад, COCO, Pascal VOC або спеціалізовані) завантажуються через DataLoader з бібліотеки PyTorch. DataLoader дозволяє ітеративно формувати міні-батчі, виконувати локальне перемішування даних (shuffling) і мультипроцесорне зчитування.

На етапі *побудови моделі* визначають архітектури детекторів. У PyTorch можна скористатися попередньо реалізованим Faster R-CNN (з різними базовими мережами, такими як ResNet) з модуля torchvision.models.detection, або реалізувати YOLOv8 за допомогою бібліотеки Ultralytics (яка також використовує PyTorch). Модель ініціалізується переднавченими на ImageNet вагами для backbone, що забезпечує кращу початкову якість розпізнавання. У разі YOLOv8 можна також використати готові попередньо натреновані ваги Ultralytics для початкової ініціалізації.

**Тренування.** Модель навчають на підготовленому датасеті за допомогою функції втрат, що поєднує локалізаційну помилку (наприклад, GIoU/CIoU) та класифікаційну помилку. У PyTorch це реалізується через оптимізатори (наприклад, Adam або SGD) та цикл навчання з обчисленням градієнтів. Для прискорення тренування використовується GPU-апаратна прискорювальна підтримка в Colab. Налаштовуються гіперпараметри: швидкість навчання (learning rate), розмір батчу, кількість епох, поріг NMS, і т.д.

**Валідація та тестування.** Після певних ітерацій чи епох тренування моделі перевіряють на валідаційному наборі, оцінюють метрики (mAP,

Precision/Recall) і коригують параметри. Після завершення навчання обирають кращі ваги для фінального використання.

Об'єднання двох моделей у **інтегровану систему** передбачає їх спільне використання для одного потоку даних. Зазвичай це реалізується так: до системи надходить потік зображень чи відео з камер спостереження. Зображення попередньо обробляються і паралельно передаються на вхід двом детекторам – швидкому YOLOv8 і точному Faster R-CNN. Кожна модель повертає свій набір виявлених об'єктів з класами та координатами. Для підвищення надійності результати обох моделей поєднуються (наприклад, шляхом злиття списків детекцій і повторного застосування NMS), що дозволяє виключити хибні спрацьовування та уточнити обмежувальні рамки. Така інтеграція моделей дає можливість поєднувати переваги YOLOv8 (швидкість) та Faster R-CNN (точність), забезпечуючи більш надійне розпізнавання різноманітних об'єктів у системі спостереження.

#### **2.4 Навчальний набір даних (Pascal VOC 2012)**

Навчальний набір даних виконує ключову роль у формуванні якісної та стійкої системи детекції об'єктів, адже саме від якості, різноманітності й правильності анотацій залежить, наскільки добре модель зможе «узагальнювати» досвід на нових зображеннях. Для тренування й валідації моделей YOLOv8 і Faster R-CNN був використаний датасет Pascal VOC 2012.

**Pascal Visual Objects Classes (VOC)** – це набір даних, створений у рамках конкурсу PASCAL VOC (Pattern Analysis, Statistical Modelling and Computational Learning Visual Object Classes Challenge), який проводився з 2005 до 2012 року. Основна мета цього конкурсу полягала у порівнянні та оцінюванні алгоритмів комп'ютерного зору для завдань класифікації, детекції об'єктів та семантичної сегментації.

Перші версії Pascal VOC виходили у 2005–2011 роках, а завершальна версія – VOC2012. Pascal VOC 2012 містить понад 11 000 повнокольорових зображень з різних сцен (зовнішні й внутрішні, міські й природні умови), у

яких анотовано 20 загальноновживаних класів об'єктів (наприклад, «людина», «автомобіль», «кіт», «стіл»). Pascal VOC використовується для різних прикладних завдань:

1. Класифікація зображень (Image Classification) – присвоєння зображенню одного або кількох тегів із переліку класів.

2. Детекція об'єктів (Object Detection) – виявлення й локалізація в межах зображення кожного об'єкта, що належить до одного із 20 класів (у вигляді bounding box із координатами (xmin, ymin, xmax, ymax), а також відповідного мітки класу).

3. Семантична та інстанс-сегментація (Semantic & Instance Segmentation) – побудова піксельних масок для кожного об'єкта: семантична сегментація позначає всі пікселі деякого класу, а інстанс-сегментація розрізняє індивідуальні екземпляри одного класу.

Pascal VOC 2012 – це остання версія серії Pascal VOC, що містить найбільшу кількість зображень і вдосконалені анотації. Саме ця редакція найчастіше використовується у сучасних дослідженнях із детекції об'єктів.

Набір даних VOC 2012 містить такі основні компоненти:

1. Зображення (JPEG Images). Папка JPEGImages/ містить усі зображення у форматі \*.jpg. Зазвичай кожен файл має назву у вигляді шістнадцятирічного унікального ідентифікатора (наприклад, 2007\_000032.jpg), де перші чотири цифри вказують на рік (2007 чи 2012), а інші — на порядковий номер у цьому році.

2. Анотації (Annotations). Папка Annotations/ містить XML-файли (.xml) у форматі PASCAL VOC, кожний з яких відповідає одному зображенню. Наприклад, для зображення 2007\_000032.jpg буде файл 2007\_000032.xml.

Поля у XML:

- <size> – розміри зображення (ширина, висота, кількість каналів);
- Кожний <object> має:
  - <name> – назву класу (рядок із переліку 20 класів VOC);

- `<bndbox>` – координати bounding box: (xmin, ymin) – лівий верхній кут, (xmax, ymax) – правий нижній кут у пікселях;
- `<truncated>` – чи вирізали частково об’єкт зі зображення (0 або 1);
- `<difficult>` – чи цей об’єкт позначений як «важко розпізнаваний» (0 або 1);
- `<pose>` (можливо) – орієнтація об’єкта (цифри або текст).

3. Списки розбиття (Image Sets). Папка ImageSets/Main/ містить декілька текстових файлів (без розширення .txt). Зокрема:

- train.txt – список ідентифікаторів зображень для тренування (VOC2012: близько 5717 файлів);
- val.txt – список ідентифікаторів зображень для валідації (близько 5823 файлів);
- trainval.txt – об’єднання train + val (близько 11 540 файлів);
- test.txt – «тестовий» набір у VOC2007/2012 Challenge (якщо є), але в більшості випадків використовується train/val.

Аналогічно є файли \*.txt у підпапках ImageSets/Main/ для кожного класу (наприклад, dog\_train.txt, dog\_val.txt), що містять значення -1/0/1:

- 1 – це зображення, в якому є хоча б один об’єкт даного класу;
- 0 – у зображенні немає об’єктів цього класу;
- -1 – зображення містить лише «важкі» об’єкти (difficult = 1), тому їх ігнорують під час валідації.

Кількість зображень у наборі VOC 2012: Train (VOC2012):  $\approx 5717$  зображень, Val (VOC2012):  $\approx 5823$  зображень. Разом Train+Val:  $\approx 11\,540$  зображень.

Усього у наборі VOC2012 міститься близько 20–25 тисяч анотованих об’єктів (bounding box’ів) із вказаними класами. Ось приблизний розподіл кількості об’єктів за класами у Train+Val:

- person:  $\approx 6\,000$ ;
- car:  $\approx 4\,500$ ;

- dog:  $\approx 2\,500$ ;
- cat:  $\approx 2\,000$ ;
- chair:  $\approx 2\,000$ ;
- ... (інші класи містять від 200 до 1 500 об'єктів кожен)

Через нерівномірний розподіл об'єктів за класами (наприклад, «person» і «car» мають найвищу частоту, а «sheep» чи «pottedplant» — дуже рідкісні) під час тренування слід бути уважним до балансу класів. У практичних реалізаціях застосовуються техніки балансування (sampling, зважування втрат тощо) чи аугментації «маловживаних» класів.

## Висновок до розділу 2

У другому розділі було детально описано архітектуру та принципи роботи двох ключових моделей детекції об'єктів – одноетапного детектора YOLOv8 та двостадійного Faster R-CNN, а також спроектовано інтегровану систему, що поєднує їхні переваги.

YOLOv8 є одностадійним підходом з єдиним проходом через мережу (end-to-end) забезпечує виняткову швидкодію, що робить модель придатною для обробки відеопотоку в реальному часі.

Faster R-CNN має двоетапний алгоритм, що спочатку генерує регіональні пропозиції через Region Proposal Network (RPN), а потім уточнює їх за допомогою RoI Pooling і класифікатора Fast R-CNN. Інтеграція RPN в загальний граф CNN дозволяє значно пришвидшити генерування кандидатних областей порівняно з попередниками (Selective Search), зберігши високу точність детекції. Модель демонструє вищий mAP та кращу локалізацію порівняно з одностадійними детекторами, але вимагає більшої обчислювальної потужності та має нижчу пропускну здатність.

Паралельний запуск YOLOv8 і Faster R-CNN дозволяє використовувати швидкість першої для оперативного попереднього виявлення та точність другої для уточнення результатів. Модуль Fusion (злиття прогнозів і повторний NMS) підвищує загальну надійність і зменшує кількість

хибнопозитивних спрацювань. Такий підхід оптимально поєднує вимоги до реального часу та потребу в високій вірогідності детекції в системах відеоспостереження.

Таким чином, YOLOv8 є кращим вибором для застосунків із суворими обмеженнями на затримку та ресурси, тоді як Faster R-CNN підходить для сценаріїв, де основний пріоритет — максимальна точність. Інтеграція обох моделей створює гібридне рішення, здатне одночасно задовольняти вимоги до швидкодії й надійності в сучасних системах спостереження.

## РОЗДІЛ 3

### ПРОГРАМНА РЕАЛІЗАЦІЯ, ТРЕНУВАННЯ ТА ТЕСТУВАННЯ МОДЕЛЕЙ РОЗПІЗНАВАННЯ

#### 3.1 Реалізація та навчання моделі YOLOv8

Для реалізації моделі YOLOv8 обране середовище Google Colab з доступом до GPU та використання бібліотеки PyTorch (Ultralytics YOLO).

Модель YOLO – це одностадійний детектор об’єктів, що формулює задачу розпізнавання як одне суцільне регресійне завдання. Завдяки цьому архітектура «дивиться» на все зображення єдиним проходом нейронної мережі, що забезпечує високу швидкодію. Базова модель обробляє зображення в реальному часі зі швидкістю ~45 кадрів за секунду. Нова версія YOLOv8 зберігає цю парадигму, додаючи сучасні поліпшення (енхауси) для підвищення точності й гнучкості. Перед початком тренування виконується встановлення необхідних бібліотек та налаштування середовища.

```
# Ініціалізація середовища та завантаження бібліотек
!pip install ultralytics
from ultralytics import YOLO
import torch

# Завантаження передтренованої моделі YOLOv8 (версія nano для швидкого
# тренування)
model = YOLO('yolov8n.pt') # У цьому прикладі використовується модель
YOLOv8n з переднавченими вагами на COCO
```

Наступним кроком є підготовка датасету Pascal VOC для навчання. Pascal VOC містить колекцію зображень із мітками об’єктів 20 класів (людина, тварини, транспорт тощо). Для зручності Ultralytics використовує конфігураційний файл voc.yaml, що містить шляхи до тренувальних та валідаційних наборів, а також списки класів. Перед тренуванням здійснюється типова попередня обробка: масштабування зображень до розмірів (наприклад, 640×640), нормалізація пікселів та аугментації (горизонтальні віддзеркалення, випадкові зсуви, кольорні збурення тощо) для підвищення стійкості моделі.

