

Cấu trúc tổng thể của CAFace

Hình: Khung tổng quan của hệ thống CAFace (Cluster and Aggregate). Đầu vào là tập N ảnh probe, lần lượt qua ba thành phần chính – Style Input Maker (SIM), Cluster Network (CN) và Aggregation Network (AGN) – để sinh ra một vectơ đặc trưng khuôn mặt hợp nhất cuối cùng ① ② .

CAFace là framework hai giai đoạn nhằm hợp nhất một tập lớn các ảnh probe thành một đặc trưng duy nhất cho nhận dạng khuôn mặt ① . Tổng quan dữ liệu luồng như sau (Hình trên): - **Đầu vào (Probe)**: Tập N ảnh (kích thước $H \times W \times 3$). Mỗi ảnh được đưa qua mạng trích xuất đặc trưng (backbone) để thu được: vectơ đặc trưng nhận dạng $f_{\text{sub}i} \in \mathbb{R}^D$ (ví dụ $D=512$) và ma trận đặc trưng trung gian $F_{\text{sub}i} \in \mathbb{R}^{C \times H' \times W'}$.

- **Style Input Maker (SIM)**: Nhận mỗi ma trận $F_{\text{sub}i}$, SIM tính vectơ thông tin style $s_{\text{sub}i} \in \mathbb{R}^d$ bằng cách lấy trung bình μ và độ lệch chuẩn σ theo từng kênh của $F_{\text{sub}i}$ ③ . Vectơ $s_{\text{sub}i}$ (chiều d) có thể được ghép thêm embedding của chuẩn (norm) của $f_{\text{sub}i}$ ④ . Kết quả là N vectơ style $S = \{s_1, \dots, s_N\}$.

- **Cluster Network (CN)**: CN nhận cặp ($f_{\text{sub}i}, s_{\text{sub}i}$) của N ảnh. Nó dùng M truy vấn toàn cục (các tâm cụm học được) để tính ma trận phân cụm $A \in \mathbb{R}^{M \times N}$ (soft assignment) bằng cách so $s_{\text{sub}i}$ với M tâm cụm rồi Softmax theo cột ① ② . Ma trận A cho biết mức độ mỗi ảnh được gán vào từng cụm. Sau đó CN tạo ra M vectơ đặc trưng cụm $C_{\text{sub}j} \in \mathbb{R}^D$ và M vectơ style cụm $X_{\text{sub}j} \in \mathbb{R}^d$, trong đó mỗi $C_{\text{sub}j}$ là tổng tuyến tính của các $f_{\text{sub}i}$ thuộc cụm j, và mỗi $X_{\text{sub}j}$ là tổng tuyến tính của các $s_{\text{sub}i}$ ① .

- **Aggregation Network (AGN)**: AGN nhận M cặp ($C_{\text{sub}j}, X_{\text{sub}j}$) từ CN và dùng một MLP-Mixer để học mối quan hệ giữa các cụm. AGN sinh ra vectơ trọng số $w \in \mathbb{R}^M$ cho mỗi cụm ④ . Cuối cùng, vectơ đặc trưng hợp nhất $f \in \mathbb{R}^D$ được tính bằng tổng có trọng số: $f = \sum_j w_j \cdot C_{\text{sub}j}$ ④ . Đầu ra của toàn hệ thống là vectơ f dùng cho nhận dạng.

Kết hợp lại, CAFace chia quy trình feature fusion thành ba thành phần riêng biệt. Luồng dữ liệu cụ thể có thể tóm tắt như sau:

- **Đầu vào**: Tập ảnh probe $\{I_1, \dots, I_N\}$. Mỗi ảnh đầu tiên qua mạng nhận dạng khuôn mặt (backbone) để lấy đặc trưng $f_i \in \mathbb{R}^D$ và các đặc trưng trung gian $F_i \in \mathbb{R}^{C \times H' \times W'}$.
- **SIM**: Nhận F_i , tính vectơ style $s_i = W \cdot [\mu(F_i); \sigma(F_i)]$ (cộng thêm embedding của $|f_i|$) ③ . Đầu ra của SIM là tập vectơ style $S = \{s_1, \dots, s_N\}$ (Mỗi $s_i \in \mathbb{R}^d$).
- **CN**: Nhận $(f_i, s_i)_{i=1..N}$, sử dụng M truy vấn Q_j để tính ma trận $A \in \mathbb{R}^{M \times N}$ (công thức attention sửa đổi như trong (2)-(3) ⑤ ⑥). Từ A, CN tạo ra hai ma trận: $C \in \mathbb{R}^{M \times D}$ (tập M vectơ C_j) và $X \in \mathbb{R}^{M \times d}$ (M vectơ X_j), với $C_j = (\sum_i A_{ij} f_i) / \sum_i A_{ij}$ và tương tự cho X_j .
- **AGN**: Nhận C, X (kích thước $M \times (D+d)$), qua MLP-Mixer sinh trọng số $w \in \mathbb{R}^M$ ④ . Tính vectơ $f_{\text{fusion}} = \sum_j w_j C_j$. Đầu ra là vectơ đặc trưng hợp nhất $f_{\text{fusion}} \in \mathbb{R}^D$.

