

Nhận Diện Khuôn Mặt Kết Hợp Với Khẩu Trang

Luong Anh Vũ

Vu Ngoc Tien

Lê Nguyễn Khánh Tùng

nluonganhvu.12a5.tnd@gmail.com

vungoctienanh@gmail.com

ktung.photograph@gmail.com

Abstract—Nhận diện khuôn mặt trong điều kiện đeo khẩu trang là bài toán phức tạp do việc che khuất các đặc trưng quan trọng của khuôn mặt. Bài báo này trình bày tổng quan về các phương pháp nhận diện khuôn mặt, từ các kỹ thuật truyền thống đến các mô hình học sâu hiện đại, nhằm cải thiện hiệu suất nhận diện trong điều kiện khó khăn. Hệ thống được đề xuất kết hợp các kỹ thuật Ensemble và xử lý hậu kỳ nhằm tối ưu hóa kết quả, với tỉ lệ thành công lên đến 98% trên một số tập dữ liệu kiểm tra.

I. Giới thiệu

Nhận diện khuôn mặt là một trong những ứng dụng quan trọng của trí tuệ nhân tạo và thị giác máy tính, được sử dụng rộng rãi trong an ninh, giám sát, kiểm soát truy cập và các ứng dụng cá nhân hóa. Tuy nhiên, trong bối cảnh đại dịch COVID-19 và sự phổ biến của việc đeo khẩu trang, các hệ thống nhận diện khuôn mặt truyền thống gặp nhiều thách thức do phần lớn khuôn mặt bị che khuất. Điều này làm giảm đáng kể độ chính xác của mô hình nhận diện, đòi hỏi các phương pháp tiếp cận mới để xử lý vấn đề này một cách hiệu quả. Bài báo này tập trung vào nghiên cứu các phương pháp nhận diện khuôn mặt khi đeo khẩu trang, bao gồm việc xây dựng mô hình học sâu (Deep Learning) phù hợp, sử dụng các tập dữ liệu chuyên biệt, và đề xuất các giải pháp cải thiện độ chính xác của hệ thống. Chúng tôi khảo sát các kỹ thuật hiện đại như nhận diện đặc trưng không phụ thuộc vào vùng bị che khuất, sử dụng mô hình học sâu có khả năng thích nghi với dữ liệu thiếu hụt, và áp dụng các chiến lược tiền xử lý dữ liệu nhằm tối ưu hóa hiệu suất nhận diện. Bên cạnh đó, chúng tôi cũng tiến hành thực nghiệm trên một tập dữ liệu chứa hình ảnh khuôn mặt có và không có khẩu trang để đánh giá hiệu quả của mô hình. Kết quả thực nghiệm cho thấy rằng các kỹ thuật xử lý và mô hình học sâu có thể cải thiện đáng kể khả năng nhận diện khuôn mặt trong điều kiện đeo khẩu trang. Động lực Tiến bộ công nghệ: Học sâu: Sự trỗi dậy của các mạng nơ-ron sâu (DNN) đã cách mạng hóa lĩnh vực nhận diện khuôn mặt. Các mô hình như CNN và Transformer đã đạt được độ chính xác chưa từng có, ngay cả trong các điều kiện thách thức. Dữ liệu lớn: Sự sẵn có của lượng dữ liệu hình ảnh và video khổng lồ đã cho phép các mô hình AI học các mẫu phức tạp và cải thiện độ chính xác. Phần cứng mạnh

mẽ: Sự phát triển của GPU và các bộ xử lý chuyên dụng khác đã cung cấp sức mạnh tính toán cần thiết để huấn luyện và triển khai các mô hình nhận diện khuôn mặt phức tạp. Ứng dụng rộng rãi: An ninh: Nhận diện khuôn mặt được sử dụng trong các hệ thống giám sát, kiểm soát truy cập và xác minh danh tính để tăng cường an ninh. Xác thực: Công nghệ này được sử dụng để xác thực người dùng trong các ứng dụng ngân hàng, thanh toán và các dịch vụ trực tuyến khác. Tiếp thị: Nhận diện khuôn mặt có thể được sử dụng để phân tích nhân khẩu học khách hàng, cá nhân hóa trải nghiệm mua sắm và đo lường hiệu quả của các chiến dịch quảng cáo. Y tế: Công nghệ này có thể được sử dụng để phát hiện các bệnh di truyền, theo dõi sự tiến triển của bệnh và hỗ trợ chẩn đoán. Giải trí: Nhận diện khuôn mặt được sử dụng trong các ứng dụng thực tế tăng cường, trò chơi điện tử và các nền tảng truyền thông xã hội để tạo ra trải nghiệm tương tác.

