

**BỘ GIÁO DỤC VÀ ĐÀO TẠO
ĐẠI HỌC QUỐC GIA TP. HỒ CHÍ MINH
TRƯỜNG ĐẠI HỌC CÔNG NGHỆ THÔNG TIN**



BÁO CÁO ĐỒ ÁN
MÔN CS114 – MÁY HỌC

ĐỀ TÀI: FACE MASK DETECTION

Giảng viên hướng dẫn : Phạm Nguyễn Trường An
Lê Đình Duy

Nhóm thực hiện :

Trần Nhật Đức - 21521968
Hồ Trung Tín - 21521536
Phạm Đức Toàn - 21521550

Lớp: CS114.O11.KHCL

Tp HCM, tháng 12 năm 2023

NHẬN XÉT CỦA GIẢNG VIÊN HƯỚNG DẪN

Tp.HCM, ngày ... tháng ... năm ...

GVHD

(Ký tên)

Mục lục

Chương 0: Update sau khi vấn đáp	3
1. mAP và các khái niệm liên quan:	3
2. Cập nhật dataset:	6
3. So sánh kết quả giữa mô hình Yolov5 với Faster R-CNN:	6
Chương 1: Tổng quan về đề án	7
1. Tóm tắt đề án:	7
2. Mô tả Input và Output của bài toán:	7
Chương 2: Xây dựng dataset	8
1. Tiền xử lý dữ liệu:	8
2. Gán Nhãn Bằng LabelImg:	8
3. Về Dataset:	9
Chương 3: Cơ sở lý thuyết, training và đánh giá model	10
1. YOLO:	10
2. Model YOLOv5:	10
3. Giai đoạn training model:	14
4. Về quá trình trích xuất đặc trưng (feature engineering):	14
5. Lý do chọn model:	15
6. Quá trình train model:	16
7. Kết quả và đánh giá:	17
8. Thực hiện dự đoán:	19
TÀI LIỆU THAM KHẢO:	21

link tới github:

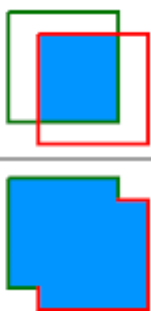
<https://github.com/jayjay2k3/CS114>

Chương 0: Update sau khi vấn đáp

1. mAP và các khái niệm liên quan:

- Tìm hiểu về IoU:

- IoU (Intersection over Union) là một phép đo được sử dụng để đo độ chính xác của việc dự đoán bounding box bằng cách so sánh giữa bounding box dự đoán và bounding box thực tế của đối tượng.
- IoU được tính bằng tỉ lệ giữa diện tích phần giao (intersection) và diện tích phần hợp (union) của hai bounding box hoặc vùng quan tâm (regions of interest - ROI).

$$IOU = \frac{\text{area of overlap}}{\text{area of union}} = \frac{\text{Intersection}}{\text{Union}}$$


- Trong đó:
 - Phần giao là diện tích của vùng mà hai bounding box hoặc vùng quan tâm chồng lấn lên nhau.
 - Phần hợp là diện tích tổng của cả hai bounding box hoặc vùng quan tâm, bao gồm cả phần giao và phần không giao.
 - Giá trị IoU càng cao, thì việc dự đoán càng chính xác. Một IoU cao hơn một ngưỡng nhất định thường được coi là một dự đoán chính xác.



- Precision, Recall và Precision Recall Curve:

- Precision (Độ Chính Xác): Precision đo lường tỉ lệ của các dự đoán Positive mà thực sự là đúng so với tổng số các dự đoán positive. Precision thể hiện khả năng của mô hình trong việc tránh việc phân loại sai các mẫu là negative thành positive.

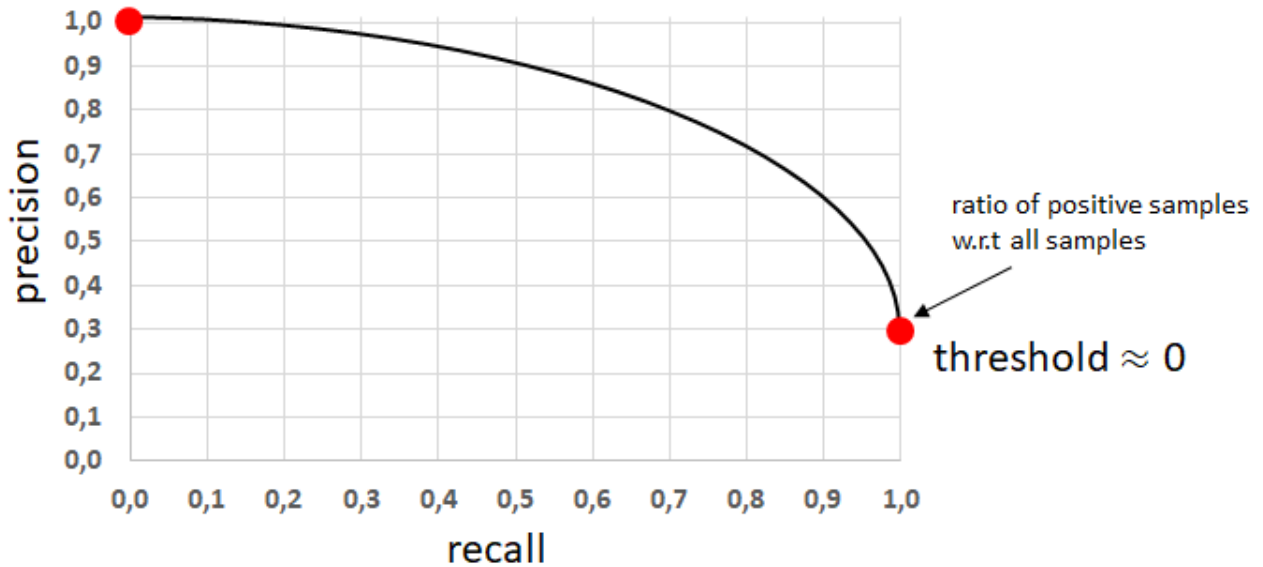
$$\begin{aligned}\text{Precision} &= \frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}} \\ &= \frac{\text{True Positive}}{\text{Total Predicted Positive}}\end{aligned}$$