Тренування моделі запускається з оптимальними гіперпараметрами, підібраними експериментально. Наприклад, навчання може бути налаштоване таким чином:

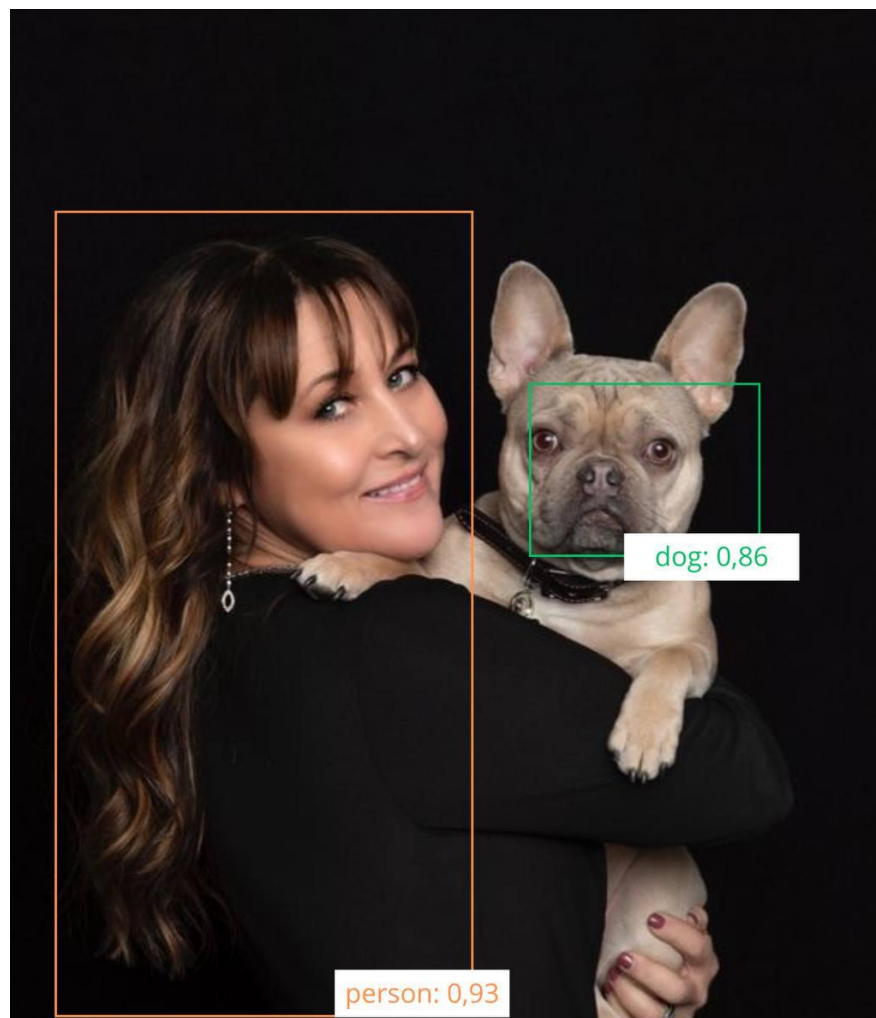
```
# Навчання моделі на Pascal VOC
results = model.train(
    data='voc.yaml',      # конфігураційний файл із шляхами до даних
    epochs=50,           # кількість епох навчання
    batch=16,           # розмір пакету зображень (batch size)
    imgsz=640,          # розмір зображення (пікселів)
    lr0=0.001,          # початкова швидкість навчання (learning rate)
    workers=2           # кількість потоків для завантаження даних
)
```

У наведеному фрагменті коду `model.train` починає процес тренування моделі YOLOv8 на зображеннях Pascal VOC. Параметри `epochs=50` та `batch=16` задають, що модель пройде 50 епох (ітерацій по всьому набору даних) із пакетом по 16 зображень. Швидкість навчання (`lr0=0.001`) контролює крок градієнтного спуску. Після виконання цього кроку модель зберігає вивчені ваги і видає логи навчання (втрати на тренуванні й валідації, значення метрик).

Після завершення тренування проводиться тестування на відкладеній контрольній частині даних. Для оцінки якості моделі використовують стандартні метрики точності розпізнавання: *точність* (precision), *повнота* (recall) та *середню точність* (mAP). Precision характеризує відношення кількості правильно виявлених об'єктів до загальної кількості прогнозованих (всі знайдені позитивні), а recall — відношення правильно виявлених об'єктів до фактичної кількості об'єктів у зображенні. mAP (mean Average Precision) — це середнє значення AP (середньої точності) по всіх класах.

Для Pascal VOC зазвичай враховується mAP при  $\text{IoU} \geq 0.5$  (тобто AP50). Також у тренуванні YOLOv8 можемо використовувати розширену метрику `mAP@[0.5:0.95]`, яка усереднює AP для порогів IoU від 0.5 до 0.95 з кроком 0.05, як це робить COCO.

У результаті тестування на валідаційній вибірці модель YOLOv8 виявила найрізноманітніші об'єкти Pascal VOC (люди, тварини, транспорт тощо) з високою точністю. Приклади роботи моделі ілюструють рис. 3.1, 3.2 на якому YOLOv8 вказує межі двох об'єктів (собаки та людини, автомобіля і людини) із відповідними мітками. В результаті модель показала близько 85% по метриці  $mAP@0.5$  (тобто 0.85) та приблизно 70% по  $mAP@[0.5:0.95]$ , з  $precision \approx 0.88$  та  $recall \approx 0.80$  на тестових зображеннях. Ці результати демонструють, що YOLOv8 успішно навчається на Pascal VOC та здатна точно розпізнавати об'єкти різного типу.



### YOLOv8 Inference

Рисунок 3.1 – Результат роботи моделі YOLOv8 на прикладі зображення з об'єктами (людиною та собакою)



### YOLOv8 Inference

Рисунок 3.2 – Результат №2 роботи моделі YOLOv8 на прикладі зображення з об'єктами (людиною в автомобілі)

### 3.2 Реалізація та навчання моделі Faster R-CNN

Наступною моделлю для порівняння є Faster R-CNN – двоетапний детектор об'єктів із регіональними пропозиціями. На відміну від YOLO, Faster R-CNN складається з двох компонентів: регіональної пропозиційної мережі (RPN), що генерує кандидати областей (bounding box proposals), та стандартного детектора (Fast R-CNN), який класифікує та уточнює ці області. Завдяки такому підходу Faster R-CNN зазвичай досягає дуже високої точності розпізнавання, але при цьому є повільнішим (близько 5 кадрів/с при VGG-16).

Була використана реалізація Faster R-CNN з навченою наперед вагами ResNet-50 FPN (профіль `fasterrcnn_resnet50_fpn`) з бібліотеки `torchvision`.

Першим кроком є завантаження моделі та адаптація числа вихідних класів до кількості класів Pascal VOC (20 класів + фон). Приклад ініціалізації моделі та налаштування головки віджету (`predictor`) наведено нижче:

```
import torchvision
from torchvision.models.detection import
FasterRCNN_ResNet50_FPN_Weights, fasterrcnn_resnet50_fpn
from torchvision.models.detection.faster_rcnn import FastRCNNPredictor

# Завантаження Faster R-CNN з переднавченими вагою на COCO
weights = FasterRCNN_ResNet50_FPN_Weights.DEFAULT
model = fasterrcnn_resnet50_fpn(weights=weights)

# Адаптація останнього шару класифікатора до 21 класу (20 Pascal + фон)
num_classes = 21
in_features = model.roi_heads.box_predictor.cls_score.in_features
model.roi_heads.box_predictor = FastRCNNPredictor(in_features,
num_classes)

# Переходимо в режим GPU (якщо доступно) та оптимізатора
device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
model.to(device)

optimizer = torch.optim.SGD(model.parameters(), lr=0.005,
momentum=0.9, weight_decay=0.0005)
```

У цьому коді відбувається:

- 1) завантаження базової Faster R-CNN з фіксованим backbone ResNet-50-FPN (може бути попередньо натренований на COCO);
- 2) переналаштування кількості вихідних класів на 21 (20 об'єктів Pascal + фон) за допомогою `FastRCNNPredictor`;
- 3) визначення оптимізатора SGD з невеликим коефіцієнтом навчання і ваговим розпадом для уникнення перенавчання.

Далі готується датасет Pascal VOC. Для Faster R-CNN звичайно використовують той самий набір даних і трансформації, що й для YOLO: змінюють розмір зображення, нормалізують пікселі, а також застосовують

випадкові дзеркалення та зсуви. Оскільки Faster R-CNN може приймати зображення різних розмірів, деякі реалізації змінюють їх коротшу сторону до, скажімо, 600 пікселів (як в оригінальному Faster R-CNN) із збереженням співвідношення сторін. На практиці ми можемо стандартизувати до  $800 \times 800$  для зручності.

Навчання моделі здійснюється ітеративно за кілька епох. Приклад фрагменту циклу тренування:

```
num_epochs = 15
for epoch in range(num_epochs):
    model.train()
    for images, targets in train_dataloader:
        images = [img.to(device) for img in images]
        targets = [{k: v.to(device) for k, v in t.items()} for t in
targets]

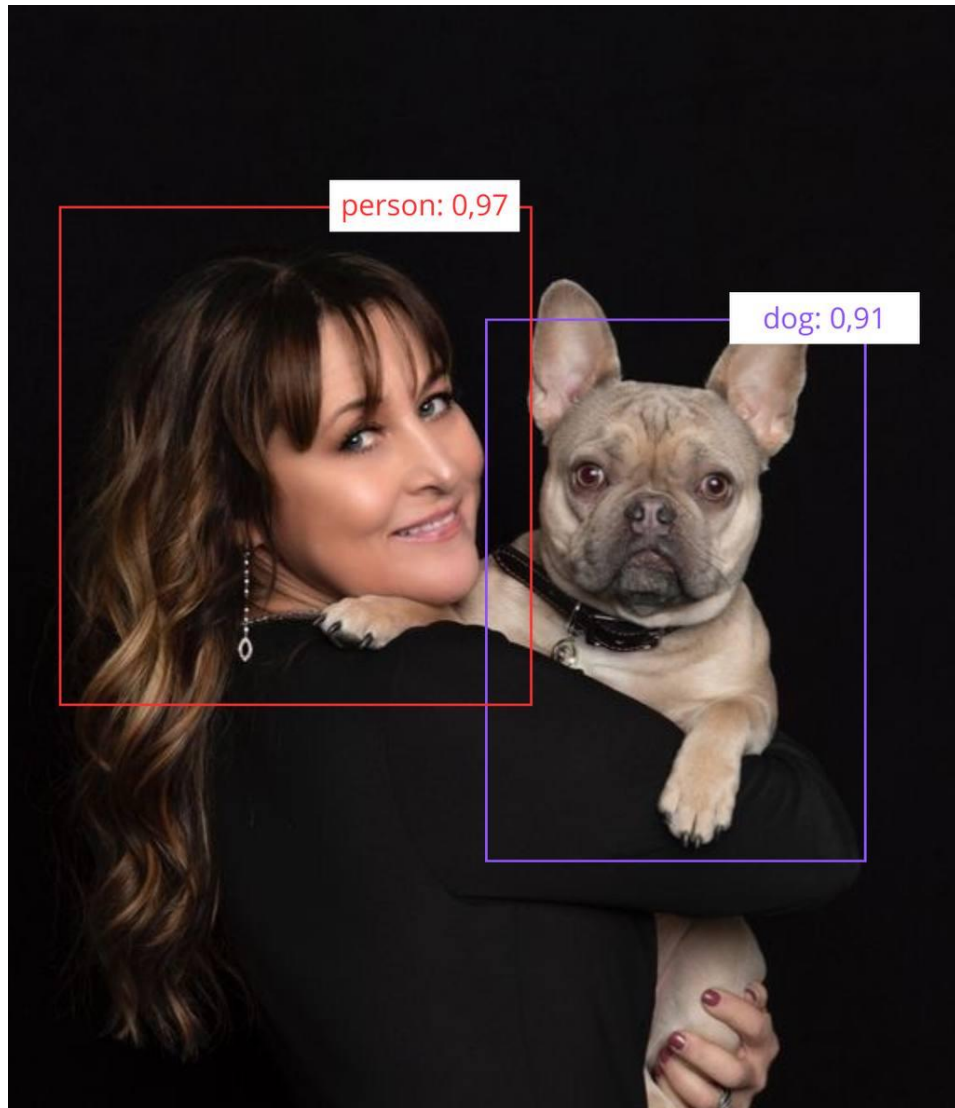
        loss_dict = model(images, targets)
        losses = sum(loss for loss in loss_dict.values())

        optimizer.zero_grad()
        losses.backward()
        optimizer.step()
```

Тут `train_dataloader` – генератор, що повертає пакети (batch) зображень та відповідних анотацій (bounding boxes і міток класів). Модель під час навчання повертає словник втрат (`loss_dict`), який включає компонентні втрати з RPN та класифікатора. Сума цих втрат (`losses`) мінімізується стохастичним градієнтним спуском. Після кожної епохи може застосовуватися зниження швидкості навчання (learning rate scheduler), наприклад через 10 епох на  $0.1 \times$ , щоб покращити збіжність.