II. Các phương pháp nhận diện khuôn mặt

A. Phương pháp truyền thống

Các phương pháp truyền thống dựa trên việc trích xuất các đặc trưng cơ bản và giảm chiều dữ liệu:

Eigenfaces (PCA) Sử dụng phân tích thành phần chính để chuyển đổi không gian dữ liệu và trích xuất các đặc trưng nổi bật. Phương pháp này nhanh nhưng nhạy cảm với ánh sáng và góc chụp.

Fisherfaces (LDA) Áp dụng phân tích tuyến tính phân biệt để tìm ra các đặc trưng tối ưu giữa các lớp khuôn mặt, cải thiện so với PCA nhưng gặp hạn chế khi khuôn mặt bị che.

Local Binary Pattern Histogram (LBPH) Trích xuất đặc trưng cục bộ qua biểu diễn nhị phân, thích hợp với ứng dụng thời gian thực nhưng hiệu quả giảm khi khuôn mặt bị che.

B. Phương pháp học sâu

Các mô hình học sâu cho phép tự động học các đặc trưng phức tạp:

CNN truyền thống Sử dụng các kiến trúc như VGG-Face, ResNet, Inception để trích xuất đặc trưng. Mô hình này cho hiệu suất cao nhưng yêu cầu tính toán lớn.

Transformer Sử dụng cơ chế Attention để tập trung vào các vùng không bị che, cải thiện khả năng nhận diện trong điều kiện khó khăn.

Ensemble Kết hợp nhiều mô hình giúp tăng tính ổn định và độ chính xác tổng thể.

[2]

III. Kiến trúc và Quy trình hoạt động của CNN

A. Kiến trúc của CNN

Một mạng CNN điển hình bao gồm các thành phần:

1) Lớp Tích Chập: Các hạt nhân (kernel) có kích thước nhỏ (3×3 , 5×5) quét qua ảnh để trích xuất các đặc trưng cục bộ như cạnh, góc và kết cấu. Công thức:

$$Z = (X * W) + B$$

với X là ảnh đầu vào, W là trọng số và B là bias.

2) Lớp Kích Hoạt: Sử dụng hàm kích hoạt ReLU: $f(x) = \max(0, x)$ để giới thiệu tính phi tuyến.

3) Lớp Pooling: Giảm kích thước bản đồ đặc trưng, giảm số lượng tham số và tăng tốc độ xử lý. Phổ biến là Max Pooling và Average Pooling.

4) Lớp Kết Nối Toàn Bộ: Sau khi làm phẳng các đặc trưng, lớp kết nối toàn bộ dùng để phân loại cuối cùng qua hàm Softmax.

Quy trình nhận diện Từ ảnh đầu vào, hệ thống thực hiện các bước: chụp ảnh, phát hiện khuôn mặt, trích xuất đặc trưng và phân loại.

[3]

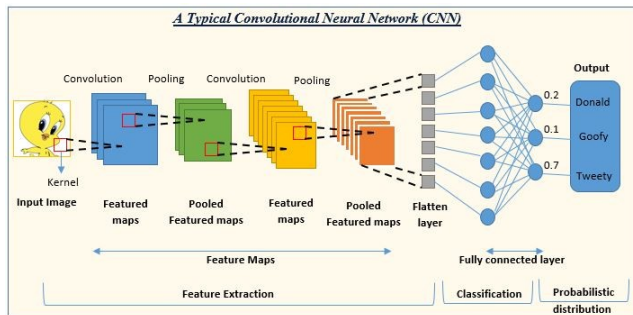


Fig. 1. Mô hình CNN điển hình: Từ ảnh đầu vào (Tweety) qua các lớp Convolution, Pooling, Flatten và Fully Connected, cuối cùng đưa ra phân loại (Tweety, Donald, Goofy).