- Recall (Độ Phủ): Recall đo lường tỉ lệ các mẫu positive thực sự được phát hiện đúng so với tổng số mẫu positive trong tập dữ liệu. Recall thể hiện khả năng của mô hình trong việc bắt được tất cả các mẫu thực sự là positive.

$$\text{Recall} = \frac{\text{True Positive}(TP)}{\text{True Positive}(TP) + \text{False Negative}(FN)}$$

- Precision-Recall Curve (Đường Cong Precision-Recall): Precision-Recall Curve là một biểu đồ mô tả mối quan hệ giữa precision và recall của một mô hình phân loại ở các ngưỡng quyết định khác nhau. Curve này thường được sử dụng để đánh giá hiệu suất của mô hình phân loại trong các trường hợp mà sự cân bằng giữa precision và recall là quan trọng. Đường cong Precision-Recall Curve thường được vẽ bằng cách thay đổi ngưỡng quyết định và ghi lại precision và recall tương ứng với mỗi ngưỡng đó.

threshold ≈ 1



- AP (Average Precision) là một chỉ số đánh giá hiệu suất trong các bài toán phân loại. Trong hình trên, AP chính là vùng diện tích nằm dưới đường cong Precision Recall.
 - AP lớn nếu vùng này lớn, suy ra đường cong có xu hướng gần góc trên bên phải và có nghĩa là tại các threshold khác nhau thì Precision và Recall đều khá cao. Từ đó suy ra model tốt.
 - AP nhỏ thì cả Precision và Recall đều khá thấp và model không tốt.

- Tìm hiểu về mAP:

mAP là viết tắt của "mean Average Precision", là một chỉ số tổng hợp để đánh giá hiệu suất của một hệ thống phát hiện đối tượng trong lĩnh vực computer vision, đặc biệt là trong các nhiệm vụ như phát hiện đối tượng và phân loại đối tượng.

Trong ngữ cảnh của bài toán object detection, mAP là sự kết hợp của các giá trị Average Precision (AP) cho từng lớp đối tượng trong tập dữ liệu. Điểm mAP là giá trị trung bình của các giá trị AP này, thường được tính toán bằng cách lấy tổng của tất cả các AP sau đó chia cho số lớp đối tượng.

Mean Average Precision Formula

$$\text{Mean Average Precision} = \frac{1}{n} \sum_{k=1}^{k=n} AP_k$$

n = the number of classes

AP_k = the average precision of class k

2. Cập nhật dataset:

- Sau khi nhận thấy dataset quá ít nên đã tiến hành gán nhãn thêm và hiện tại tổng dataset thu được 996 bức hình.

3. So sánh kết quả giữa mô hình Yolov5 với Faster R-CNN:

Đặc điểm	Yolov5	Faster R-CNN
Bộ dữ liệu	custom dataset	CoCo
map50	khoảng 81%	khoảng 99%
Kết quả	Tốc độ nhận diện chậm hơn Faster R-CNN khi có nhiều người trong khung hình	Tốc độ nhận diện nhanh hơn khi có nhiều người trong khung hình

Chương 1: Tổng quan về đề án

1. Tóm tắt đề án:

- Mô hình máy học này được phát triển nhằm mục đích phát hiện và phân loại người trong hình ảnh hoặc video liệu họ có đang đeo khẩu trang hay không.
- Đề án này có thể được ứng dụng trong nhiều lĩnh vực thực tế:
 - Đại dịch COVID-19 trong quá khứ đã làm nổi bật tầm quan trọng của việc đeo khẩu trang như một biện pháp phòng ngừa quan trọng trong bối cảnh dịch bệnh.
 - Các mô hình phát hiện khẩu trang có thể được sử dụng trong các trung tâm y tế để đảm bảo rằng bệnh nhân, người thăm và nhân viên tuân thủ các giao thức đeo khẩu trang. Điều này rất quan trọng trong việc ngăn chặn sự lây lan của các bệnh truyền nhiễm trong bệnh viện và phòng khám.
 - Trong các ngành công nghiệp nơi việc đeo khẩu trang đóng vai trò rất quan trọng, như các nhà máy sản xuất, các công trường xây dựng, thì việc triển khai mô hình phát hiện người đeo khẩu trang có thể đóng góp vào việc duy trì một môi trường làm việc an toàn.

2. Mô tả Input và Output của bài toán:

- INPUT:

1 hình ảnh hoặc 1 video từ camera an ninh chứa các khuôn mặt người cần được phân tích để xác định xem người đó có đeo khẩu trang hay không, nếu có thì họ có đeo đúng cách không.

- OUTPUT:

- Mỗi khuôn mặt người được phát hiện trong ảnh hoặc video sẽ được 1 bounding box bao quanh.
- Nếu người đó đang đeo khẩu trang, một nhãn "Đeo khẩu trang" sẽ được gán cho bounding box của họ.
- Nếu người đó không đeo khẩu trang, một nhãn "Không đeo khẩu trang" sẽ được gán cho bounding box của họ.

- Nếu người đó đeo khẩu trang không đúng cách (khẩu trang không che mũi, không che miệng), một nhãn "Đeo không đúng cách" sẽ được gán cho bounding box của họ.

Chương 2: Xây dựng dataset

- Dataset được nhóm thu thập và xây dựng thủ công.
- Nguyên nhân: các bộ dữ liệu có sẵn chưa phù hợp với yêu cầu và mục tiêu của bài toán là nhận diện người đeo khẩu trang từ camera an ninh. Mô hình này sẽ được áp dụng tại các camera an ninh ở nhiều nơi như sân bay, ga tàu, trạm xe buýt, trường học, và cửa hàng, vì vậy bộ dữ liệu được nhóm xây dựng thủ công bằng cách thu thập các hình ảnh và video từ camera an ninh ở nhiều nơi. Mục tiêu là có một bộ dữ liệu đa dạng với nhiều người, tình huống ánh sáng và phong cách đeo khẩu trang khác nhau.