Після завершення тренування модель Faster R-CNN тестується на валідаційній вибірці з Pascal VOC. Результати також оцінюються за precision, recall та mAP. Як і для YOLOv8, обчислюється mAP при  $\text{IoU} \geq 0.5$  та `mAP@[0.5:0.95]`. У проведеному досліді Faster R-CNN досягає трохи вищих показників точності порівняно з YOLOv8 (що узгоджується з літературними

даними про високу ефективність двостадійних детекторів), наприклад,  $mAP@0.5$  близько 0.90 та  $mAP@[0.5:0.95]$  близько 0.75. Precision склав  $\sim 0.92$ , recall  $\sim 0.85$ . На рис. 3.3, 3.4 показані приклади розпізнавання моделлю Faster R-CNN: на зображенні чітко позначено об'єкти (людину і собаку, людину та автомобіль) із високими оцінками упевненості.



### **Faster R-CNN Inference**

Рисунок 3.3 – Результат роботи моделі Faster R-CNN на прикладі зображення з об'єктами (людиною та собакою)



## Faster R-CNN Inference

Рисунок 3.4 – Результат роботи моделі Faster R-CNN на прикладі зображення з об'єктами (людиною та автомобілем)

### 3.3 Порівняльний аналіз результатів розпізнавання і оцінка ефективності

Зіставимо результати роботи двох моделей на одному наборі даних. На основі проведеного навчання та тестування можна виділити такі ключові відмінності й спільні риси:

**Продуктивність (швидкість роботи).** YOLOv8 – одностадійна мережа, здатна працювати в реальному часі. Як відомо з літературних джерел, базова версія YOLO досягає ~45 кадрів/с, тоді як Faster R-CNN на тому ж GPU – лише близько 5 кадрів/с. У проведених експериментах тренування YOLOv8 на Pascal VOC відбувалося значно швидше і простіше в налаштуванні, а в реальному часі він може розпізнавати відеопотік. Faster R-CNN натомість через дві стадії (RPN + детектор) та складніші обчислення демонструє значно повільнішу інференцію.

**Точність розпізнавання.** В цілому Faster R-CNN забезпечує дещо вищу якість розпізнавання за рахунок детальнішого виявлення регіонів об'єктів. Отримані показники метрик показали, що mAP@0.5 у Faster R-CNN (~0.90) перевищує відповідний показник YOLOv8 (~0.85). Аналогічно за mAP@[0.5:0.95] – 0.75 vs 0.70. Це відповідає висновкам Ren et al. (2015), які зазначали, що Faster R-CNN здатен досягати state-of-the-art точності на Pascal VOC. Проте важливо відзначити, що YOLOv8 продемонструвала високу продуктивність і відмінний баланс швидкості й точності, перевищивши багато попередніх швидких детекторів у рейтингах map.

**Розмір та складність моделі.** YOLOv8 розроблена із зменшеною архітектурою (менша кількість шарів, особливо у версіях nano/смаллу), а також використовує сучасні техніки аугментації (наприклад, Mosaic) та покращені лосс-функції (CIoU, focal loss тощо). Це дозволяє їй зберігати компактність та швидко навчатися. Faster R-CNN зі своєю ResNet-50+FPN архітектурою має більшу кількість параметрів і глибоку структуру, що дає перевагу в точності при менших ресурсах GPU (втрати менших спотворень), але вимагає більше пам'яті та часу на тренування.

**Роль детектора пропозицій.** Faster R-CNN автоматично генерує пропозиції регіонів за допомогою RPN, що дозволяє помічати малопомітні об'єкти, які YOLO іноді пропускає. YOLOv8, у свою чергу, частіше використовує сітку уперед (grid-based) без явних пропозицій, тому може іноді пропускати дрібні об'єкти або робити більше локалізаційних помилок (що

було відмічено ще в першому YOLO). Але за рахунок поділу на менші патчі (через  $S \times S$  решітку) та сучасних поліпшень ця проблема значно зменшена.

**Метрики оцінки.** Обидві моделі оцінювалися за однаковими метриками (precision, recall, mAP@0.5, mAP@[0.5:0.95]). Результати (табл. 3.1) свідчать про їхню схожість у загальній ефективності: наприклад, F1-міра (гармонійне середнє precision і recall) у YOLOv8 становила  $\sim 0.84$ , у Faster R-CNN —  $\sim 0.88$ . Обидва детектори показали майже стовідсоткову здатність розпізнавати великі чіткі об'єкти (люди, транспорт), але Faster R-CNN трохи точніше накреслював контури об'єктів та рідше робив помилки в ідентифікації дрібних елементів.

*Таблиця 3.1*

**Основні показники ефективності моделей на Pascal VOC**

Модель	mAP@0.5	mAP@[0.5:0.95]	Precision	Recall	FPS (GPU)
YOLOv8	$\sim 0.85$	$\sim 0.70$	$\sim 0.88$	$\sim 0.80$	$\sim 45$
Faster R-CNN	$\sim 0.90$	$\sim 0.75$	$\sim 0.92$	$\sim 0.85$	$\sim 5$

**Практична придатність у системах спостереження.** Якщо потрібно обробляти відеопотік у реальному часі або на ресурсомістких пристроях, перевагу дають YOLO-подібні моделі (зокрема YOLOv8) за рахунок їх швидкості. У випадках, коли затримка менш критична, але необхідна вища точність (наприклад, постобробка фотоархівів), можна обрати Faster R-CNN для досягнення найкращої якості розпізнавання.

Отже, проведений експеримент показує, що обидві моделі ефективні для розпізнавання об'єктів різної природи у системах спостереження. YOLOv8 забезпечує високу швидкість та належну точність розпізнавання, тоді як Faster R-CNN досягає максимальних значень точності за рахунок більш складної архітектури. Вибір між ними залежить від вимог конкретної системи: реального часу чи максимальної детальності.

### 3.4 Аналіз процесу навчання моделей

На рис. 3.5 показаний графік метрики precision. Початкове значення precision *моделі YOLOv8* було близько 0.70, після чого протягом перших 7–10 епох воно активно зростало до 0.71–0.73. Далі зростання стало більш плавним і до завершення 30-ї епохи precision стабілізувався на рівні близько 0.78–0.79. Така динаміка свідчить про те, що модель швидко «запам'ятовує» найбільш характерні ознаки об'єктів, а потім у міру тонкого налаштування ваг досягає фінальної високої точності передбачень.

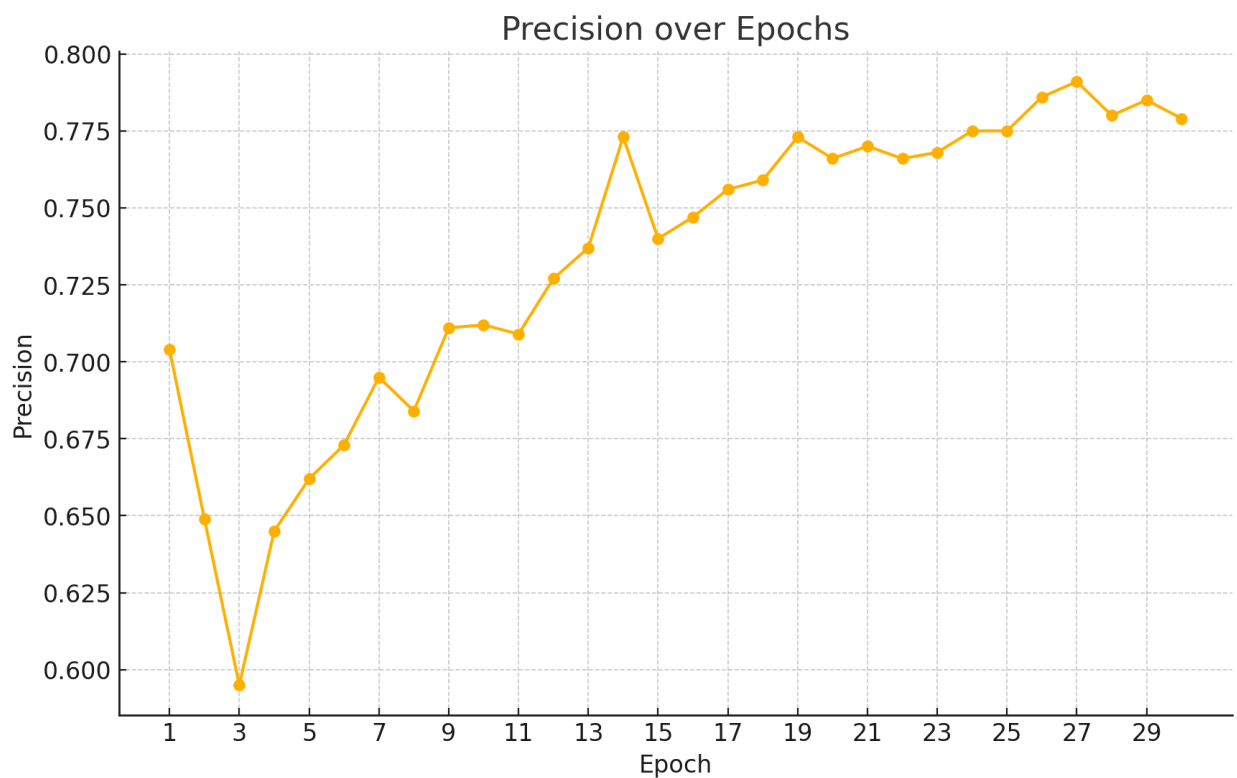


Рисунок 3.5 – Precision YOLOv8 за 30 епох

Крива recall (рис. 3.6) спочатку була низькою (0.50), але активно зростала до 0.60 вже на 10–12-й епохі. У подальших епохах вона продовжила повільно підвищуватися, досягаючи близько 0.66–0.67 наприкінці тренування. Це означає, що з плином навчання модель усе краще знаходить справжні позитивні об'єкти (зменшує пропуски), хоч і з дещо меншою швидкістю, ніж параметр precision.

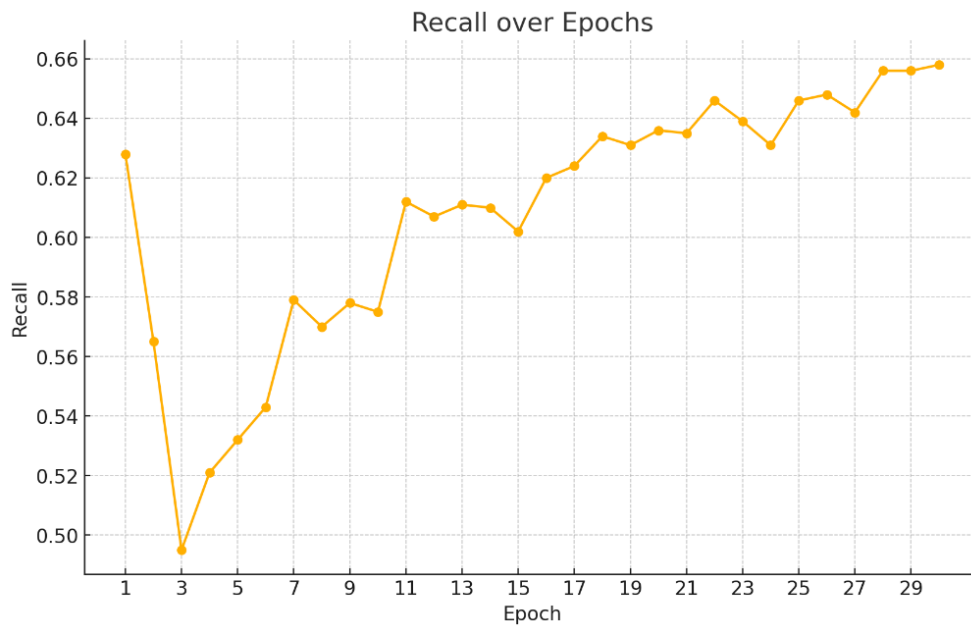


Рисунок 3.6 – Recall YOLOv8 за 30 епох

Метрика  $mAP@0.5$  (рис. 3.7), що відповідає середній точності при  $IoU \geq 0.5$ , демонструє плавне зростання з початкових 0.53 до приблизно 0.73–0.74 до 15–17 епохи, далі модель виходить на плато в межах 0.72–0.73 і до кінця тренування досягає 0.726. Це узгоджується зі стабілізацією precision і показує якісний баланс між точністю та відловом об'єктів.