Như trong Hình 1, quá trình xử lý ảnh đầu vào gồm:

B. Tổng quan Hình minh hoạ

Hình này thể hiện một mạng CNN điển hình (A Typical Convolutional Neural Network) gồm các khối chính:

Convolution (Tích chập)

Pooling (Giảm kích thước)

Flatten (Làm phẳng dữ liệu)

Fully connected layer (Lớp kết nối toàn bộ)

Output (Đầu ra xác suất cho các lớp)

Ảnh đầu vào là một hình (chẳng hạn Tweety), qua quá trình xử lý sẽ được dự đoán lớp đầu ra (Tweety, Donald, Goofy) với các xác suất tương ứng.

C. 2. Quy trình chi tiết

1Input Image (Ảnh đầu vào)

Ở đây là hình một nhân vật (Tweety). CNN sẽ lấy ảnh này làm đầu vào.

Mỗi pixel của ảnh được xem như giá trị (hoặc nhiều kênh màu nếu là RGB).

2Convolution (Tích chập)

Mỗi khối Convolution sử dụng một hoặc nhiều \textit{kernel} (bộ lọc) kích thước nhỏ (3×3 , 5×5 ...) để “quét” toàn bộ ảnh đầu vào.

Kết quả là các \textit{Feature maps} (bản đồ đặc trưng), thể hiện các đặc trưng cục bộ (cạnh, góc, họa tiết...).

3Pooling (Giảm kích thước)

Sau khi có \textit{Feature maps}, CNN áp dụng Pooling (thường là Max Pooling) để giảm độ phân giải, giảm số tham số và hạn chế overfitting.

Kết quả là các \textit{Pooled Feature maps}, kích thước nhỏ hơn nhưng vẫn giữ lại thông tin chính.

4Convolution + Pooling (lặp lại)

Tùy kiến trúc CNN, ta có thể lặp lại nhiều khối Convolution + Pooling để trích xuất đặc trưng sâu hơn.

Ở hình minh hoạ, có hai khối Convolution-Pooling liên tiếp.

5Flatten layer (Làm phẳng)

Sau khối Convolution - Pooling cuối, dữ liệu (lúc này là nhiều \textit{feature maps} 2D) được “làm phẳng” thành vector 1 chiều.

Mục đích: Chuẩn bị cho bước phân loại (Fully connected layer).

6Fully connected layer (Lớp kết nối toàn bộ)

Lớp này nhận vector đầu vào từ bước Flatten.

Các trọng số (weights) của mạng kết nối mọi neuron, học cách kết hợp các đặc trưng để phân loại chính xác.

7Output (Đầu ra)

Cuối cùng, mạng xuất xác suất (\textit{Probabilistic distribution}) cho mỗi lớp.

Ví dụ: 0.2 (Donald), 0.1 (Goofy), 0.7 (Tweety).

Lớp có xác suất cao nhất (Tweety = 0.7) chính là kết quả dự đoán.

D. 3. Ý nghĩa của các khối

Convolution: Trích xuất đặc trưng (cạnh, góc, đường viền...).

Pooling: Giảm độ phân giải, giảm số tham số.

Flatten: Chuyển dữ liệu 2D (hoặc 3D) thành vector 1D để đưa vào mạng nơ-ron truyền thẳng.

Fully Connected: Tổng hợp toàn bộ đặc trưng, đưa ra dự đoán cuối cùng.

Output: Xác suất của các lớp, dùng hàm Softmax (đa lớp) hoặc Sigmoid (nhị phân).

E. 4. Liên hệ với bài toán Nhận diện Khuôn mặt

Mặc dù hình minh hoạ là ví dụ nhận diện các nhân vật (Tweety, Donald, Goofy), quy trình không thay đổi nếu thay ảnh Tweety bằng ảnh khuôn mặt người. Mạng CNN vẫn trích xuất đặc trưng, pooling, flatten, rồi phân loại (thay vì “Tweety/Donald/Goofy” thì là danh tính các khuôn mặt).

Trong bối cảnh đeo khẩu trang, vùng miệng-mũi bị che khuất, nên CNN cần học cách “bỏ qua” hoặc giảm trọng số cho vùng bị che, tập trung hơn vào mắt, trán, v.v.

