

Reconnaissance d'Émotions Faciales en Temps Réel

Rapport Technique - Deep Learning

Cyril

Novembre 2025

Résumé

Ce rapport présente le développement d'un système de reconnaissance d'émotions faciales en temps réel utilisant des réseaux de neurones convolutifs (CNN). Nous détaillons l'architecture initiale, les optimisations apportées, ainsi que les choix techniques effectués pour améliorer les performances de détection. Le système est capable de classifier huit émotions (colère, dégoût, peur, joie, tristesse, surprise, neutralité, et mépris) à partir d'un flux vidéo webcam, en utilisant le dataset Balanced AffectNet.

Table des matières

1 Introduction

La reconnaissance automatique des émotions faciales est un domaine en pleine expansion avec des applications dans l’interaction homme-machine, la santé mentale, le marketing, et l’éducation. Ce projet implémente un système complet de reconnaissance d’émotions en temps réel, depuis l’entraînement du modèle jusqu’à l’inférence via webcam.

1.1 Objectifs

- Développer un modèle CNN capable de classifier 8 émotions faciales (7 de FER2013 + Contempt de FER+)
- Implémenter une application temps réel avec détection de visage
- Optimiser les performances pour une utilisation fluide
- Améliorer la précision par des techniques avancées de deep learning

2 Architecture Initiale

2.1 Dataset : Balanced AffectNet

Le dataset Balanced AffectNet est une version équilibrée et pré-traitée du dataset AffectNet original. Il contient :

- 41 008 images de visages en couleur (RGB)
- Résolution : 75×75 pixels
- 8 classes d’émotions parfaitement équilibrées ($\sim 5\ 126$ images par classe)
- Split prédéfini : train (29 526), val (7 382), test (4 100)

Avantages par rapport à FER2013 :

- Images RGB offrant plus d’informations (couleur de peau, rougeurs)
- Résolution supérieure (75×75 vs 48×48)
- Dataset parfaitement équilibré (pas de biais de classe)
- Classe Contempt (mépris) incluse nativement

2.2 Modèle Initial (model.py)

L’architecture initiale est un CNN simple avec 3 blocs convolutifs :

```
1 class FaceEmotionCNN(nn.Module):  
2     def __init__(self, num_classes=8, in_channels=3, input_size=75):  
3         super(FaceEmotionCNN, self).__init__()  
4         # Bloc Convolutif 1 (75x75 -> 37x37)  
5         self.conv1a = nn.Conv2d(in_channels, 64, kernel_size=3, padding  
6 =1)  
7         self.bn1a = nn.BatchNorm2d(64)  
8         self.conv1b = nn.Conv2d(64, 64, kernel_size=3, padding=1)  
9         self.bn1b = nn.BatchNorm2d(64)  
10        self.pool1 = nn.MaxPool2d(2, 2)  
11        self.dropout1 = nn.Dropout(0.1)  
12  
13        # ... Blocs 2, 3, 4 similaires ...  
14  
15        # Global Average Pooling  
16        self.global_avg_pool = nn.AdaptiveAvgPool2d(1)
```

```

17     # Couches Fully Connected
18     self.fc1 = nn.Linear(512, 256)
19     self.fc2 = nn.Linear(256, 128)
20     self.fc3 = nn.Linear(128, num_classes) # 8 emotions

```

Listing 1 – Architecture CNN pour Balanced AffectNet

2.2.1 Analyse de l'architecture initiale

TABLE 1 – Dimensions à travers le réseau (Balanced AffectNet)

Couche	Entrée	Sortie	Paramètres
Block1 (2×Conv) + Pool	$3 \times 75 \times 75$	$64 \times 37 \times 37$	38,592
Block2 (2×Conv) + Pool	$64 \times 37 \times 37$	$128 \times 18 \times 18$	147,712
Block3 (2×Conv) + Pool	$128 \times 18 \times 18$	$256 \times 9 \times 9$	590,336
Block4 (2×Conv) + Pool	$256 \times 9 \times 9$	$512 \times 4 \times 4$	2,360,320
Global Avg Pool	$512 \times 4 \times 4$	512	0
FC1 + BN	512	256	131,584
FC2 + BN	256	128	33,024
FC3	128	8	1,032
Total			$\approx 3.3M$