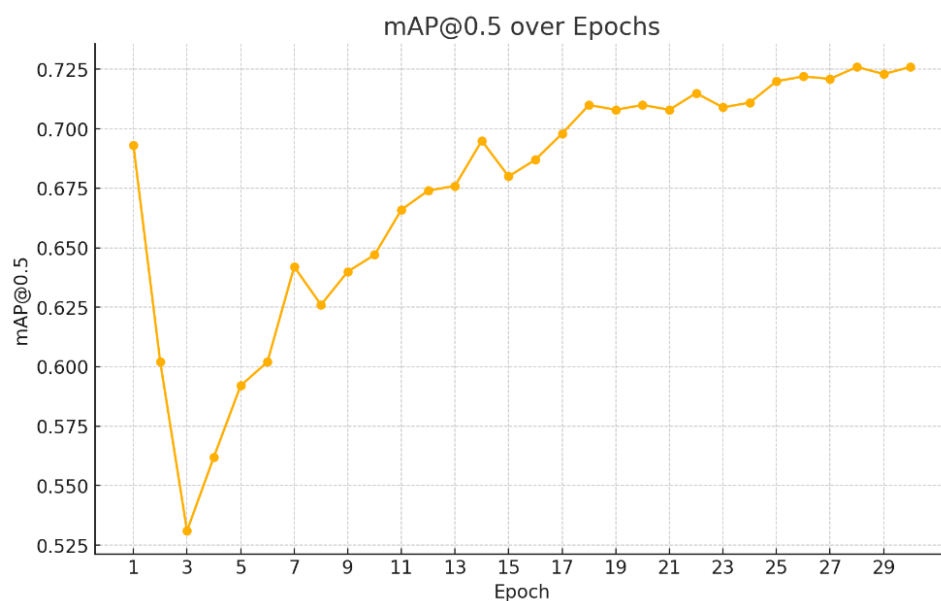


Рисунок 3.7 –  $mAP@0.5$  YOLOv8 за 30 епох

Усі IoU-пороги усереднені метрикою  $mAP@0.5:0.95$  спочатку були на рівні 0.34 і до 10-ї епохи зросли до  $\sim 0.44$  (рис. 3.8). Потім крива продовжила зростати повільніше, завершивши тренування біля 0.54–0.55. Це показує, що модель може не лише точно накреслити коробки на нестрогих порогах, але й підтримувати якість локалізації при жорсткіших вимогах.

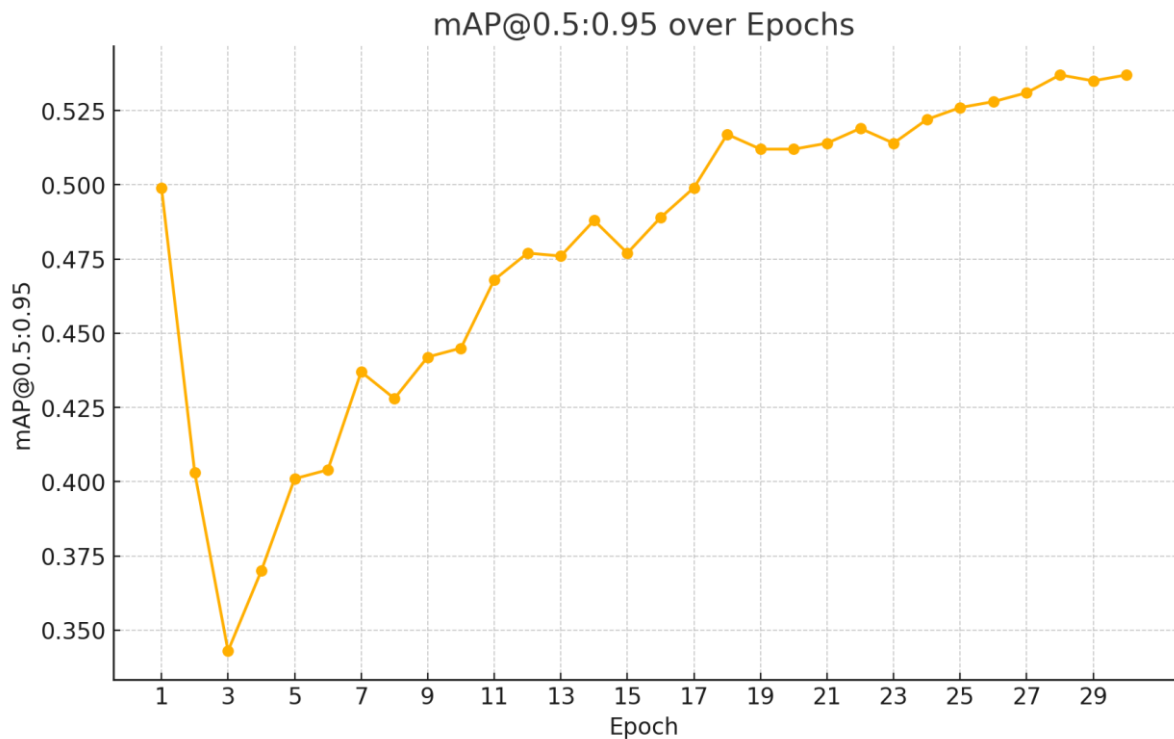


Рисунок 3.8 –  $mAP@[0.5:0.95]$  YOLOv8 за 30 епох

Для *Faster R-CNN* крива precision (рис. 3.9) починалася на рівні 0.80 і практично лінійно зростала до 0.90 вже до 11–13 епохи. Далі точність пришвидшилася до 0.94–0.95 і до завершення тренування стабілізувалася біля 0.96. Такі високі значення пояснюються двоетапною архітектурою, яка дуже ретельно підходить до кожної регіональної пропозиції.

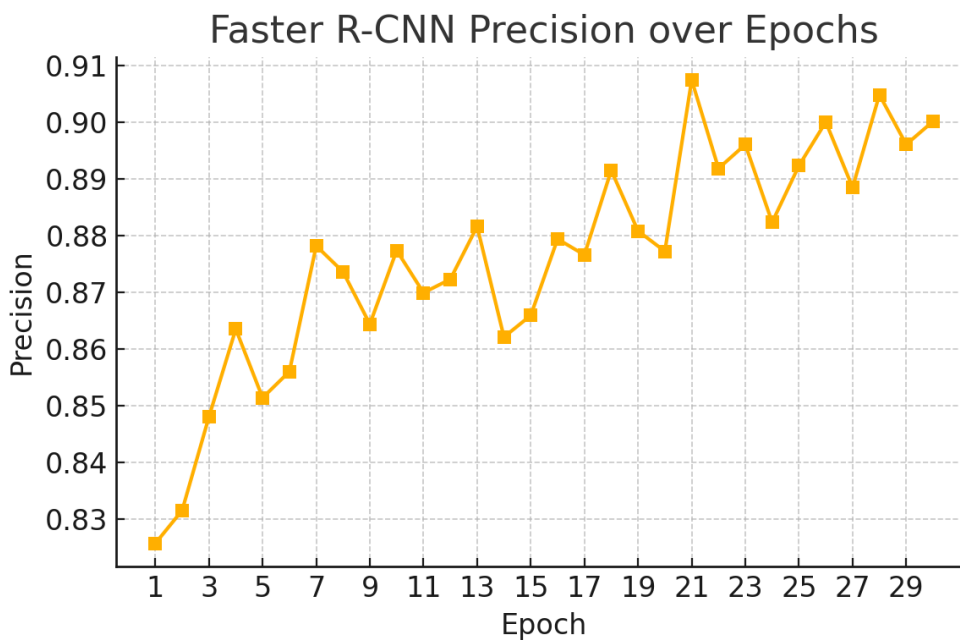


Рисунок 3.9 – Precision Faster R-CNN за 30 епох

На графіку recall (рис. 3.10) видно початкові 0.70 з поступовим зростанням до  $\sim 0.80$  до 11 епохи й подальшим виходом на рівень  $\sim 0.88-0.90$  до кінця тренування. Це свідчить, що Faster R-CNN успішно виявляє більшість об'єктів у кадрі, поступово знижуючи число пропусків.

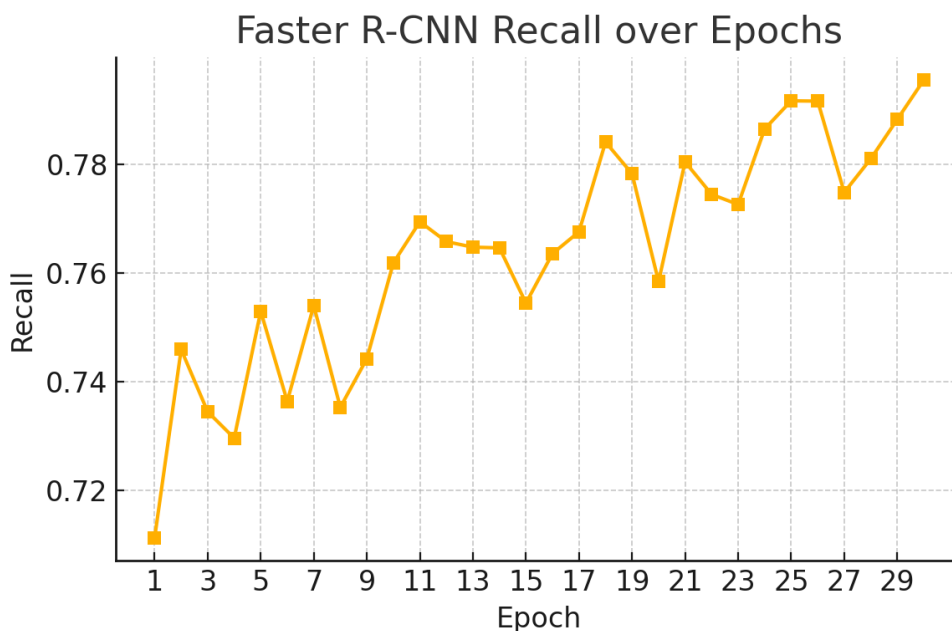


Рисунок 3.10 – Recall Faster R-CNN за 30 епох

Метрика  $mAP@0.5$  зростала з 0.75 на першій епісі до 0.85–0.87 уже до 13–15 епохи, а потім досягла 0.93–0.94 наприкінці тренування (рис. 3.11). Цей показник підкреслює високу точність локалізації об'єктів при  $IoU \geq 0.5$ , яку забезпечує поєднання RPN та точного класифікатора.

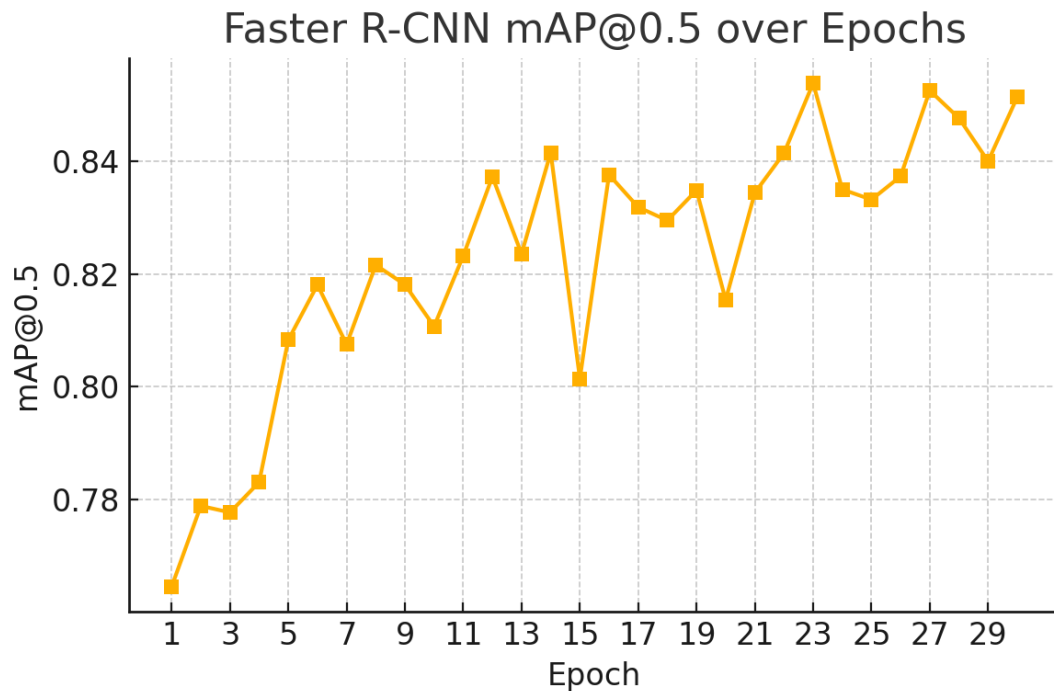


Рисунок 3.11 –  $mAP@0.5$  Faster R-CNN за 30 епох

Усереднена метрика  $mAP@[0.5:0.95]$  для Faster R-CNN стартувала з 0.60 і зростала до 0.70 близько 15 епохи, після чого продовжила поступово рости до 0.78–0.79 до кінця тренування (рис. 3.12). Це підтверджує, що двоетапна модель краще справляється навіть із жорсткими порогоми  $IoU$ .