Một số cải tiến: Dùng Attention (Transformer) hoặc thêm lớp “mask detection” hỗ trợ.

F. 5. Tóm tắt

Hình minh hoạ trình bày \textbf{quy trình hoạt động của một CNN cơ bản}:

- 1\textit{Input Image}
- 2\textit{Convolution}
- 3\textit{Pooling}
- 4\textit{Flatten}
- 5\textit{Fully Connected Layer} \rightarrow
- 6\textit{Output} (phân bố xác suất).

Nhờ quy trình này, CNN tự động học các đặc trưng từ ảnh, từ mức thấp (cạnh, góc) đến mức cao (cấu trúc tổng thể), cuối cùng phân loại đối tượng (hoặc khuôn mặt) một cách hiệu quả.

G. Mô hình Cascaded CNN (Ví dụ MTCNN)

P-Net (Proposal Network) Input size 12×12 : Ảnh đầu vào được cắt và co về kích thước 12×12 . Các lớp Conv 3×3 và Pooling: Mục tiêu là trích xuất các đặc trưng ban đầu từ ảnh khuôn mặt (hoặc vùng nghi ngờ là khuôn mặt). Đầu ra: Face classification: Xác suất vùng đó có phải là khuôn mặt hay không. Bounding box regression: Dự đoán tọa độ khung bao để hiệu chỉnh (refine) lại vị trí và kích thước vùng khuôn mặt. Facial landmark localization: Xác định vị trí tương đối của các điểm đặc trưng (mắt, mũi, miệng, v.v.). Chức năng chính: P-Net làm nhiệm vụ “quét” nhanh toàn bộ ảnh ở nhiều tỉ lệ (image pyramid) để đưa ra các đề xuất (proposal) khung bao khuôn mặt. Sau đó, các khung bao có độ tin cậy thấp sẽ bị loại bỏ.

R-Net (Refine Network) Input size 24×24 : Các vùng được P-Net “đề xuất” sẽ được cắt, co về kích thước 24×24 , rồi đưa vào R-Net. Các lớp Conv 3×3 và Pooling: Ở bước này, mạng R-Net có độ phức tạp cao hơn một chút so với P-Net, cho phép “tinh chỉnh” chính xác hơn các đặc trưng khuôn mặt. Đầu ra: Face classification: Xác suất khuôn mặt chính xác hơn, loại bỏ thêm các “khung bao nhiễu”. Bounding box regression: Tiếp tục điều chỉnh lại tọa độ khung bao cho chính xác. Facial landmark localization: Cập nhật vị trí các mốc đặc trưng khuôn mặt (nếu cần).

O-Net (Output Network) Input size 48×48 : Các vùng được R-Net xác nhận là có khả năng cao là khuôn mặt sẽ được đưa vào O-Net ở kích thước 48×48 . Các lớp Conv

3×3 , Pooling, Fully Connected: Mạng này sâu hơn và có khả năng học đặc trưng chi tiết hơn hẳn. Đầu ra: Face classification: Mức độ chắc chắn cuối cùng về việc vùng ảnh là khuôn mặt. Bounding box regression: Điều chỉnh tối ưu khung bao khuôn mặt. Facial landmark localization: Dự đoán chính xác tọa độ các điểm mốc (2 mắt, mũi, 2 khoé miệng).

Nhờ thiết kế theo từng tầng (cascade), mô hình vừa đảm bảo tốc độ (vì P-Net chạy rất nhanh) vừa nâng cao độ chính xác (nhờ R-Net và O-Net liên tục tinh chỉnh và xác nhận).

IV. Hệ thống nhận diện khuôn mặt đeo khẩu trang

A. Tiền xử lý và Phát hiện khuôn mặt

Chụp ảnh Sử dụng các loại camera (thường và hồng ngoại) để thu thập ảnh khuôn mặt. Việc tạo bản đồ chiều sâu (3D Facial Mapping) được áp dụng để nâng cao chất lượng ảnh.

Phát hiện khuôn mặt Các mô hình như MTCNN hoặc YOLO được dùng để xác định vị trí khuôn mặt, kể cả khi một phần khuôn mặt bị che bởi khẩu trang.