Points positifs :

- Utilisation de Batch Normalization pour stabiliser l'entraînement
- Dropout (50%) pour régulariser et éviter l'overfitting
- Architecture simple et rapide à entraîner

Limitations :

- Peu de couches convolutives (faible capacité d'extraction de features)
- Nombre de filtres limité (32-64-128)
- Pas de régularisation dans les couches convolutives
- Couche FC1 très large (2.3M paramètres) créant un goulet d'étranglement

2.3 Script d'entraînement initial (train.py)

```

1 # Hyperparametres
2 BATCH_SIZE = 64
3 LEARNING_RATE = 0.001
4 EPOCHS = 25
5
6 # Transformations simples
7 transform = transforms.Compose([
8     transforms.ToPILImage(),
9     transforms.ToTensor(),
10])
11
12 # Chargement du dataset
13 dataset = FER2013Dataset('./data/fer2013/fer2013.csv', transform=transform)
14 train_loader = DataLoader(dataset, batch_size=BATCH_SIZE, shuffle=True)
15

```

```

16 # Modele et optimiseur
17 model = FaceEmotionCNN().to(device)
18 criterion = nn.CrossEntropyLoss()
19 optimizer = optim.Adam(model.parameters(), lr=LEARNING_RATE)
20
21 # Boucle d'entraînement
22 for epoch in range(EPOCHS):
23     running_loss = 0.0
24     for inputs, labels in train_loader:
25         inputs, labels = inputs.to(device), labels.to(device)
26         optimizer.zero_grad()
27         outputs = model(inputs)
28         loss = criterion(outputs, labels)
29         loss.backward()
30         optimizer.step()
31         running_loss += loss.item()
32
33     print(f"Epoch {epoch+1}, Loss: {running_loss/len(train_loader)}")
34
35 torch.save(model.state_dict(), 'emotion_model.pth')

```

Listing 2 – Script d’entraînement initial

Limitations de l’entraînement initial :

- Pas de data augmentation (risque d’overfitting)
- Pas de séparation train/validation (impossible de détecter l’overfitting)
- Learning rate fixe (convergence sous-optimale)
- Pas d’early stopping (gaspillage de ressources)
- Pas de gestion du déséquilibre des classes

2.4 Application initiale (app.py)

L’application initiale effectue :

1. Capture vidéo via webcam
2. Détection de visage avec Haar Cascade
3. Prétraitement de l’image (redimensionnement 48×48 , grayscale)
4. Inférence avec le modèle CNN
5. Affichage de l’émotion détectée avec emoji

```

1 while True:
2     ret, frame = cap.read()
3     gray_frame = cv2.cvtColor(frame, cv2.COLOR_BGR2GRAY)
4     faces = face_cascade.detectMultiScale(gray_frame,
5                                           scaleFactor=1.3,
6                                           minNeighbors=5)
7
8     for (x, y, w, h) in faces:
9         cv2.rectangle(frame, (x, y), (x + w, y + h), (255, 0, 0), 2)
10        roi_gray = gray_frame[y:y+h, x:x+w]
11
12        # Preprocessing
13        roi_tensor = data_transform(roi_gray).unsqueeze(0).to(device)
14
15        # Prediction
16        with torch.no_grad():

```

```

17     outputs = model(roi_tensor)
18     probabilities = torch.nn.functional.softmax(outputs, dim=1)
19     max_prob, predicted_idx = torch.max(probabilities, 1)
20
21     idx = predicted_idx.item()
22     confidence = max_prob.item() * 100
23     emotion_text = emotion_dict[idx]

```

Listing 3 – Boucle principale de l’application initiale

Limitations :