1. Tiền xử lý dữ liệu:

- Sau khi thu thập, dữ liệu được tiền xử lý để chuẩn hóa và làm sạch. Chúng em sẽ loại bỏ các ảnh hoặc khung hình không phù hợp, thực hiện các biện pháp xử lý ảnh như cắt tỉa hoặc chỉnh sửa kích thước để làm giảm nhiễu và tối ưu hóa chất lượng dữ liệu.

2. Gán Nhãn Bằng LabelImg:

- Công cụ LabelImg được sử dụng để vẽ bounding box và gán nhãn cho các khuôn mặt người trong bộ dữ liệu. Quá trình này được thực hiện bằng cách mở mỗi hình ảnh trong công cụ LabelImg và vẽ các hộp giới hạn xung quanh khuôn mặt của mỗi người trong ảnh. Sau đó, mỗi hộp giới hạn được gán một nhãn tương ứng để chỉ ra trạng thái đeo khẩu trang của người đó: "Mask" (có khẩu trang), "Wearing Improperly" (đeo khẩu trang không đúng cách), hoặc "No Mask" (không đeo khẩu trang).

3. Về Dataset:

- Bộ dữ liệu bao gồm tổng cộng 874 tấm ảnh, mỗi tấm ảnh chứa 1 hoặc nhiều khuôn mặt người, mỗi khuôn mặt người sẽ thuộc một trong ba class: Mask, No Mask, Wearing Improperly.
- Mỗi tấm ảnh sẽ có tương ứng một file text riêng, file này sẽ được đặt tên giống với tấm ảnh, nó chứa thông tin về các đối tượng được gắn nhãn trong ảnh, bao gồm vị trí của bounding box và nhãn của chúng.

<label_1> <x_min_1> <y_min_1> <width_1> <height_1>

<label_2> <x_min_2> <y_min_2> <width_2> <height_2>

...

- <label_i>: Nhãn của bounding box thứ i (0 là có đeo khẩu trang, 1 là không đeo khẩu trang, 2 là đeo không đúng cách).
 - <x_min_i>, <y_min_i>: Tọa độ của góc trái phía trên của bounding box thứ i, được biểu diễn bằng tỉ lệ so với chiều rộng và chiều cao của hình ảnh.
 - <width_i>, <height_i>: Chiều rộng và chiều cao của bounding box thứ i, được biểu diễn bằng tỉ lệ so với chiều rộng và chiều cao của hình ảnh.
- Bộ dữ liệu được chia train và test theo tỉ lệ 80:20. Ngoài ra, model còn được test qua một số video dài từ 10 giây đến 1 phút của nhóm.

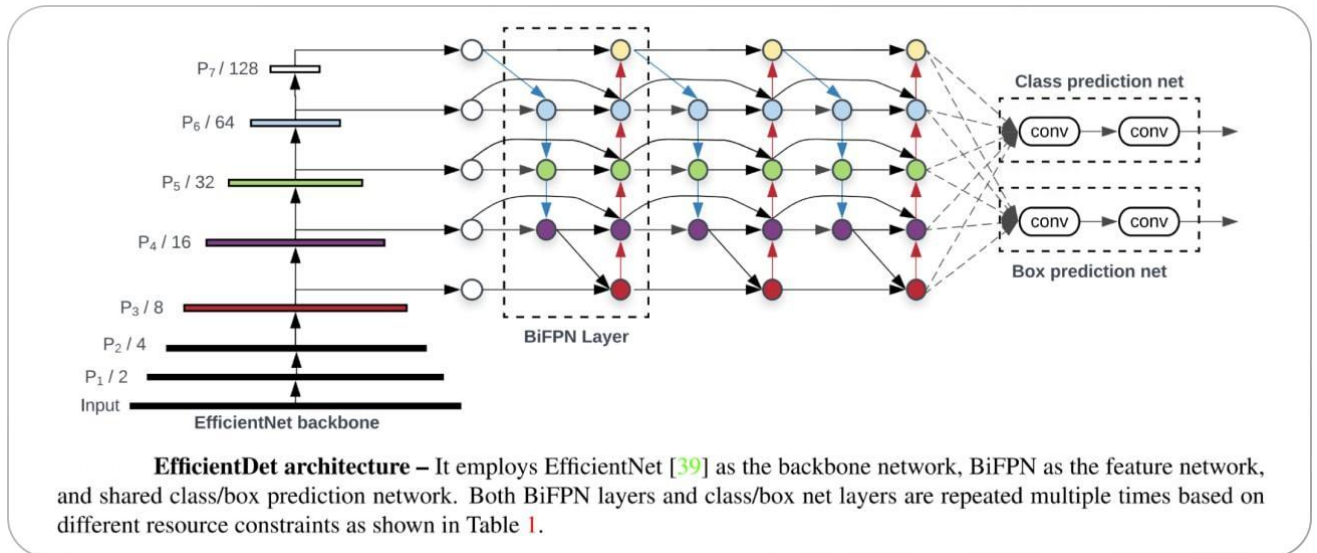
Chương 3: Cơ sở lý thuyết, training và đánh giá model

1. YOLO:

- "YOLO", viết tắt của "You Only Look Once", là một đại gia đình của các mô hình phát hiện đối tượng được giới thiệu bởi Joseph Redmon trong một bài báo năm 2016 có tiêu đề "You Only Look Once: Unified, Real-Time Object Detection".
- Kể từ đó, các biến thể của mô hình phát hiện đối tượng YOLO đã nhanh chóng phát triển, với việc phát hành phiên bản YOLO-v8 mới nhất vào tháng 1 năm 2023. Các biến thể của YOLO dựa trên nguyên tắc của hiệu suất phân loại thời gian thực và cao, dựa trên các tham số tính toán hạn chế nhưng hiệu quả.
- YOLO là một phương pháp phát hiện đối tượng thời gian thực trong ảnh và video bằng cách sử dụng một mạng neural network end-to-end để dự đoán các hộp giới hạn (bounding boxes) và xác suất lớp (class probabilities) cùng một lúc.