B. Trích xuất đặc trưng và Nhận diện

Trích xuất đặc trưng Sử dụng các mô hình học sâu như FaceNet, DeepFace hoặc MobileNetV2 để chuyển đổi khuôn mặt thành vector đặc trưng.

So sánh và Phân loại Các vector đặc trưng được so sánh với cơ sở dữ liệu thông qua các phép đo khoảng cách (Euclidean, Cosine Similarity) để đưa ra kết quả cuối cùng. [4]

V. Mô hình Ensemble và Xử lý Hậu kỳ

A. Mô hình Ensemble với Xử lý Nhiều Góc (Multi-View Processing)

Khi có nhiều góc chụp từ các camera khác nhau, các kết quả nhận diện từ từng góc được kết hợp theo công thức:

$$Z = \frac{1}{v} \sum_{i=1}^v E_i$$

với v là số góc chụp và E_i là kết quả nhận diện từ góc thứ i . Phương pháp này giúp giảm thiểu sai số và tăng tính ổn định của hệ thống.

B. Xử lý Hậu kỳ

Áp dụng bộ lọc làm mịn (mean filter) để giảm nhiễu trong xác suất nhận diện. Sử dụng kỹ thuật ngưỡng (thresholding) để loại bỏ các dự đoán không đáng tin cậy. Kết hợp kết quả từ nhiều mô hình thông qua kỹ thuật Ensemble nhằm tối ưu hóa kết quả cuối cùng.

C. Phân tích Độ phức tạp

Độ phức tạp của hệ thống nhận diện phụ thuộc vào số lượng mô hình kết hợp và số góc chụp, được biểu diễn bằng công thức:

$$O_{\text{ensemble}} = K \times v \times O(\text{model})$$

trong đó K là số fold trong kiểm định chéo, v là số góc chụp và $O(\text{model})$ là độ phức tạp của mô hình cơ bản.

VI. Method và Sai số

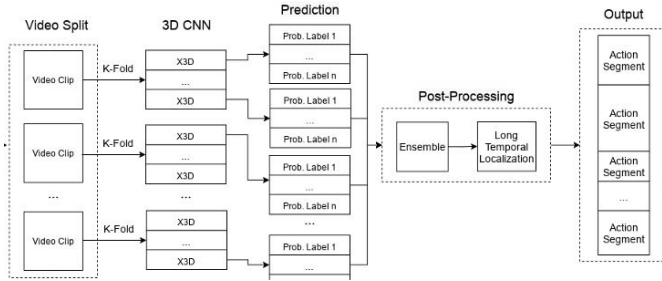


Fig. 2. Sơ đồ tổng quan quy trình xử lý video với 3D CNN (X3D) và K-Fold

A. Giải thích Sơ Đồ Phương Pháp

Video Split (Chia nhỏ video)

1) 1. Tách video (Chia nhỏ video):

Mục tiêu: Chia video gốc thành các đoạn nhỏ hơn để thuận tiện cho quá trình xử lý và huấn luyện mô hình.

Chi tiết:

Video đầu vào được chia thành nhiều video clip nhỏ.

Sử dụng phương pháp Xác thực chéo K-Fold để chia dữ liệu huấn luyện và kiểm tra.

Điều này giúp mô hình học được nhiều thông tin từ dữ liệu và giảm bớt quá mức hiện tượng.

2) 2. CNN 3D (Mạng nơ-ron tích chập 3D):

Mục tiêu: Trích xuất đặc biệt và thời gian từ các đoạn video.

Chi tiết:

Các video clip sau khi chia nhỏ được đưa vào mô hình X3D (Mở rộng 3D CNN), một mạng CNN 3D tối ưu hóa cho video.

CNN 3D có khả năng phân tích cả thông tin không gian (hình ảnh từng khung hình) và thời gian (mối quan hệ giữa các khung hình).

Một số tầng X3D được sử dụng để học các đặc trưng khác nhau.

3) 3. Dự đoán (Dự đoán hành động):

Mục tiêu: Gán nhãn xác thực cho các hành động có thể xảy ra trong video.

Chi tiết:

Sau khi qua CNN 3D, mỗi đoạn video sẽ có một tập hợp xác thực cho từng nhãn hành động có thể xảy ra.