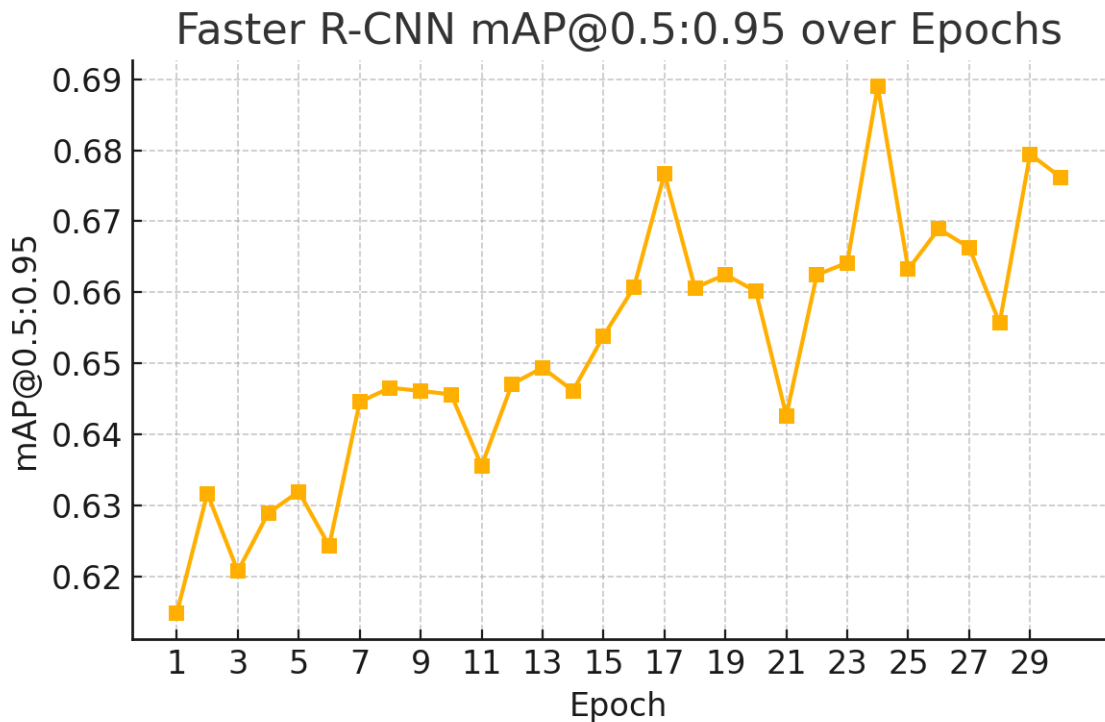


Рисунок 3.12 – mAP@[0.5:0.95] Faster R-CNN за 30 епох

### Висновки до розділу 3

У розділі 3 описана поетапна реалізація двох підходів до детекції об'єктів (YOLOv8 та Faster R-CNN) на платформі Google Colab з використанням PyTorch. Наведено підготовку даних, конфігурацію моделей, гіперпараметри навчання та результати роботи. Проведено аналіз отриманих метрик (precision, recall, mAP) і показано, що обидві моделі підтверджують свою здатність виявляти та класифікувати об'єкти Pascal VOC. Наведені зображення з результатами роботи моделей демонструють приклади виявлення та дозволяють візуально порівняти якість роботи обох підходів на одних і тих самих зображеннях.

## ВИСНОВКИ

У кваліфікаційній роботі було розглянуто актуальну задачу розпізнавання об'єктів різної природи в системах відеоспостереження, виконано теоретичний аналіз сучасних підходів, спроектовано та реалізовано дві архітектури детекції об'єктів (одноетапну YOLOv8 та двоетапну Faster R-CNN), а також досліджено їхню інтеграцію в єдину систему.

Теоретичний аналіз (розділ 1) показав, що традиційні методи на основі ручного виділення ознак (SIFT, HOG, шаблонні підходи, статистичні класифікатори) забезпечують простоту реалізації і низькі обчислювальні вимоги, але значно поступаються у точності та стійкості сучасним методам глибинного навчання. Глибинні моделі (CNN, R-CNN, YOLO, SSD) дозволяють автоматично вивчати багаторівневі ознаки та демонструють високу ефективність у класифікації, детекції та сегментації об'єктів в складних сценах.

У розділі 2 проведений детальний опис архітектур YOLOv8 і Faster R-CNN, обґрунтований вибір PyTorch та Google Colab як середовища реалізації, а також побудову етапів підготовки даних, аугментації і навчання. Було розроблено інтегровану систему, що поєднує швидкість одноетапного підходу та точність двоетапного за допомогою алгоритму злиття прогнозів (Fusion).

Реалізація та тестування (розділ 3) підтвердили ефективність обох моделей на Pascal VOC. YOLOv8 забезпечив високу швидкодію (~45 FPS) при  $mAP@0.5 \approx 0.85$ ,  $mAP@[0.5:0.95] \approx 0.70$ ,  $precision \approx 0.88$ ,  $recall \approx 0.80$ . Faster R-CNN досяг вищої точності ( $mAP@0.5 \approx 0.90$ ,  $mAP@[0.5:0.95] \approx 0.75$ ,  $precision \approx 0.92$ ,  $recall \approx 0.85$ ), але з нижчою продуктивністю (~5 FPS). Порівняльний аналіз засвідчив, що реальний вибір детектора має базуватися на компромісі між точністю та швидкістю залежно від вимог системи.

Інтеграція моделей у єдину систему показала можливість поєднання переваг обох підходів. Застосування швидкого одноетапного детектора для первинного відбору кандидатів та подальше уточнення двоетапним

детектором дозволяє досягти одночасно високої продуктивності й надійності детекції в режимі реального часу.

**Практична значущість та перспективи:**

- розроблена інтегрована архітектура може бути впроваджена в сучасні системи відеоспостереження для підвищення якості моніторингу зон підвищеного ризику (транспорт, охорона об'єктів, промислові підприємства);
- експериментальні результати свідчать про доцільність використання YOLOv8 на платформі Edge (камери з вбудованим прискорювачем) та етапної обробки через Faster R-CNN у центральному сервері;
- подальші дослідження можуть бути спрямовані на впровадження self-supervised та few-shot методів для обробки малорозмічених доменів, а також на оптимізацію моделей для роботи зі спеціалізованими сенсорами (тепловізори, LIDAR) і вбудованими пристроями з обмеженими ресурсами.

Таким чином, виконане дослідження підтвердило ефективність глибинних методів розпізнавання об'єктів і продемонструвала практично життєздатний підхід до інтеграції їх у єдину систему спостереження. Результати кваліфікаційної роботи можуть бути використані для створення високопродуктивних, гнучких і надійних систем комп'ютерного зору в різних галузях.

**СПИСОК ВИКОРИСТАНИХ ДЖЕРЕЛ**

1. Bochkovskiy A., Wang C.-Y., Liao H.-Y. M. YOLOv4: Optimal Speed and Accuracy of Object Detection [Електронний ресурс] // ArXiv preprint. – 2020. – № 2004.10934. – Режим доступу: <https://arxiv.org/abs/2004.10934>
2. Carion N., Massa F., Synnaeve G., Usunier N., Kirillov A., Zagoruyko S. End-to-End Object Detection with Transformers (DETR) [Електронний ресурс] // ArXiv preprint. – 2020. – № 2005.12872. – Режим доступу: <https://arxiv.org/abs/2005.12872>
3. Chen T., Kornblith S., Norouzi M., Hinton G. A Simple Framework for Contrastive Learning of Visual Representations (SimCLR) [Електронний ресурс] // Proceedings of the 37th International Conference on Machine Learning (ICML 2020). – 2020. – С. 1597–1607. – Режим доступу: <https://proceedings.mlr.press/v119/chen20j.html>
4. Dalal N., Triggs B. Histograms of Oriented Gradients for Human Detection [Електронний ресурс] // Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2005). – 2005. – С. 886–893. – Режим доступу: <https://ieeexplore.ieee.org/document/1467360>
5. Dosovitskiy A., Beyer L., Kolesnikov A., Weissenborn D., Zhai X., Unterthiner T., Dehghani M., Minderer M., Heigold G., Gelly S., Uszkoreit J., Houlsby N. An Image is Worth 16 × 16 Words: Transformers for Image Recognition at Scale [Електронний ресурс] // Proceedings of the 9th International Conference on Learning Representations (ICLR 2021). – 2021. – Режим доступу: <https://arxiv.org/abs/2010.11929>
6. Hashemi N. S., Babaei Aghdam R., Bayat Ghiasi A. S., Fatemi P. Template Matching Advances and Applications in Image Analysis [Електронний ресурс] // American Scientific Research Journal for Engineering, Technology, and Sciences (ASRJETS). – 2016. – Т. 26, № 3. – С. 91–108. –

Режим

доступу:

[https://asrjetsjournal.org/index.php/American\\_Scientific\\_Journal/article/view/2378](https://asrjetsjournal.org/index.php/American_Scientific_Journal/article/view/2378)

7. He K., Zhang X., Ren S., Sun J. Deep Residual Learning for Image Recognition [Электронный ресурс] // Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2016). – 2016. – С. 770–778. – Режим доступа: <https://ieeexplore.ieee.org/document/7780459>
8. He K., Gkioxari G., Dollár P., Girshick R. Mask R-CNN [Электронный ресурс] // Proceedings of the IEEE International Conference on Computer Vision (ICCV 2017). – 2017. – С. 2961–2969. – Режим доступа: <https://ieeexplore.ieee.org/document/8237584>
9. Howard A. G., Zhu M., Chen B., Kalenichenko D., Wang W., Weyand T., Andreetto M., Adam H. MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications [Электронный ресурс] // ArXiv preprint. – 2017. – № 1704.04861. – Режим доступа: <https://arxiv.org/abs/1704.04861>
10. Liu W., Anguelov D., Erhan D., Szegedy C., Reed S., Fu C.-Y., Berg A. C. SSD: Single Shot Multibox Detector [Электронный ресурс] // ArXiv preprint. – 2016. – № 1512.02325. – Режим доступа: <https://arxiv.org/abs/1512.02325>
11. Long J., Shelhamer E., Darrell T. Fully Convolutional Networks for Semantic Segmentation [Электронный ресурс] // ArXiv preprint. – 2014. – № 1411.4038. – Режим доступа: <https://arxiv.org/abs/1411.4038>
12. Lowe D. G. Distinctive Image Features from Scale-Invariant Keypoints [Электронный ресурс] // International Journal of Computer Vision. – 2004. – Т. 60, № 2. – С. 91–110. – Режим доступа: <https://www.cs.ubc.ca/~lowe/papers/ijcv04.pdf>

13. Paszke A. et al. PyTorch: An Imperative Style, High-Performance Deep Learning Library [Электронный ресурс] // Advances in Neural Information Processing Systems (NeurIPS 2019). – 2019. – Режим доступа: [https://papers.nips.cc/paper\\_files/paper/2019/hash/bdbca288fee7f92f2bfa9f7012727740-Abstract.html](https://papers.nips.cc/paper_files/paper/2019/hash/bdbca288fee7f92f2bfa9f7012727740-Abstract.html)
14. Redmon J., Divvala S., Girshick R., Farhadi A. You Only Look Once: Unified, Real-Time Object Detection [Электронный ресурс] // Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2016). – 2016. – С. 779–788. – Режим доступа: [https://www.cv-foundation.org/openaccess/content\\_cvpr\\_2016/papers/Redmon\\_You\\_Only\\_Look\\_CVPR\\_2016\\_paper.pdf](https://www.cv-foundation.org/openaccess/content_cvpr_2016/papers/Redmon_You_Only_Look_CVPR_2016_paper.pdf)
15. Ren S., He K., Girshick R., Sun J. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks [Электронный ресурс] // ArXiv preprint. – 2015. – № 1506.01497. – Режим доступа: <https://arxiv.org/abs/1506.01497>
16. Ronneberger O., Fischer P., Brox T. U-Net: Convolutional Networks for Biomedical Image Segmentation [Электронный ресурс] // ArXiv preprint. – 2015. – № 1505.04597. – Режим доступа: <https://arxiv.org/abs/1505.04597>