- Prédictions instables (oscillations entre frames)
- Pas de normalisation de l’éclairage
- Paramètres de détection de visage sous-optimaux
- Interface utilisateur basique

3 Optimisations de l’Architecture

3.1 Nouvelle Architecture CNN

L’architecture améliorée comporte plusieurs optimisations majeures :

```

1 class FaceEmotionCNN(nn.Module):
2     def __init__(self, num_classes=7):
3         super(FaceEmotionCNN, self).__init__()
4
5         # Bloc 1 (48x48 -> 24x24) - Double convolution
6         self.conv1a = nn.Conv2d(1, 64, kernel_size=3, padding=1)
7         self.bn1a = nn.BatchNorm2d(64)
8         self.conv1b = nn.Conv2d(64, 64, kernel_size=3, padding=1)
9         self.bn1b = nn.BatchNorm2d(64)
10        self.pool1 = nn.MaxPool2d(2, 2)
11        self.dropout1 = nn.Dropout(0.25)
12
13        # Bloc 2 (24x24 -> 12x12)
14        self.conv2a = nn.Conv2d(64, 128, kernel_size=3, padding=1)
15        self.bn2a = nn.BatchNorm2d(128)
16        self.conv2b = nn.Conv2d(128, 128, kernel_size=3, padding=1)
17        self.bn2b = nn.BatchNorm2d(128)
18        self.pool2 = nn.MaxPool2d(2, 2)
19        self.dropout2 = nn.Dropout(0.25)
20
21        # Bloc 3 (12x12 -> 6x6)
22        self.conv3a = nn.Conv2d(128, 256, kernel_size=3, padding=1)
23        self.bn3a = nn.BatchNorm2d(256)
24        self.conv3b = nn.Conv2d(256, 256, kernel_size=3, padding=1)
25        self.bn3b = nn.BatchNorm2d(256)
26        self.pool3 = nn.MaxPool2d(2, 2)
27        self.dropout3 = nn.Dropout(0.25)
28
29        # Bloc 4 (6x6 -> 3x3)
30        self.conv4a = nn.Conv2d(256, 512, kernel_size=3, padding=1)
31        self.bn4a = nn.BatchNorm2d(512)
32        self.conv4b = nn.Conv2d(512, 512, kernel_size=3, padding=1)
33        self.bn4b = nn.BatchNorm2d(512)
34        self.pool4 = nn.MaxPool2d(2, 2)
35        self.dropout4 = nn.Dropout(0.25)

```

```

36
37     # Global Average Pooling
38     self.global_avg_pool = nn.AdaptiveAvgPool2d(1)
39
40     # Fully Connected avec dropout progressif
41     self.fc1 = nn.Linear(512, 256)
42     self.bn_fc1 = nn.BatchNorm1d(256)
43     self.dropout_fc1 = nn.Dropout(0.5)
44
45     self.fc2 = nn.Linear(256, 128)
46     self.bn_fc2 = nn.BatchNorm1d(128)
47     self.dropout_fc2 = nn.Dropout(0.4)
48
49     self.fc3 = nn.Linear(128, num_classes)

```

Listing 4 – Architecture CNN améliorée

3.1.1 Justification des choix architecturaux

1. Double convolution par bloc (style VGG)

Inspiré de VGGNet, deux convolutions 3×3 successives permettent un champ réceptif de 5×5 avec moins de paramètres qu'une convolution 5×5 unique :

$$\text{Paramètres } 3 \times 3 \times 2 = 2 \times (3^2 \times C^2) = 18C^2 \quad (1)$$

$$\text{Paramètres } 5 \times 5 = 5^2 \times C^2 = 25C^2 \quad (2)$$

2. Augmentation progressive des filtres (64→128→256→512)

Suit le principe que les premières couches extraient des features simples (bords, textures) tandis que les couches profondes capturent des concepts complexes (parties du visage, expressions).