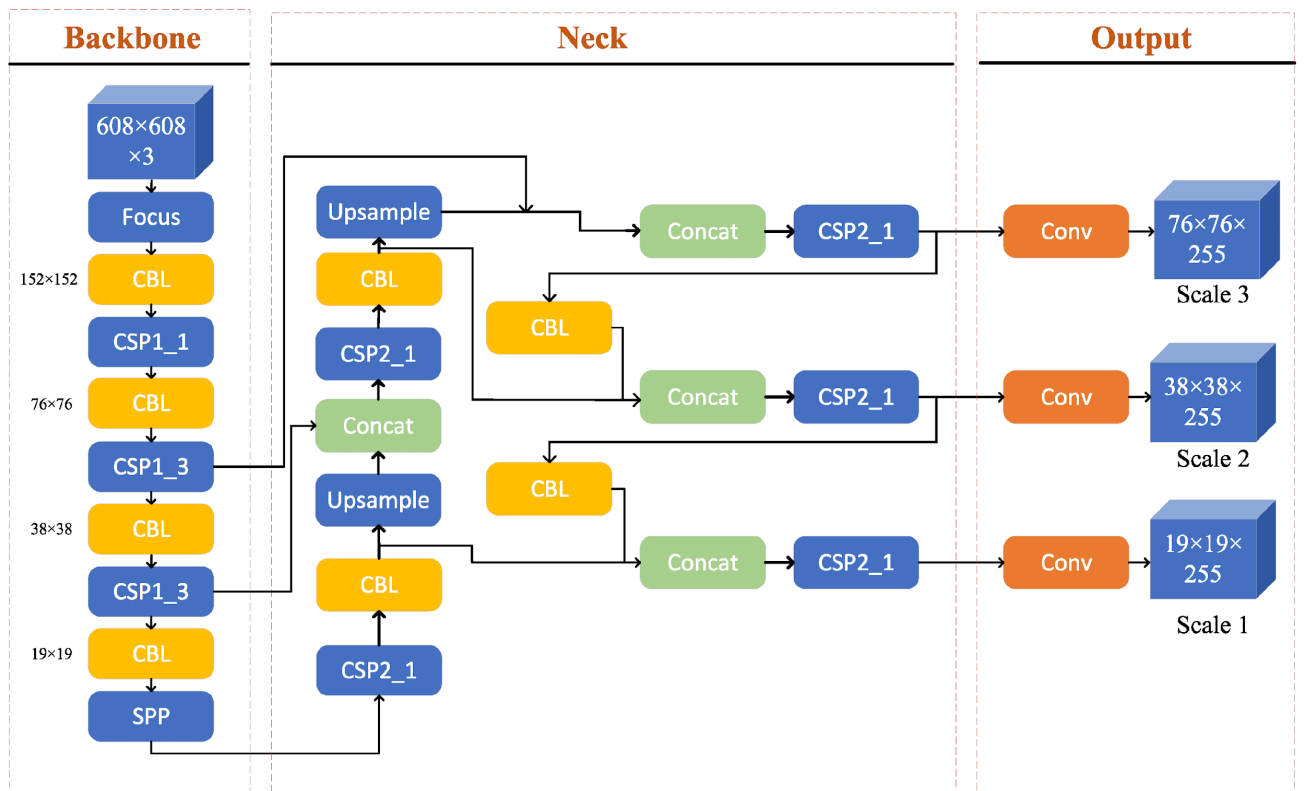
2. Model YOLOv5:

- YOLOv5 là một sự cải tiến lớn so với các phiên bản trước đó, với nhiều cải tiến về hiệu suất và tính linh hoạt.
- Mô hình YOLOv5 được thiết kế để đơn giản hóa cả quá trình huấn luyện và triển khai, với việc tối ưu hóa mô hình cho tốc độ và hiệu suất cao.
- YOLO v5 được giới thiệu vào năm 2020 bởi cùng một nhóm đã phát triển thuật toán YOLO dưới dạng một dự án mã nguồn mở và được duy trì bởi Ultralytics. YOLO v5 xây dựng dựa trên sự thành công của các phiên bản trước đó và thêm vào đó một số tính năng và cải tiến mới.
- Khác với YOLO, YOLO v5 sử dụng một kiến trúc phức tạp hơn được gọi là EfficientDet, dựa trên kiến trúc mạng EfficientNet. Việc sử dụng một kiến trúc phức tạp hơn trong YOLO v5 cho phép nó đạt được độ chính xác cao hơn và khả năng tổng quát hóa tốt hơn đối với một loạt các loại đối tượng.