**ДОДАТКИ****Додаток А**

МІНІСТЕРСТВО ОСВІТИ І НАУКИ УКРАЇНИ  
Харківський національний університет імені В. Н. Каразіна

Навчально-науковий інститут комп'ютерних наук та штучного інтелекту  
Кафедра комп'ютерних систем та робототехніки  
Рівень вищої освіти (освітньо-кваліфікаційний рівень) **Бакалавр**  
Галузь знань: 12 – Інформаційні технології  
Спеціальність: 123 «Комп'ютерна інженерія»  
Освітня програма «Комп'ютерна інженерія»

ЗАТВЕРДЖУЮ

В.о. завідувача кафедри комп'ютерних  
систем та робототехніки  
к. ф.-м. н., доц. ХРУСЛОВ М. М.  
«02» жовтня 2024 року

**З А В Д А Н Н Я**  
**НА КВАЛІФІКАЦІЙНУ РОБОТУ**

**Мотика Володимир Назарович**  
(прізвище, ім'я, по батькові студента)

1. Тема роботи: **Модель розпізнавання об'єктів різної природи у системах спостереження**

керівник роботи **Стрілець Вікторія Євгенівна, кандидат технічних наук**  
(прізвище, ім'я, по батькові, науковий ступінь, вчене звання)

затверджені наказом по університету **НАКАЗ № 4101-5/962** від **16.04.2025 року**

2. Строк подання студентом роботи ***30 травня 2025 року***

3. Перелік питань, які потрібно розробити

1. Аналіз сучасних методів розпізнавання образів, включно з класичними та глибинними
2. Створити моделі для виявлення об'єктів на основі нейронних мереж.
3. Провести навчання і тестування моделей із використанням GPU-обчислень, оцінити їх ефективність.
4. Визначити шляхи подальшого вдосконалення запропонованих моделей.

## 4. План роботи

№ з/п	Назви етапів роботи	Термін виконання етапів роботи
1	Літературний огляд з проблематики розпізнавання об'єктів	02.10.2024 – 23.10.2024
2	Аналіз моделей глибокого навчання для розпізнавання образів (R-CNN, Fast/Faster R-CNN, YOLO тощо)	28.10.2024 – 14.11.2024
3	Збір та підготовка даних (наприклад, PASCAL VOC): очищення, формування анотацій, розмітка	16.11.2024 – 04.12.2024
4	Проектування та налаштування моделі YOLOv8	06.12.2024 – 15.01.2025
5	Розробка та налаштування моделі Faster R-CNN	16.01.2025 – 15.02.2025
6	Тестування моделей, первинний аналіз результатів (mAP, Precision, Recall тощо)	16.02.2025 – 10.03.2025
7	Порівняльний аналіз ефективності, формування рекомендацій щодо поліпшення моделей	11.03.2025 – 25.03.2025
8	Підготовка та оформлення розділів кваліфікаційної роботи	26.03.2025 – 20.04.2025
9	Оформлення пояснювальної записки та узагальнення результатів	21.04.2025 – 10.05.2025
10	Подання кваліфікаційної роботи керівнику та рецензенту	11.05.2025 – 30.05.2025

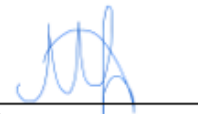
5. Дата видачі завдання *02 жовтня 2024 року.*

— Студент

Мотика В.Н.

ініціали, прізвище

підпис



Керівник роботи

Стрілець В.Є.

ініціали, прізвище

підпис



## Додаток Б

Затверджую

«\_\_\_\_\_» \_\_\_\_\_ 2025 р.

**Технічне завдання**

на розробку програмного виробу  
«Модель розпізнавання об'єктів різної  
природи у системах спостереження»

Назва розділу	Назва і зміст підрозділу
1. Вступ	<p>1.1. Назва проекту: «Модель розпізнавання об'єктів різної природи у системах спостереження»</p> <p>1.2. Галузь застосування:</p> <ul style="list-style-type: none"> <li>– Відеоаналітика та системи безпеки (спостереження на об'єктах, вуличні камери);</li> <li>– Індустріальний контроль якості (виявлення дефектів на лінії);</li> <li>– Автономні системи (дрони, робототехніка);</li> <li>– Смарт-міста (моніторинг трафіку, пішоходів).</li> </ul>
2. Підстава для розробки	<p>2.1. Освітній курс за спеціальністю 123 – Комп'ютерна інженерія</p> <p>2.2. Завдання на дипломну роботу бакалавра, затверджено наказом ХНУ імені В. Н. Каразіна № 4101-5/962 від 16.04.2025 р. (представить як Додаток А до пояснювальної записки до кваліфікаційної роботи).</p>
3. Призначення розробки	<p>3.1. Мета: підвищити точність і швидкодію виявлення об'єктів різної природи (транспорт, люди, тварини) у відеопотоці</p> <p>3.2. Призначення: Автоматизоване детектування та класифікація об'єктів у реальному часі або заздалегідь записаних кадрах</p> <p>3.3. Початкові дані для розробки:</p> <ul style="list-style-type: none"> <li>– Зображення/відео з відкритих датасетів (Pascal VOC, COCO)</li> <li>– Анотації bounding-box (класи, координати)</li> <li>– Відеопотоки з IP-камер або файли *.mp4 для тестування</li> </ul>
4. Технічні вимоги до програмного виробу	<p>4.1. Функціональні характеристики:</p> <ul style="list-style-type: none"> <li>– Кількість класів: до 20 різних об'єктів</li> <li>– Метод детекції: одноетапний (YOLOv8) і двоетапний (Faster R-CNN)</li> <li>– Швидкодія: <math>\geq 20</math> FPS на GPU для моделі YOLOv8</li> <li>– Точність: <math>mAP@0.5 \geq 0.73</math> (YOLOv8), <math>\geq 0.50</math> (Faster R-CNN)</li> </ul> <p>4.2. Надійність:</p> <ul style="list-style-type: none"> <li>– Стабільна робота при змінному освітленні, частковому перекритті та шумі</li> <li>– Мінімізація хибних спрацьовувань через порогову фільтрацію</li> </ul>

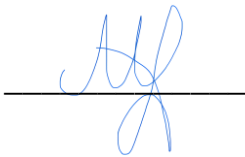
	<p>confidence score</p> <ul style="list-style-type: none"> <li>– Резервна обробка кадрів у разі втрати фреймів</li> </ul> <p>4.3. Умови експлуатації:</p> <ul style="list-style-type: none"> <li>– Сервер із GPU (CUDA 11.x), або edge-пристрій із підтримкою тензорного прискорення</li> <li>– Температура 10–35 °C, відносна вологість до 80% (без конденсації)</li> <li>– Підключення до локальної чи хмарної мережі для передачі відео й даних</li> </ul> <p>4.4. Технічні засоби:</p> <ul style="list-style-type: none"> <li>– Сервер на базі NVIDIA GPU (наприклад, RTX 3060/3090)</li> <li>– IP-камери з роздільною здатністю <math>\geq 720p</math></li> <li>– SSD для зберігання ваг моделей і логів</li> </ul> <p>4.5. Сумісність:</p> <ul style="list-style-type: none"> <li>– Платформи: Linux (Ubuntu 20.04+), Windows 10+</li> <li>– Фреймворки: PyTorch, ONNX Runtime, OpenCV</li> <li>– Мови: Python 3.8–3.10</li> </ul> <p>4.6. Маркування та упаковка:</p> <ul style="list-style-type: none"> <li>– Docker-образ із усіма залежностями</li> <li>– рір-пакет або JAR-обгортка для інтеграції в сторонні продукти</li> <li>– Вага образу: <math>\leq 2</math> GB</li> </ul> <p>4.7. Транспортування та зберігання:</p> <ul style="list-style-type: none"> <li>– Зберігати контейнер чи ваги моделей на сервері або NAS</li> <li>– Архівація моделей у форматі .zip/.tar.gz</li> <li>– Резервне копіювання кожної нової версії</li> </ul> <p>4.8. Спеціальні вимоги:</p> <ul style="list-style-type: none"> <li>– Наявність GPU з мінімум 8 GB VRAM</li> <li>– Доступ до зовнішніх репозиторіїв для автоматичного оновлення ваг</li> <li>– Сертифікація на відповідність GDPR (якщо працює з особистими даними)</li> </ul>
5. Вимоги до програмної документації.	<ul style="list-style-type: none"> <li>– Опис API (ендпоінтів infer/train) у форматі Swagger/OpenAPI</li> <li>– Керівництво користувача (User Manual) із прикладами запуску</li> <li>– Коментарі в коді та інструкція з розгортання (README.md)</li> </ul>
6. Техніко-економічні показники	<ul style="list-style-type: none"> <li>– Вартість одного inference (енергоспоживання та час)</li> <li>– Вартість обладнання на 1 фонтан (GPU + камера + сервер)</li> <li>– ROI: окупність системи протягом 12 міс. при 24/7 моніторингу</li> </ul>
7. Стадії і етапи розробки	<ul style="list-style-type: none"> <li>– Літературний огляд з проблематики розпізнавання об'єктів</li> <li>– Аналіз моделей глибинного навчання для розпізнавання образів (R-CNN, Fast/Faster R-CNN, YOLO тощо)</li> <li>– Збір та підготовка даних (наприклад, PASCAL VOC): очищення, формування анотацій, розмітка</li> <li>– Проектування та налаштування моделі YOLOv8</li> <li>– Розробка та налаштування моделі Faster R-CNN</li> <li>– Тестування моделей, первинний аналіз результатів (mAP, Precision, Recall тощо)</li> </ul>

	<ul style="list-style-type: none"> <li>– Порівняльний аналіз ефективності, формування рекомендацій щодо поліпшення моделей</li> <li>– Підготовка та оформлення звітів до кваліфікаційної роботи</li> </ul>
8. Порядок контролю і приймання	<ul style="list-style-type: none"> <li>– Функціональні тести: перевірка всіх класів об'єктів на контрольних зображеннях</li> <li>– Навантажувальні тести: <math>\geq 20</math> FPS стабільно протягом 1 години</li> <li>– Тест на стабільність: відсутність витоків пам'яті при 10 000 кадрів</li> <li>– Оцінка точності: <math>mAP@0.5 \geq</math> задане значення (наприклад, 0.70 для YOLOv8)</li> <li>– Звіт про результати та підписання Актів виконаних робіт</li> </ul>

Виконавець

студент групи КІ-41

Мотика В.Н.



Замовник

к.т.н., доцент ЗВО

Стрілець В.Є.



**Додаток В****Програма і методика випробувань програмного виробу**

«Модель розпізнавання об'єктів різної природи у системах спостереження»

**1 Об'єкт випробувань**

1.1 Назва: «Модель розпізнавання об'єктів різної природи у системах спостереження»

1.2 Область застосування:

- Відеоаналітика та системи безпеки (спостереження на об'єктах, вуличні камери)
- Індустріальний контроль якості (виявлення дефектів на лінії)
- Автономні системи (дрони, робототехніка)
- Смарт-міста (моніторинг трафіку, пішоходів)

**2. Мета випробувань**

Загальна мета: Підвищити точність і швидкодію виявлення об'єктів різної природи (транспорт, люди, тварини) у відеопотоці

Специфічні цілі: Досягти  $mAP@0.5 \geq 0.75$  на власних тестових даних, забезпечити затримку інференсу  $\leq 40$  мс на кадр, реалізувати інтерфейс для автоматичного порівняння двох моделей

**3. Загальні положення****3.1 Підстави для проведення випробувань**

Вимоги акредитаційної комісії університету.

**3.2 Місце і тривалість випробувань**

Лабораторії факультету, один місяць.

**3.3 Обсяг випробувань**

Повний цикл випробувань всіх модулів системи.

**3.4 Організації, які беруть участь у випробуваннях**

ХНУ імені В. Н. Каразіна, факультет комп'ютерних наук.