Style Input Maker (SIM)

SIM là mạng phụ trích xuất thông tin kiểu ảnh (style) từ mỗi ảnh mà không liên quan đến nhận dạng. Đầu vào của SIM là ma trận đặc trưng trung gian $F_{\text{sub}} \in \mathbb{R}^{C \times H' \times W'}$ thu được từ backbone. SIM tính toán giá trị trung bình μ và độ lệch chuẩn σ của F_{sub} theo từng kênh, rồi dùng ma trận trọng số học được để biến đổi thành vectơ style $s_{\text{sub}} \in \mathbb{R}^d$ ³. Kết quả bao gồm thông tin về độ sáng, độ tương phản, chất lượng ảnh, độ xoay, v.v., vốn đã bị loại bỏ khỏi vectơ nhận dạng f_{sub} do mạng FR được huấn luyện nhận dạng chỉ quan tâm đến ID. Ngoài ra, SIM còn nhúng chuẩn $|f_{\text{sub}}|$ (theo dạng sinusoidal) vào s_{sub} để đánh giá độ tin cậy của ảnh ³.

- **Đầu vào:** Ma trận $F_{\text{sub}} \in \mathbb{R}^{C \times H' \times W'}$ (kiểu float) kích thước $C \times H' \times W'$.

- **Đầu ra:** Vectơ style $s_{\text{sub}} \in \mathbb{R}^d$ (độ dài d).

Cluster Network (CN)

CN dùng cơ chế attention sửa đổi để gom nhóm N vectơ đầu vào vào M cụm cố định. Đầu vào của CN là bộ N vectơ nhận dạng $f_{\text{sub}} \in \mathbb{R}^D$ và N vectơ style $s_{\text{sub}} \in \mathbb{R}^d$. Bên trong, CN có M truy vấn toàn cục $Q_{\text{sub}} \in \mathbb{R}^{M \times d}$ (học được chung cho mọi probe) ⁷. CN tính ma trận affinities $Q \cdot S^T$ rồi áp dụng Softmax theo cột để thu ma trận phân cụm $A \in \mathbb{R}^{M \times N}$ ⁷ ¹. Từ A , CN tạo ra M vectơ cụm $C_{\text{sub}} \in \mathbb{R}^D$ và M vectơ style cụm $X_{\text{sub}} \in \mathbb{R}^d$ bằng tổng cộng tuyến tính: mỗi $C_{\text{sub}} = (\sum_i A_{i,j} f_{\text{sub}})/\sum_i A_{i,j}$, tương tự với X_{sub} từ s_{sub} ¹. Do đó, số lượng đầu ra của CN luôn cố định M cụm, bất kể N thay đổi ².

- **Đầu vào:** N vectơ $f_{\text{sub}} \in \mathbb{R}^D$ và N vectơ $s_{\text{sub}} \in \mathbb{R}^d$.

Aggregation Network (AGN)

AGN kết hợp M vectơ cụm thành một vectơ duy nhất. Đầu vào của AGN là tập M cặp (C_{sub} , X_{sub}) (mảng kích thước $M \times (D+d)$). AGN sử dụng cấu trúc MLP-Mixer để lan truyền thông tin giữa các cụm. Đầu tiên, nó biểu diễn M vectơ này thành một ma trận rồi xuất ra vectơ $w \in \mathbb{R}^M$ biểu diễn “tầm quan trọng” của mỗi cụm ⁴. Cuối cùng, vectơ đầu ra $f \in \mathbb{R}^D$ được tính bằng tổng có trọng số: $f = \sum_{j=1..M} w_{\text{sub}} \cdot C_{\text{sub}}$. Giá trị w được xác định sao cho $\sum_j w_{\text{sub}} = 1$ (qua Softmax) để f là trung bình có trọng số của các cụm. Kết quả f là vectơ đặc trưng cuối cùng phục vụ tác vụ nhận dạng khuôn mặt.

- **Đầu vào:** M vectơ $C_{\text{sub}} \in \mathbb{R}^D$ và M vectơ $X_{\text{sub}} \in \mathbb{R}^d$.

- **Đầu ra:** Vectơ hợp nhất $f \in \mathbb{R}^D$ và vectơ trọng số cụm $w \in \mathbb{R}^M$ ⁴.

Mỗi thành phần trên đóng vai trò rõ ràng trong pipeline: SIM trích xuất thông tin style (độc lập với ID), CN gom nhóm đầu vào thành M đặc trưng trung gian cố định, và AGN tổng hợp các đặc trưng đó thành một biểu diễn cuối cùng. Cách thiết kế này cho phép CAFace xử lý hiệu quả các tập probe rất lớn và duy trì bất biến với thứ tự đầu vào ⁸ ⁹.

Tài liệu tham khảo: Các thông tin trên được trích từ bài báo gốc CAFace ¹ ² ³ ⁴.

¹ ² ³ ⁴ ⁵ ⁶ ⁷ ⁸ ⁹ [2210.10864] Cluster and Aggregate: Face Recognition with Large Probe Set

<https://arxiv.labs.arxiv.org/html/2210.10864>