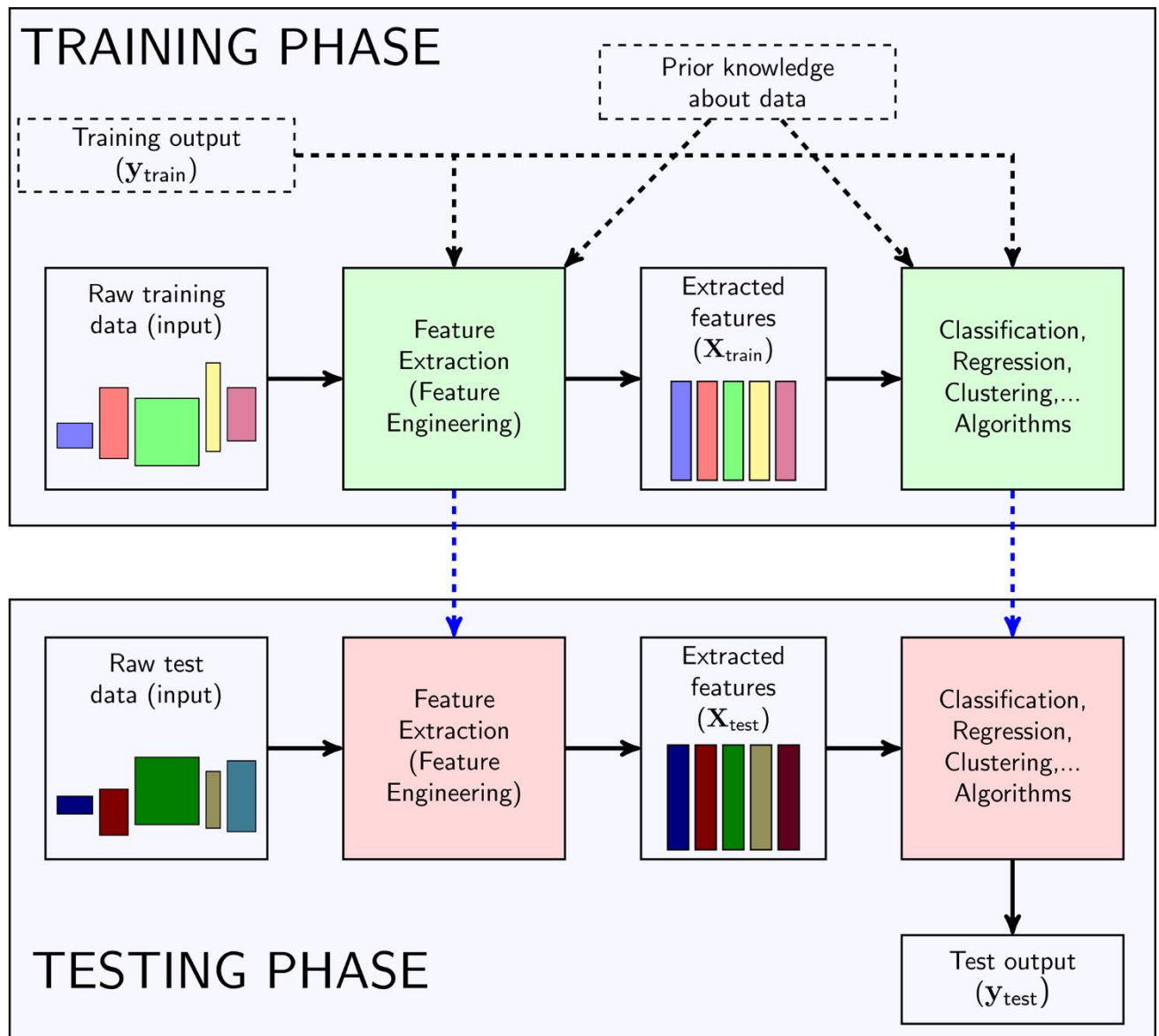
- Dưới đây là một phân tích chi tiết về sự khác biệt giữa kiến trúc YOLOv4 và YOLOv5:
 - Backbone (Lõi):
 - YOLOv4: Sử dụng CSPResidualBlock, một kiến trúc residual block có sự kết hợp của CSP (Cross-Stage Partial connections) và các residual connections.
 - YOLOv5: Sử dụng C3 module, một kiến trúc convolutional layer mới được thiết kế, có thể tối ưu hóa hiệu quả hơn so với CSPResidualBlock.
 - Neck (Cổ):
 - YOLOv4: Sử dụng SPP (Spatial Pyramid Pooling) và PAN (Path Aggregation Network) để kết hợp thông tin đa phân cấp và cải thiện mạng.
 - YOLOv5: Sử dụng SPPF (Fused Spatial Pyramid Pooling) và PAN, một biến thể của SPP kết hợp thông tin cụ thể từ các đối tượng nhỏ hơn và PAN để tối ưu hóa sự kết hợp thông tin.
 - Head (Đầu):

- YOLOv4: Giữ nguyên từ YOLOv3, sử dụng một loạt các convolutional layers để dự đoán bounding box và xác suất lớp.
- YOLOv5: Tiếp tục sử dụng đầu từ YOLOv3 mà không có thay đổi đáng kể.
- Các Thay Đổi Khác:
 - Data Augmentation (Tăng cường dữ liệu): YOLOv5 thêm vào các phương pháp tăng cường dữ liệu mới như Mosaic Augmentation, Copy-paste Augmentation, và MixUp Augmentation để cải thiện sự đa dạng và chất lượng của dữ liệu huấn luyện.
 - Loss Function (Hàm mất mát): YOLOv5 thêm hệ số scale cho Objectness Loss để cân bằng giữa các loại lỗi và cải thiện độ chính xác của mô hình.
 - Anchor Box (Hộp Neo): YOLOv5 sử dụng Auto Anchor sử dụng Genetic Algorithm để tự động tối ưu hóa các anchor box cho dữ liệu huấn luyện cụ thể.
 - Loại bỏ Grid Sensitivity: YOLOv5 loại bỏ Grid Sensitivity và thay đổi công thức để cải thiện độ nhạy của mô hình đối với các đối tượng nhỏ.
 - EMA Weight: Cân nhắc sử dụng Exponential Moving Average (EMA) Weight để cải thiện ổn định và hiệu suất của mô hình trong quá trình huấn luyện.
 - Những thay đổi này giúp YOLOv5 đạt được hiệu suất và linh hoạt tốt hơn so với các phiên bản tiền nhiệm, với khả năng huấn luyện và triển khai dễ dàng hơn.



YOLOv5 architecture

3. Giai đoạn training model:



- Từ tập dataset đã có từ trước, dữ liệu đó sẽ được tiền xử lý trước khi đưa vào training.
- Sau đó, chúng sẽ được trích xuất các đặc trưng tạo ra các feature map tương đương.
- Những feature map này cuối cùng được đưa vào những thuật toán máy học để phân loại.

4. Về quá trình trích xuất đặc trưng (feature engineering):

Trong YOLO, quá trình trích xuất đặc trưng diễn ra thông qua một deep neural network được huấn luyện trên dữ liệu ảnh.