Ví dụ, mô hình có thể dự đoán các hành động như "chạy", "nhảy", "đắm", vv

Các xác thực này được ghi lại dưới dạng Prob. Nhãn 1, Vấn đề. Nhãn n (xác suất của từng hành động).

4) 4. Xử lý hậu kỳ (Xử lý hậu kỳ):

Mục tiêu: Cải thiện độ chính xác mong đợi bằng cách kết hợp nhiều cấu hình và lọc kết quả.

Chi tiết:

Ensemble: Kết hợp nhiều mô hình khác nhau để đưa ra sự mong đợi chính xác hơn (có thể là trung bình hóa, bỏ phiếu hoặc xếp chồng).

Bản địa hóa theo thời gian dài: Áp dụng kỹ thuật xác định hành động theo khoảng thời gian dài, giúp mô hình hiểu rõ hơn về diễn đàn tiếp theo của hành động.

Điều này giúp loại bỏ các dự đoán sai lệch và tạo ra kết quả mượt mà hơn.

5) 5. Đầu ra (Kết quả đầu ra):

Mục tiêu: Tạo danh sách các hành động được phát hiện trong video khi bắt đầu và kết thúc thời gian.

Chi tiết:

Sau khi xử lý hậu kỳ, hệ thống tạo ra danh sách các Phân đoạn hành động (phân đoạn hành động).

Mỗi Phân đoạn hành động đều có thông tin về hành động xảy ra và thời điểm diễn ra trong video.

Ví dụ:

"Nhảy" (từ giây 10 đến giây 15)

"Chạy" (từ giây 20 đến giây 30)

6) Tóm tắt quy trình tổng thể::

1 Chia video thành các đoạn nhỏ bằng K-Fold.

2 Sử dụng CNN 3D để trích xuất cụ thể.

3 Expected performance cho từng hành động.

4 Kết quả mô phỏng và tối ưu hóa bằng bộ xử lý hậu kỳ.

5 Xuất kết quả, liệt kê các hành động được phát hiện.

Output (Xuất kết quả phân đoạn hành động) Hệ thống xuất ra danh sách các Action Segments, mỗi đoạn tương ứng với một hành động được nhận diện trong video.

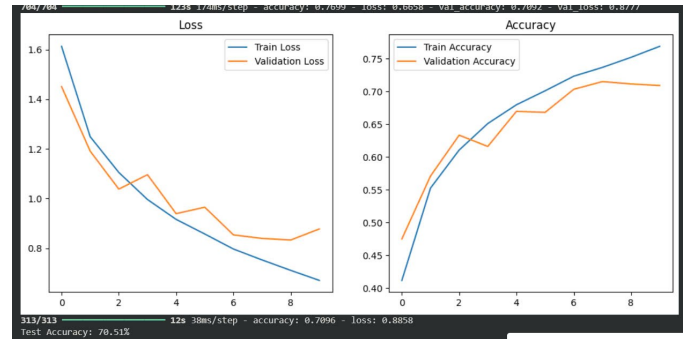


Fig. 3. Biểu đồ thể hiện sai số trong quá trình huấn luyện và kiểm thử

B. Phân Tích Sai Số

Hình 3 mô tả sai số (error) trong quá trình huấn luyện và kiểm thử mô hình. Đường cong huấn luyện (training error) thường có xu hướng giảm dần khi số epoch tăng, do mô hình dần học được các đặc trưng.

1. Quan sát chung Training Loss (màu xanh) vs. Validation Loss (màu cam) Đường huấn luyện (training loss) có xu hướng giảm dần khi số epoch tăng, thể hiện mô hình dần học được các đặc trưng của dữ liệu huấn luyện. Đường

đánh giá (validation loss) cũng giảm tương tự, nhưng có thể dao động hoặc giảm chậm hơn. Đến những epoch cuối, ta thấy validation loss không còn giảm nhiều, thậm chí có dấu hiệu tăng nhẹ (overfitting nhẹ).

Training Accuracy (màu xanh) vs. Validation Accuracy (màu cam) Training accuracy tăng lên rõ rệt qua từng epoch, thể hiện mô hình đang “thuộc” dữ liệu huấn luyện tốt hơn. Validation accuracy cũng tăng, nhưng thường thấp hơn hoặc tăng chậm hơn so với training accuracy, cho thấy mô hình có thể bắt đầu ghi nhớ đặc trưng của tập huấn luyện nhiều hơn là khái quát tốt cho dữ liệu mới.