3. Dropout progressif (0.25 dans conv, 0.5→0.4 dans FC)

Le dropout dans les couches convolutives (0.25) régularise sans trop perturber l'apprentissage spatial. Un dropout plus fort dans les FC (0.5) combat l'overfitting où le risque est le plus élevé.

4. Global Average Pooling (GAP)

Remplace le flatten traditionnel. Avantages :

- Réduit drastiquement les paramètres ($512 \times 3 \times 3 = 4608 \rightarrow 512$)
- Agit comme régularisateur
- Invariance spatiale accrue

5. Initialisation Kaiming

Initialisation adaptée aux fonctions d'activation ReLU :

$$W \sim \mathcal{N} \left(0, \sqrt{\frac{2}{n_{in}}} \right) \quad (3)$$

TABLE 2 – Comparaison des architectures

Caractéristique	Initial	Amélioré
Blocs convolutifs	3	4
Convolutions par bloc	1	2
Filtres max	128	512
Paramètres totaux	$\approx 2.46M$	$\approx 4.8M$
Global Average Pooling	Non	Oui
Dropout conv	Non	Oui (0.25)
Batch Norm FC	Non	Oui

4 Optimisations de l’Entraînement

4.1 Data Augmentation

La data augmentation est cruciale pour un petit dataset comme FER2013. Elle augmente artificiellement la diversité des données d’entraînement.

```

1 train_transform = transforms.Compose([
2     transforms.ToPILImage(),
3     transforms.RandomHorizontalFlip(p=0.5),
4     transforms.RandomRotation(10),
5     transforms.RandomAffine(
6         degrees=0,
7         translate=(0.1, 0.1),
8         scale=(0.9, 1.1)
9     ),
10    transforms.ColorJitter(brightness=0.2, contrast=0.2),
11    transforms.ToTensor(),
12 ])

```

Listing 5 – Transformations de data augmentation

TABLE 3 – Transformations de data augmentation

Transformation	Justification
RandomHorizontalFlip ($p=0.5$)	Les expressions sont symétriques. Double effectivement le dataset.
RandomRotation ($\pm 10^\circ$)	Simule les légères inclinaisons de tête naturelles.
RandomAffine (translate)	Compense les variations de position du visage dans le cadre.
RandomAffine (scale 0.9-1.1)	Simule différentes distances caméra-visage.
ColorJitter (brightness, contrast)	Robustesse aux variations d’éclairage.

4.2 Séparation Train/Validation

```

1 VALIDATION_SPLIT = 0.15 # 15% pour validation
2

```

```

3 val_size = int(len(full_dataset) * VALIDATION_SPLIT)
4 train_size = len(full_dataset) - val_size
5
6 train_dataset, val_dataset = random_split(
7     full_dataset,
8     [train_size, val_size],
9     generator=torch.Generator().manual_seed(42)
10 )

```

Listing 6 – Séparation du dataset

Importance : La validation permet de :

- Déetecter l’overfitting (train loss ↓ mais val loss ↑)
- Sélectionner le meilleur modèle
- Ajuster les hyperparamètres

4.3 Early Stopping

L’early stopping arrête l’entraînement quand la validation ne s’améliore plus :

```

1 PATIENCE = 7
2 best_val_acc = 0.0
3 patience_counter = 0
4
5 for epoch in range(EPOCHS):
6     # ... entraînement ...
7     val_loss, val_acc = validate(model, val_loader, criterion, device)
8
9     if val_acc > best_val_acc:
10         best_val_acc = val_acc
11         patience_counter = 0
12         torch.save(model.state_dict(), 'emotion_model_best.pth')
13     else:
14         patience_counter += 1
15         if patience_counter >= PATIENCE:
16             print("Early stopping triggered!")
17             break

```

Listing 7 – Implémentation de l’early stopping

4.4 Learning Rate Scheduler

Le scheduler réduit le learning rate quand la validation stagne :

```

1 scheduler = optim.lr_scheduler.ReduceLROnPlateau(
2     optimizer,
3     mode='min',      # Surveille la loss
4     factor=0.5,      # Divise LR par 2
5     patience=3       # Attends 3 epochs sans amélioration
6 )
7
8 # Dans la boucle d'entraînement
9 scheduler.step(val_loss)

```

Listing 8 – Learning rate scheduler

Principe : Un LR élevé au début permet une convergence rapide, puis un LR plus faible permet un affinage précis.