- Tiền Xử Lý Ảnh: Trước khi được đưa vào mạng neural network, ảnh đầu vào được tiền xử lý để chuẩn hóa kích thước và giá trị pixel.
- Mạng Neural Network Convolutional (CNN): YOLO sử dụng một mạng neural network tích chập (CNN) để trích xuất các đặc trưng từ ảnh đầu vào. Mạng CNN này có thể bao gồm nhiều lớp convolutional, pooling, và các lớp kích hoạt phi tuyến tính như ReLU.
- Lớp Global Average Pooling (GAP): Sau khi đi qua các lớp convolutional, đầu ra của mạng CNN thường là một bản đồ đặc trưng (feature map) có kích thước lớn. Để giảm số lượng tham số và tính toán, YOLO sử dụng lớp GAP để chuyển đổi mỗi feature map thành một giá trị đặc trưng duy nhất bằng cách tính trung bình của tất cả các giá trị pixel trên feature map.
- Lớp Fully Connected Layer (FC): Sau lớp GAP, các giá trị đặc trưng được đưa vào một hoặc nhiều lớp fully connected layer để ánh xạ từ không gian đặc trưng sang không gian dự đoán. Các lớp này có thể liên kết các đặc trưng với các lớp phân loại và dự đoán các bounding box.
- Lớp Output: Cuối cùng, các lớp output của mạng neural network sẽ dự đoán các bounding box và xác suất của các lớp đối tượng trong ảnh. Các bounding box thường được dự đoán bằng cách sử dụng các hàm kích hoạt như sigmoid để dự đoán tọa độ (vị trí và kích thước) của bounding box, cùng với các xác suất của các lớp đối tượng thông qua softmax hoặc sigmoid.

5. Lý do chọn model:

Lý do chúng em chọn mô hình YOLOv5 cho đề án này là dựa trên các ưu điểm sau:

- Hiệu Suất Cao: YOLOv5 là một trong những mô hình object detection hiện đại nhất, với hiệu suất cao và tốc độ xử lý nhanh, phù hợp với yêu cầu của dự án.
- Dễ Triển Khai: Mô hình YOLOv5 có cấu trúc đơn giản và dễ dàng triển khai, giúp giảm bớt thời gian và công sức trong việc triển khai hệ thống.
- Tính Linh Hoạt: YOLOv5 hỗ trợ các ứng dụng trên nhiều nền tảng và môi trường khác nhau, từ các ứng dụng di động đến các hệ thống nhúng, giúp dễ dàng tích hợp vào các ứng dụng thực tế.

6. Quá trình train model:

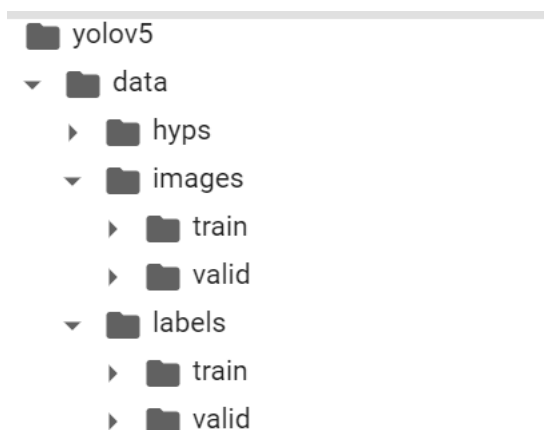
- Upload bộ dữ liệu thu thập được lên google drive.
- Tải các thư viện và framework cần thiết để train model.

```
git clone https://github.com/ultralytics/yolov5
```

```
cd yolov5
```

```
pip install -r requirements.txt
```

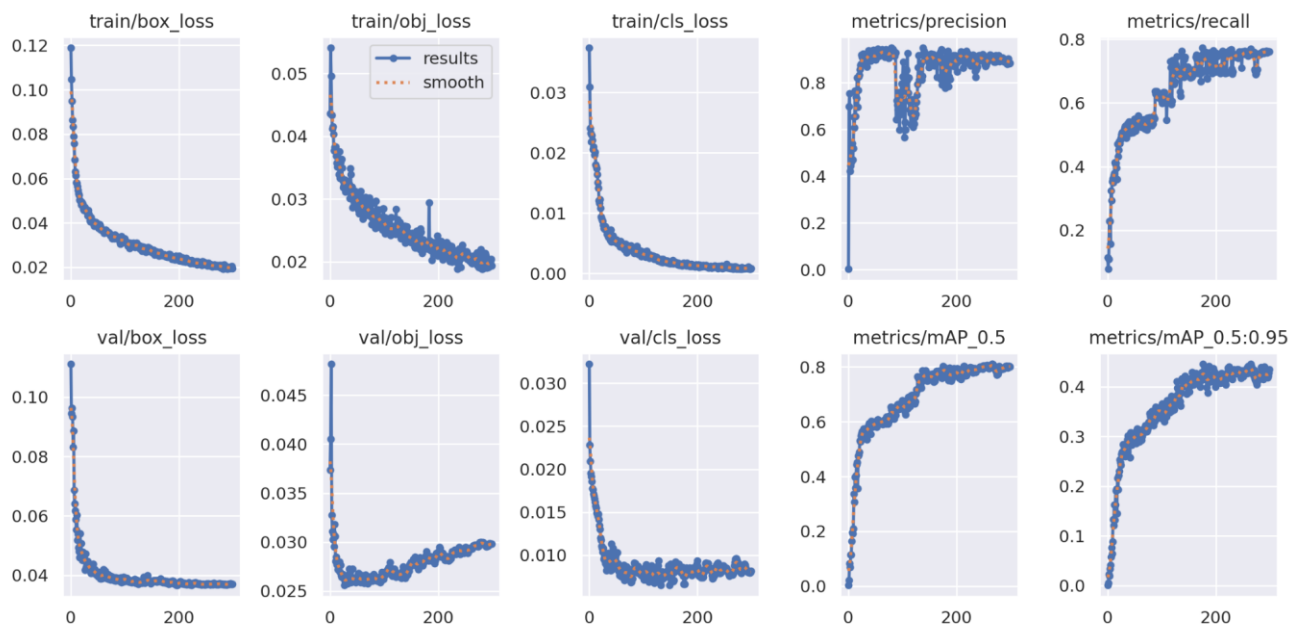
- Giải nén bộ dữ liệu, chia bộ dữ liệu thành 2 tập train và test. Tỷ lệ là 80:20. Kiểm tra lại sau khi giải nén thì folder images/ và labels/ phải nằm trong yolov5/data như sau:



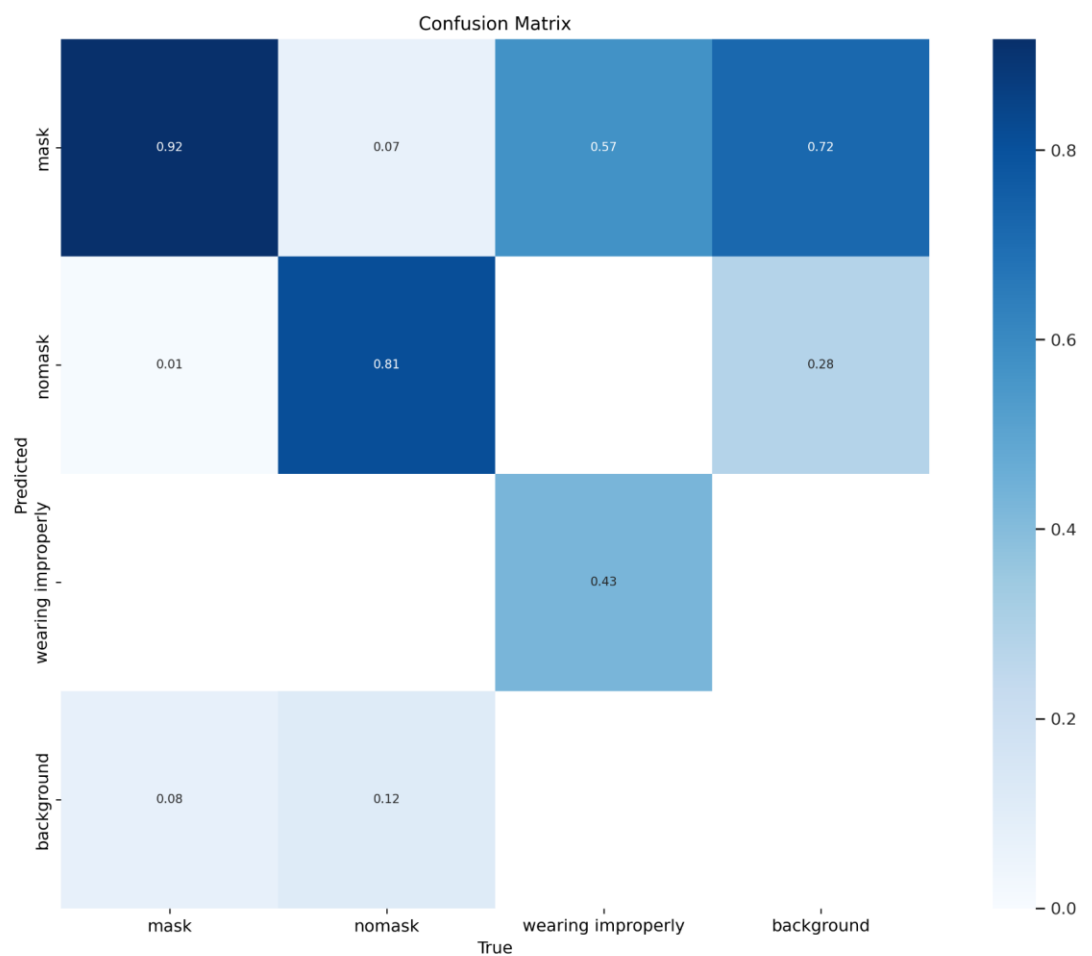
- Để huấn luyện mô hình, chúng em chạy train.py, với các đối số như sau:
 - img: kích thước ảnh đầu vào – 416.
 - batch: kích thước batch – 16.
 - epochs: số lượng epochs – 300.
 - data: đường dẫn đến tệp dataset.yaml
 - weights: đường dẫn trọng số ban đầu, mặc định là yolov5s.pt