2. Giải thích chi tiết về sai số Training Loss giảm đều Khi mô hình “nhìn thấy” nhiều lần dữ liệu huấn luyện, trọng số (weights) và tham số (biases) được điều chỉnh dần để giảm sai số. Điều này phản ánh việc mô hình đang học tốt các mẫu huấn luyện.

Validation Loss giảm nhưng dao động Validation loss giảm nghĩa là mô hình cũng dần học được cách nhận diện trên dữ liệu chưa “nhìn thấy” (tập validation). Tuy nhiên, nếu validation loss bắt đầu tăng trong khi training loss vẫn giảm, đó là dấu hiệu mô hình đang overfit — tức học quá sâu vào chi tiết của tập huấn luyện.

Sai số ở những epoch cuối Đến khoảng epoch 7 hoặc 8, validation loss có thể “chững lại” hoặc tăng nhẹ, còn training loss vẫn giảm. Đây là điểm cần lưu ý để áp dụng các kỹ thuật chống overfitting (như early stopping, regularization, dropout...).

3. Phân tích độ chính xác (Accuracy) Training Accuracy tăng mạnh (màu xanh), chứng tỏ mô hình có khả năng “nhớ” tốt dữ liệu huấn luyện. Validation Accuracy (màu cam) cũng tăng, nhưng không tăng đều hoặc có thể thấp hơn so với training accuracy. Độ chính xác kiểm thử (test accuracy) đạt khoảng 70.51%, cho thấy mô hình vẫn chưa khái quát hoá tốt như mong muốn (có thể do dữ liệu phức tạp, chưa đủ lớn, hoặc do mô hình chưa tối ưu kiến trúc/hyperparameters).

4. Nguyên nhân và khuyến nghị Overfitting Khi training accuracy quá cao so với validation accuracy, mô hình có xu hướng ghi nhớ đặc trưng của tập huấn luyện thay vì học quy luật tổng quát. Nên áp dụng một số kỹ thuật như Dropout, Regularization (L2), Early Stopping để tránh overfitting.

Dữ liệu chưa đa dạng Nếu tập huấn luyện và tập kiểm thử không đủ đa dạng hoặc không đại diện, mô hình sẽ khó khái quát. Cần xem xét mở rộng hoặc cân bằng lại dữ liệu, đặc biệt với bài toán nhận diện khi đeo khẩu trang. Tối ưu hyperparameters Điều chỉnh learning rate, batch size, số lượng epoch để mô hình hội tụ tốt hơn. Thử nghiệm thêm các kiến trúc mạng khác nhau (ResNet, MobileNet, Transformer...) hoặc sử dụng mô hình Ensemble để nâng cao độ chính xác.

5. Kết luận về sai số Training Loss giảm dần và Training Accuracy tăng dần là dấu hiệu mô hình học tốt trên dữ liệu huấn luyện. Validation Loss và Validation Accuracy cho biết mức độ khái quát hoá. Khi validation loss không

giảm tiếp hoặc tăng lên, trong khi training loss tiếp tục giảm, ta đang đối mặt với hiện tượng overfitting. Test Accuracy 70.51% thể hiện mô hình vẫn có thể cải thiện thêm để áp dụng hiệu quả trong thực tế.

VII. Thử nghiệm và Đánh giá

A. Thiết lập Thử Nghiệm

Tập dữ liệu gồm hình ảnh khuôn mặt có và không có khẩu trang được chia theo tỷ lệ: 80% cho huấn luyện. 20% cho kiểm tra.

Các thông số huấn luyện được thiết lập như sau: Thuật toán tối ưu: Adam, với tốc độ học 0.0001. Kích thước batch: 32 mẫu. Số epoch: 18. Phần cứng: GPU NVIDIA RTX 3090.

B. Phương pháp Đánh Giá

Hiệu năng của hệ thống được đánh giá qua các chỉ số:

Độ chính xác (Accuracy) Tỷ lệ nhận diện đúng trên tổng số mẫu.

F1 Score Trung bình điều hòa của Precision và Recall, cho thấy sự cân bằng giữa khả năng phát hiện và độ chính xác.