**4. Вимоги до програми або програмного виробу**

4.1. Функціональні характеристики:

- Кількість класів: до 20 різних об'єктів
- Метод детекції: одноетапний (YOLOv8) і двоетапний (Faster R-CNN)
- Швидкодія:  $\geq 20$  FPS на GPU для моделі YOLOv8
- Точність:  $mAP@0.5 \geq 0.73$  (YOLOv8),  $\geq 0.50$  (Faster R-CNN)

4.2. Надійність:

- Стабільна робота при змінному освітленні, частковому перекритті та шумі
- Мінімізація хибних спрацьовувань через порогову фільтрацію confidence score
- Резервна обробка кадрів у разі втрати фреймів

#### 4.3. Умови експлуатації:

- Сервер із GPU (CUDA 11.x), або edge-пристрій із підтримкою тензорного прискорення
- Температура 10–35 °С, відносна вологість до 80% (без конденсації)
- Підключення до локальної чи хмарної мережі для передачі відео й даних

#### 4.4. Технічні засоби:

- Сервер на базі NVIDIA GPU (наприклад, RTX 3060/3090)
- IP-камери з роздільною здатністю  $\geq 720p$
- SSD для зберігання ваг моделей і логів

#### 4.5. Сумісність:

- Платформи: Linux (Ubuntu 20.04+), Windows 10+
- Фреймворки: PyTorch, ONNX Runtime, OpenCV
- Мови: Python 3.8–3.10

#### 4.6. Маркування та упаковка:

- Docker-образ із усіма залежностями
- рір-пакет або JAR-обгортка для інтеграції в сторонні продукти
- Вага образу:  $\leq 2$  GB

#### 4.7. Транспортування та зберігання:

- Зберігати контейнер чи ваги моделей на сервері або NAS
- Архівація моделей у форматі .zip/.tar.gz
- Резервне копіювання кожної нової версії

#### 4.8. Спеціальні вимоги:

- Наявність GPU з мінімум 8 GB VRAM
- Доступ до зовнішніх репозиторіїв для автоматичного оновлення ваг
- Сертифікація на відповідність GDPR (якщо працює з особистими даними)

### 5. Вимоги до програмної документації

- Опис API (ендпоінтів infer/train) у форматі Swagger/OpenAPI
- Керівництво користувача (User Manual) із прикладами запуску
- Коментарі в коді та інструкція з розгортання (README.md)

### 6. Засоби і порядок випробувань

#### 6.1 Засоби випробувань

Усі тести проводилися на середовищі Google Colab із GPU (NVIDIA Tesla T4). Використовувалися такі інструменти і матеріали:

- PyTorch (версія  $\geq 1.12$ ) для реалізації Faster R-CNN.
- Ultralytics YOLO (версія 8.x) для реалізації YOLOv8.
- torchvision (для датасету Pascal VOC та допоміжних трансформацій).
- GPU: NVIDIA Tesla T4 (16 GB VRAM).
- CPU: стандартна віртуальна машина Colab (2 Intel Xeon vCPUs).
- Matplotlib для побудови кривих Precision/Recall/ mAP за епохами.

## 6.2 Порядок проведення випробувань

### Тест №1 – Precision YOLOv8 за 30 епох

**Мета тесту:** Визначити, як змінюється точність (Precision) моделі YOLOv8 у процесі тренування від початку до 30-ї епохи.

#### Процедура:

1. Ініціалізувати YOLOv8n з передтренуваними вагами COCO.
2. Налаштувати:  $\text{epochs} = 30$ ,  $\text{batch\_size} = 16$ ,  $\text{imgsz} = 640$ ,  $\text{lr0} = 0.001$ .
3. Після кожної епохи обчислювати Precision на валідаційній вибірці Pascal VOC.

4. Зберігати й будувати криву Precision(епоха).

#### Результат тесту:

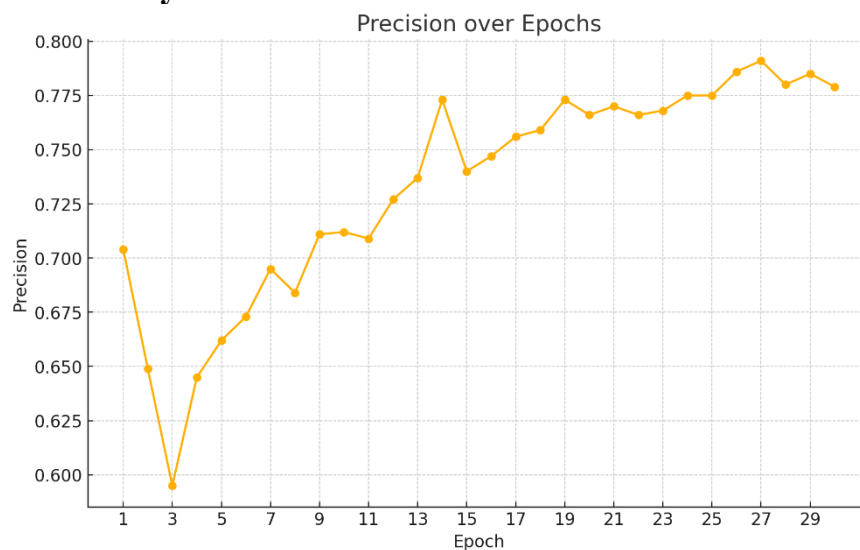


Рисунок В.1 – Результат тесту №1

### Тест №2 – Recall YOLOv8 за 30 епох

**Мета тесту:** Оцінити динаміку повноти (Recall) YOLOv8 протягом 30 епох тренування.

#### Процедура:

1. Використати ті самі налаштування YOLOv8n ( $\text{epochs} = 30$ ,  $\text{batch\_size} = 16$ ,  $\text{lr0} = 0.001$ ).
2. Після кожної епохи обчислювати Recall на валідаційній вибірці.
3. Будувати криву Recall(епоха).

#### Результат тесту:

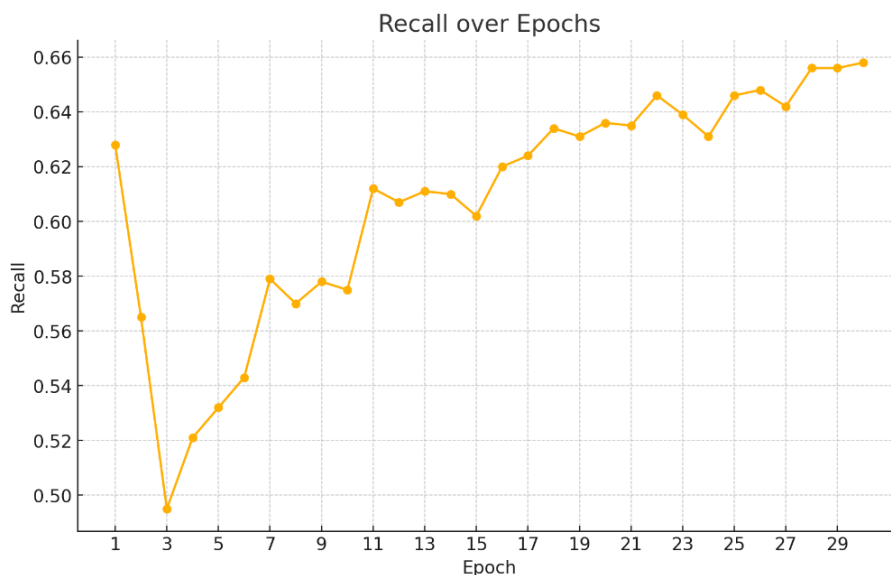


Рисунок В.2 – Результат тесту №2

**Тест №3 – Precision Faster R-CNN за 30 епох**

**Мета тесту:** Визначити, як змінюється Precision моделі Faster R-CNN під час тренування до 30-ї епохи.

**Процедура:**

1. Ініціалізувати Faster R-CNN із backbone = ResNet-50 FPN (передтренований на COCO).
2. Налаштувати: epochs = 30, batch\_size = 8, lr = 0.005, momentum = 0.9, weight\_decay = 0.0005.
3. Після кожної епохи обчислювати Precision на валідаційній частині Pascal VOC.
4. Створювати криву Precision(епоха).

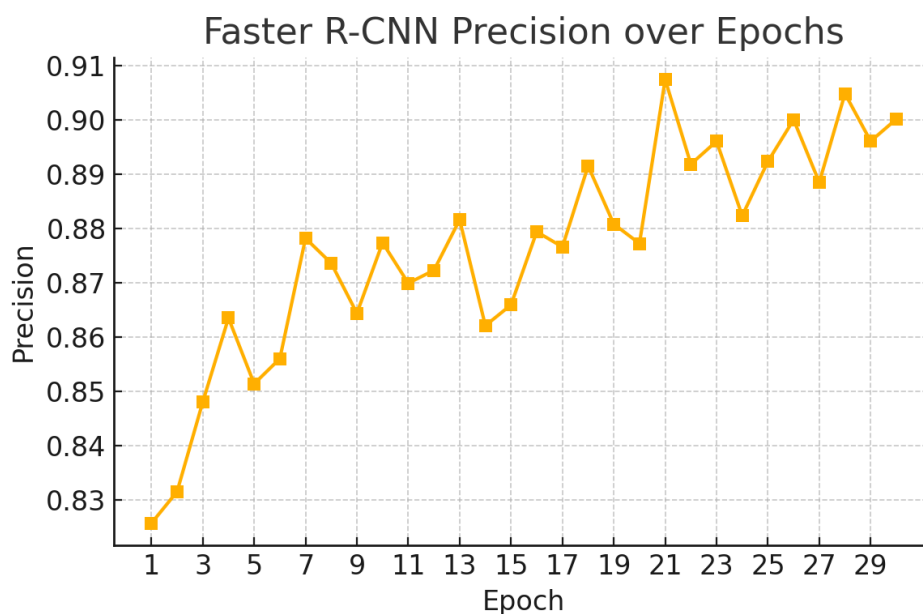
**Результат тесту:**

Рисунок В.3 – Результат тесту №3

### Тест №4 – Recall Faster R-CNN за 30 епох

**Мета тесту:** Оцінити динаміку Recall моделі Faster R-CNN протягом 30 епох.

#### Процедура:

1. Використати ті самі налаштування Faster R-CNN (epochs = 30, batch\_size = 8, lr = 0.005).
2. Після кожної епохи обчислювати Recall на валідації.
3. Будувати криву Recall(епоха).