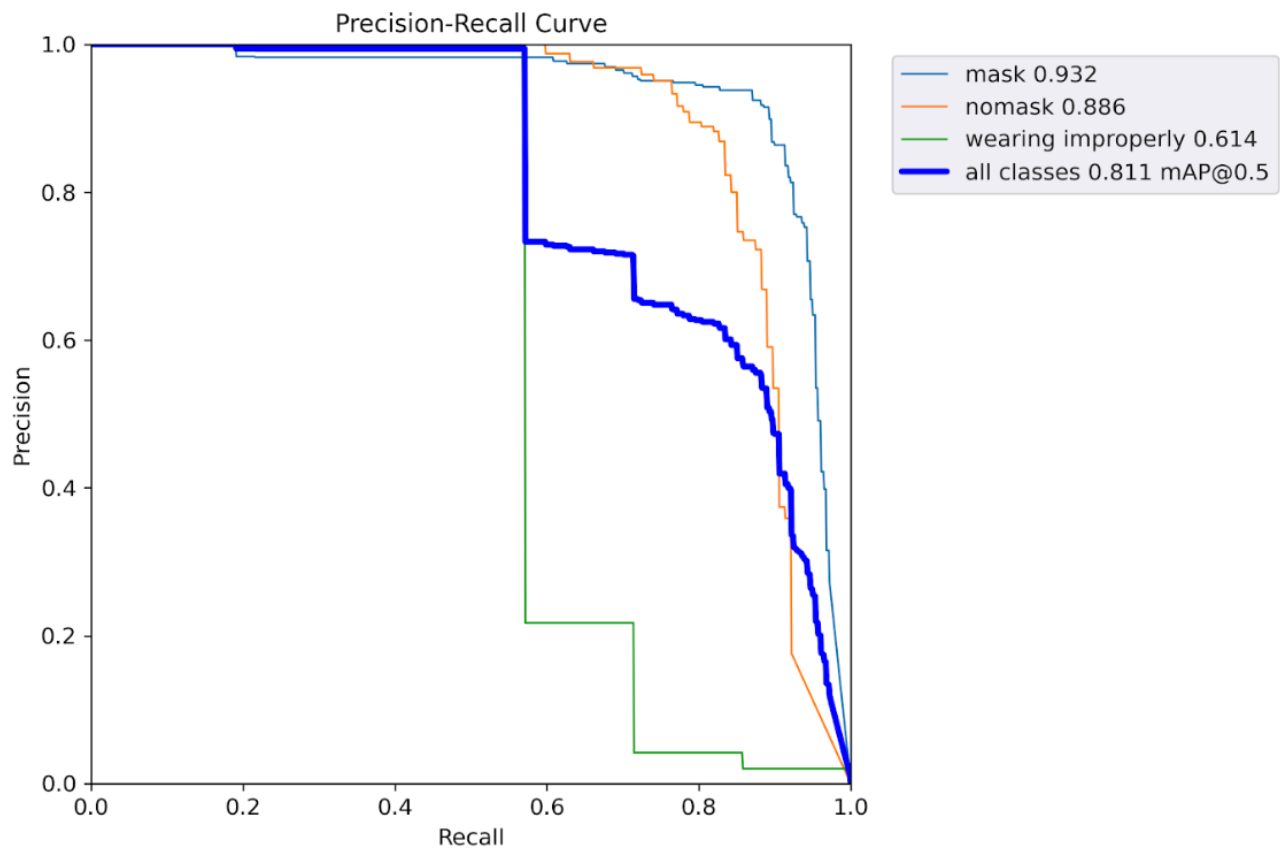
7. Kết quả và đánh giá:

(training loss)



(Confusion Matrix)





(đánh giá IoU và mAP)

Class	Images	Instances	P	R	mAP50	mAP50-95:
all	78	412	0.924	0.782	0.837	0.494
mask	78	278	0.87	0.914	0.923	0.541
nomask	78	127	0.902	0.865	0.903	0.468
wearing improperly	78	7	1	0.566	0.684	0.473

Kết luận chung:

- mAP50 (0,837) tương đối cao, cho thấy độ chính xác IoU tổng thể tốt.
- mAP50-95 (0,494) thấp hơn, cho thấy cần cải thiện trong việc xử lý kích thước đối tượng khác nhau và ngưỡng IoU.

Hiệu suất theo lớp:

- With mask: Hiệu suất tốt nhất với mAP50 là 0,923, cho thấy độ chính xác cao trong việc phát hiện khẩu trang được đeo đúng cách.
- Without mask: Đạt được mAP50 là 0,903, thể hiện hiệu suất mạnh mẽ trong việc xác định các cá nhân không đeo khẩu trang.

- Wear incorrectly: Có mAP50 thấp nhất là 0,473, do dữ liệu bị mất cân bằng tập trung vào 2 dữ liệu ở trên

8. Thực hiện dự đoán:

Để thực hiện dự đoán trên ảnh và luồng video, chúng em gọi detect.py và điều chỉnh các tham số sau:

- weights: trọng số của mô hình đã huấn luyện
- source: tệp/thư mục đầu vào để thực hiện dự đoán, 0 để sử dụng webcam
- output: thư mục để lưu kết quả
- conf-thres: ngưỡng độ tin cậy của đối tượng





Nhận xét:

- Mô hình nhìn chung hoạt động khá tốt trong việc phát hiện người đeo hoặc không đeo khẩu trang trên ảnh và video.
- Mô hình chưa phát hiện tốt các trường hợp đeo khẩu trang không đúng cách, điều này là do bộ dữ liệu chưa được cân bằng số lượng giữa các nhãn, nhãn “Mask” và “No mask” có số lượng lớn hơn nhiều so với nhãn “Wear Improperly”. Sự đa dạng trong dữ liệu huấn luyện có thể cần được tăng cường thêm để cải thiện hiệu suất của mô hình.
- Mô hình gặp khó khăn trong việc phát hiện khẩu trang dưới một số điều kiện cụ thể, ví dụ, nó có xu hướng nhầm lẫn giữa râu dài và một chiếc khẩu trang. Điều này có thể được giảm thiểu bằng cách thêm nhiều sự phong phú trong bộ dữ liệu huấn luyện và có thể thực hiện thêm các kỹ thuật data augmentation.
- Có thể cần thực hiện các bước tinh chỉnh mô hình và thử nghiệm với các siêu tham số khác nhau để cải thiện hiệu suất và đảm bảo rằng mô hình hoạt động hiệu quả trong các điều kiện đa dạng.

TÀI LIỆU THAM KHẢO:

1. <https://towardsdatascience.com/face-mask-detection-using-yolov5-3734ca0d60d8>
2. <https://github.com/iAmEthanMai/mask-detection-dataset>
3. <https://colab.research.google.com/github/ultralytics/yolov5/blob/master/tutorial.ipynb#scrollTo=N3qM6T0W53gh>
4. [You Only Look Once: Unified, Real-Time Object Detection](#)
5. <https://medium.com/analytics-vidhya/covid-19-face-mask-detection-using-yolov5-8687e5942c81>
6. <https://github.com/spacewalk01/yolov5-face-mask-detection>
7. <https://www.v7labs.com/blog/yolo-object-detection#:~:text=YOLO%20v5%20uses%20a%20new,clusters%20as%20the%20anchor%20boxes.>
8. <https://blog.roboflow.com/yolov5-improvements-and-evaluation/>
9. <https://iq.opengenus.org/yolov5/>
10. <https://viblo.asia/p/tong-hop-kien-thuc-tu-yolov1-den-yolov5-phan-1-naQZRRj0Zvx>