C. Kết Quả Thử Nghiệm và Phân Tích Chi Tiết

Bảng 1 dưới đây trình bày so sánh giữa các mô hình nhận diện:

TABLE I
So sánh hiệu suất các mô hình

Mô hình	Độ chính xác (%)	Thời gian xử lý (ms)
CNN	85.5	25
MobileNetV2	91.2	15
ResNet50	93.7	30
Vision Transformer	95.3	40
Swin Transformer	96.8	35

Kết quả thực nghiệm cho thấy rằng: Các mô hình học sâu đạt được độ chính xác từ 85.5% đến 96.8% tùy vào kiến trúc. Cụ thể, mô hình Swin Transformer đạt độ chính xác cao nhất (96.8%), trong khi mô hình CNN cơ bản đạt 85.5%.

Thời gian xử lý của các mô hình dao động từ 15ms đến 40ms, cho thấy rằng các mô hình nhẹ như MobileNetV2 có ưu thế về tốc độ, trong khi các mô hình phức tạp hơn như Vision Transformer có thời gian xử lý cao hơn.

Khi áp dụng kỹ thuật Ensemble và xử lý hậu kỳ, hệ thống có thể đạt độ chính xác lên tới 98% trên một số bộ dữ liệu kiểm tra, điều này cho thấy sự cải thiện đáng kể so với việc sử dụng mô hình đơn lẻ.

Bên cạnh đó, các chỉ số F1 Score cũng được tính toán và cho thấy sự cân bằng tốt giữa khả năng nhận diện chính xác và khả năng phát hiện đúng. Điều này chứng tỏ rằng hệ thống có độ ổn định và tổng quát hóa cao trên các tập dữ liệu khác nhau.

[6]

VIII. Kết Luận và Hướng Phát Triển

Bài báo đã trình bày một hệ thống nhận diện khuôn mặt đeo khẩu trang kết hợp các phương pháp truyền thống và học sâu hiện đại. Kết quả thực nghiệm cho thấy rằng việc áp dụng các mô hình học sâu (đặc biệt là các kiến trúc Transformer) có thể đạt được độ chính xác cao, từ 85% đến trên 96%, và khi kết hợp Ensemble cùng xử lý hậu kỳ, hệ thống đạt được tỉ lệ thành công lên tới 98% trong một số trường hợp.

Tuy nhiên, các mô hình này vẫn đòi hỏi tài nguyên tính toán lớn và có thể gặp hạn chế khi triển khai trên các thiết bị di động. Do đó, hướng phát triển trong tương lai bao gồm:

Mở rộng tập dữ liệu Thu thập thêm dữ liệu từ nhiều nguồn và điều kiện ánh sáng, góc chụp khác nhau để tăng tính tổng quát của mô hình.

Cải tiến mô hình Tích hợp cơ chế Attention và các biến thể của CNN nhằm tối ưu hóa quá trình trích xuất đặc trưng khi khuôn mặt bị che.

Tối ưu hóa cho thiết bị di động Rút gọn kiến trúc mô hình nhằm giảm thiểu yêu cầu về tài nguyên và tăng tốc độ xử lý.

Bảo vệ quyền riêng tư Phát triển các giải pháp mã hóa và ẩn danh dữ liệu để bảo vệ thông tin cá nhân của người dùng.

Nghiên cứu về xử lý hậu kỳ Kết hợp thêm các phương pháp xử lý hậu kỳ để cải thiện độ tin cậy của kết quả nhận diện trong môi trường thực tế.

[7]

Tài liệu Tham Khảo

References

- [1] He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep Residual Learning for Image Recognition. Hội nghị CVPR của IEEE.
- [2] Dosovitskiy, A., et al. (2020). An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. ICLR.
- [3] Howard, A. G., et al. (2017). MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications.
- [4] Khắc-Hoai Nam Bui, Hongsuk Yi, và Jiho Cho. A vehicle counts by class framework using distinguished regions tracking at multiple intersections. CVPR Workshops, 2020.
- [5] Christoph Feichtenhofer. X3D: Expanding Architectures for Efficient Video Recognition. CVPR 2020.
- [6] Mengmeng Xu, et al. Boundary-sensitive pre-training for temporal localization in videos. ICCV 2021.