#### Результат тесту:

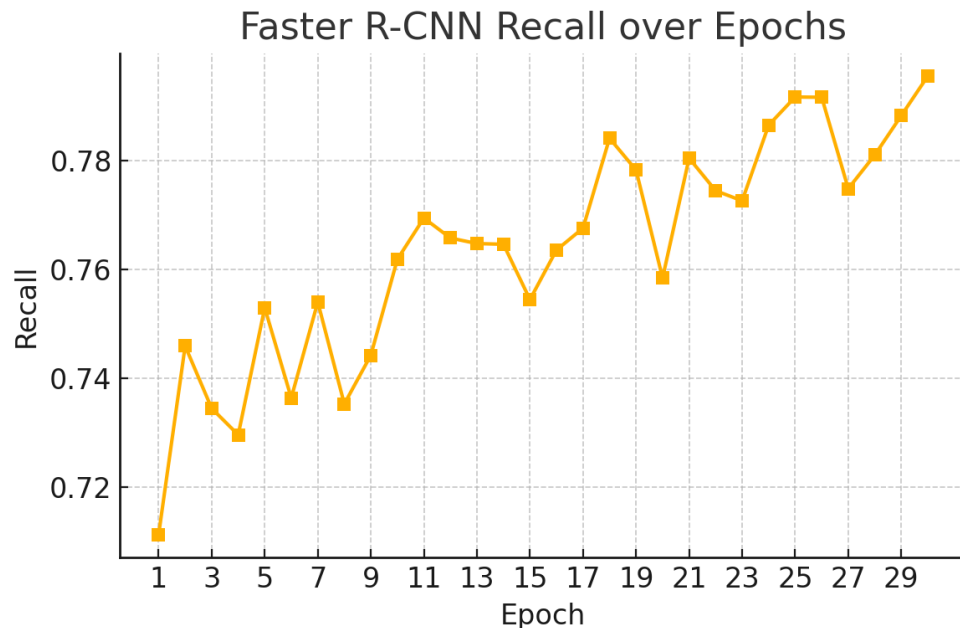


Рисунок В.4 – Результат тесту №4

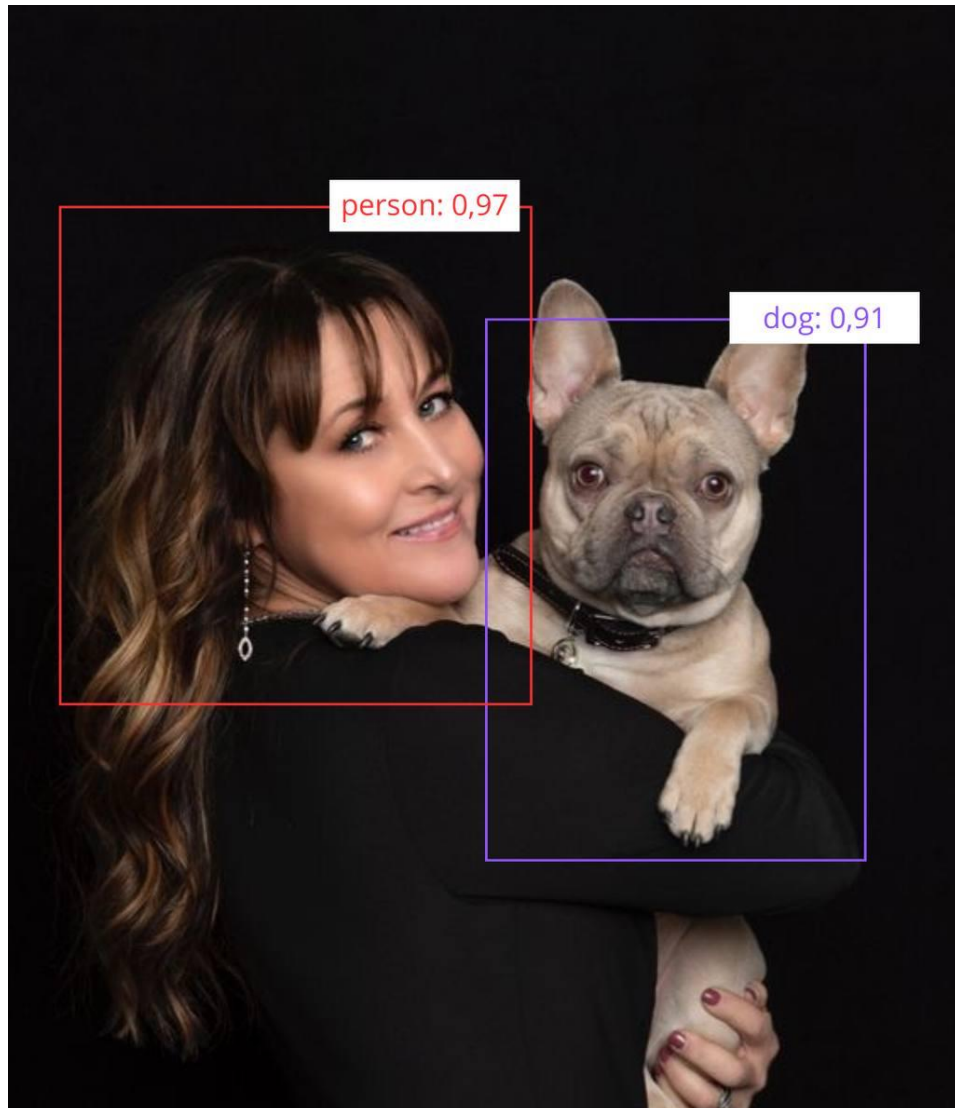
### Тест №5 – Інференс Faster R-CNN на прикладі (людина + собака)

**Мета тесту:** Порівняти якість локалізації та впевненість (confidence) Faster R-CNN в інференсі на контрольному зображенні.

#### Процедура:

1. Взяти модель Faster R-CNN із вагами з останньої (30-ї) епохи.
2. Подати одне зображення (жінка з собакою) без додаткових аугментацій.
3. Отримати прогноз: список bounding-box'ів із класами й confidence.
4. Візуалізувати результат: червона рамка – “person: 0,97”; фіолетова рамка – “dog: 0,91”.

#### Результат тесту:



## Faster R-CNN Inference

Рисунок В.5 – Результат тесту №5

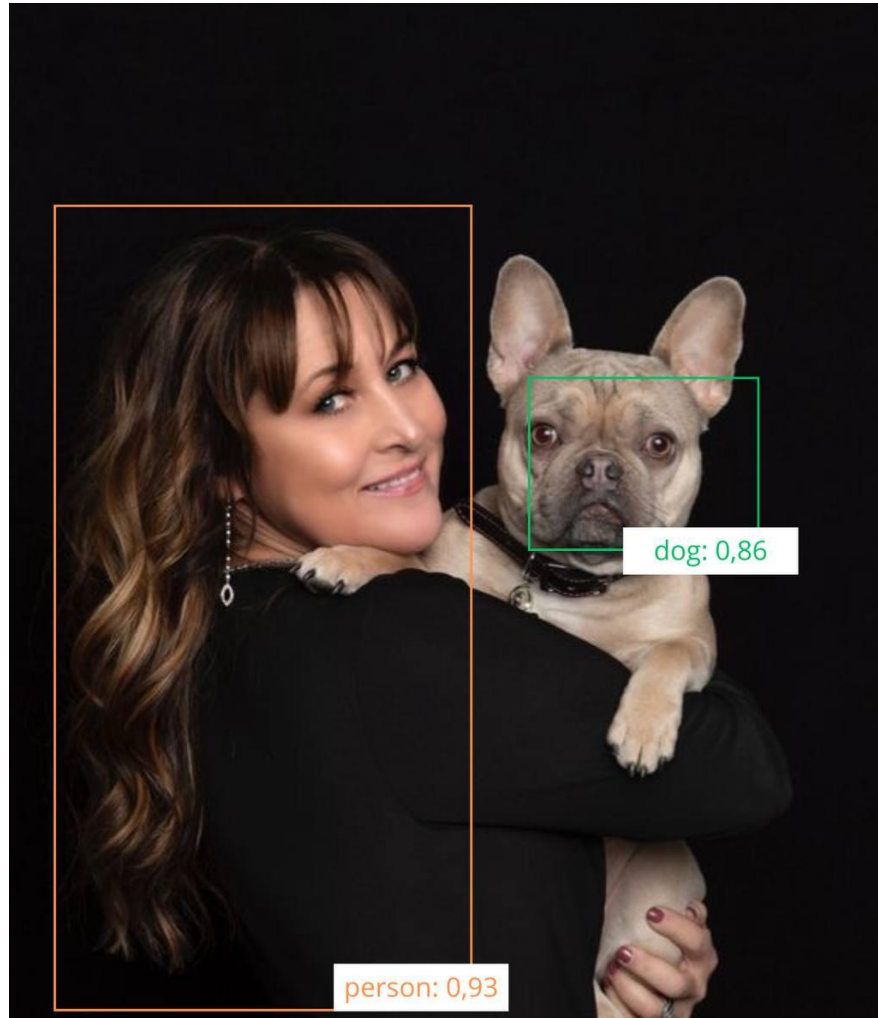
### Тест №6 – Інференс YOLOv8 на прикладі (людина + собака)

**Мета тесту:** Порівняти якість локалізації та впевненість YOLOv8 в інференсі на тому ж контрольному зображенні.

#### Процедура:

1. Взяти модель YOLOv8n із вагами з останньої (30-ї) епохи.
2. Масштабувати зображення до 640×640 (зберігаючи співвідношення сторін і нормалізацією).
3. Отримати прогноз: список bounding-box'ів із класами й confidence.
4. Візуалізувати результат: помаранчева рамка – “person: 0,93”; зелена рамка – “dog: 0,86”.

#### Результат тесту:



## YOLOv8 Inference

Рисунок В.6 – Результат тесту №6

Виконавець

студент групи КІ-41

Мотика В.Н.

A handwritten signature in blue ink, consisting of stylized letters, positioned over a horizontal line.

## Лістинг програми

```
!pip uninstall -y torch torchvision torchaudio
!pip install torch==2.0.1 torchvision==0.15.2 torchaudio==2.0.2 --extra-
index-url https://download.pytorch.org/whl/cu117
!pip install --upgrade ultralytics

import os
import torchvision
from ultralytics import YOLO
from google.colab import files
from PIL import Image
import torch

class_names = [
    "aeroplane", "bicycle", "bird", "boat", "bottle", "bus", "car", "cat", "chair"
, "cow",
    "diningtable", "dog", "horse", "motorbike", "person", "pottedplant", "sheep"
, "sofa", "train", "tvmonitor"
]

dataset_root = os.path.join(os.getcwd(), "voc_data")
images_dir_train = os.path.join(os.getcwd(), "VOC", "images", "train")
labels_dir_train = os.path.join(os.getcwd(), "VOC", "labels", "train")
images_dir_val = os.path.join(os.getcwd(), "VOC", "images", "val")
labels_dir_val = os.path.join(os.getcwd(), "VOC", "labels", "val")

os.makedirs(images_dir_train, exist_ok=True)
os.makedirs(labels_dir_train, exist_ok=True)
os.makedirs(images_dir_val, exist_ok=True)
os.makedirs(labels_dir_val, exist_ok=True)

ds_train = torchvision.datasets.VOCDetection(
    root=dataset_root, year="2012", image_set="train", download=True
)
ds_val = torchvision.datasets.VOCDetection(
    root=dataset_root, year="2012", image_set="val", download=True
)

def convert_bbox(b, w, h):
    xmin, ymin, xmax, ymax = b
    x = (xmin + xmax) / 2 / w
    y = (ymin + ymax) / 2 / h
    ww = (xmax - xmin) / w
    hh = (ymax - ymin) / h
    return x, y, ww, hh

i = 0
```

```

for img, anno in ds_train:
    w, h = img.size
    objs = anno["annotation"]["object"]
    img_path = os.path.join(images_dir_train, f"{i}.jpg")
    img.save(img_path)
    lbl_path = os.path.join(labels_dir_train, f"{i}.txt")
    with open(lbl_path, "w") as f:
        if isinstance(objs, dict):
            objs = [objs]
        for o in objs:
            cname = o["name"]
            if cname not in class_names:
                continue
            cid = class_names.index(cname)
            bb = o["bndbox"]
            x, y, ww, hh = convert_bbox(
                [float(bb["xmin"]), float(bb["ymin"]), float(bb["xmax"]),
float(bb["ymax"])],
                w, h
            )
            f.write(f"{cid} {x} {y} {ww} {hh}\n")
        i += 1
i = 0
for img, anno in ds_val:
    w, h = img.size
    objs = anno["annotation"]["object"]
    img_path = os.path.join(images_dir_val, f"{i}.jpg")
    img.save(img_path)
    lbl_path = os.path.join(labels_dir_val, f"{i}.txt")
    with open(lbl_path, "w") as f:
        if isinstance(objs, dict):
            objs = [objs]
        for o in objs:
            cname = o["name"]
            if cname not in class_names:
                continue
            cid = class_names.index(cname)
            bb = o["bndbox"]
            x, y, ww, hh = convert_bbox(
                [float(bb["xmin"]), float(bb["ymin"]), float(bb["xmax"]),
float(bb["ymax"])],
                w, h
            )
            f.write(f"{cid} {x} {y} {ww} {hh}\n")
        i += 1

with open("voc.yaml", "w") as f:
    f.write(
        f"train: {images_dir_train}\n"
        f"val: {images_dir_val}\n"

```

```
        f"nc: {len(class_names)}\n"  
        f"names: {class_names}\n"  
        "augment: True\n"  
    )  
  
print("CUDA Available:", torch.cuda.is_available())  
print("CUDA Device Name:",  
      torch.cuda.get_device_name(0) if torch.cuda.is_available() else "No  
GPU")  
  
model = YOLO("yolov8s.pt")  
model.train(  
    data="voc.yaml",  
    epochs=50,  
    batch=12,  
    imgsz=512,  
    name="voc_exp",  
    device=0,  
)  
  
model.val()  
  
uploaded = files.upload()  
for name in uploaded.keys():  
    results = model.predict(source=name, conf=0.5, iou=0.45, save=True)  
    for r in results:  
        for b in r.bboxes:  
            c = int(b.cls[0])  
            cf = float(b.conf[0])  
            print(f"{class_names[c]} {cf:.2f}")
```