4.5 Optimiseur AdamW

```
1 optimizer = optim.AdamW(
2     model.parameters(),
3     lr=LEARNING_RATE,
4     weight_decay=1e-4
5 )
```

Listing 9 – Optimiseur AdamW avec weight decay

AdamW corrige un problème de Adam : le weight decay est appliqué directement aux poids plutôt qu'au gradient, ce qui donne une meilleure régularisation L2.

4.6 Gradient Clipping

```
1 loss.backward()
2 torch.nn.utils.clip_grad_norm_(model.parameters(), max_norm=1.0)
3 optimizer.step()
```

Listing 10 – Gradient clipping pour la stabilité

Empêche les gradients explosifs qui peuvent déstabiliser l'entraînement.

5 Optimisations de l'Application Temps Réel

5.1 Lissage Temporel des Prédictions

Le problème principal des prédictions frame-par-frame est l'instabilité : l'émotion affichée peut changer rapidement entre frames successives, créant un effet de *flickering* désagréable.

```
1 from collections import deque
2
3 SMOOTHING_WINDOW = 5
4 prediction_history = deque(maxlen=SMOOTHING_WINDOW)
5
6 def get_smoothed_prediction(current_probs):
7     prediction_history.append(current_probs.cpu().numpy())
8
9     if len(prediction_history) < 2:
10         return current_probs
11
12     # Moyenne pondérée (frames récentes = plus de poids)
13     weights = np.linspace(0.5, 1.0, len(prediction_history))
14     weights = weights / weights.sum()
15
16     smoothed = np.zeros(7)
17     for i, probs in enumerate(prediction_history):
18         smoothed += weights[i] * probs
19
20     return torch.tensor(smoothed)
```

Listing 11 – Lissage temporel avec moyenne pondérée

Principe : Au lieu d'utiliser uniquement la prédiction de la frame courante, on calcule une moyenne pondérée sur les N dernières frames. Les frames récentes ont plus de poids pour maintenir la réactivité.

5.2 Égalisation d’Histogramme

```
1 # Avant preprocessing  
2 roi_gray = cv2.equalizeHist(roi_gray)
```

Listing 12 – Normalisation de l’éclairage

L’égalisation d’histogramme normalise la distribution des niveaux de gris, rendant le modèle plus robuste aux variations d’éclairage.

Image sombre → Histogramme égalisé

FIGURE 1 – L’égalisation redistribue les intensités sur toute la plage [0, 255]

5.3 Paramètres de Détection Optimisés

```
1 faces = face_cascade.detectMultiScale(  
2     gray_frame,  
3     scaleFactor=1.1,      # Plus précis (etait 1.3)  
4     minNeighbors=5,  
5     minSize=(48, 48),    # Taille minimum  
6     flags=cv2.CASCADE_SCALE_IMAGE  
7 )
```

Listing 13 – Paramètres de détection de visage optimisés

TABLE 4 – Impact des paramètres de détection

Paramètre	Valeur basse	Valeur haute
scaleFactor	Plus précis, plus lent	Moins précis, plus rapide
minNeighbors	Plus de faux positifs	Moins de detections
minSize	Déetecte petits visages	Ignore petits visages

5.4 Interface Utilisateur Améliorée

- **Couleurs par émotion** : Chaque émotion a une couleur distinctive
- **Barres de progression** : Affichage des probabilités de toutes les classes
- **Fond pour le texte** : Améliore la lisibilité

```
1 emotion_colors = {  
2     0: (0, 0, 255),      # Rouge - Angry  
3     1: (0, 128, 0),      # Vert foncé - Disgust  
4     2: (128, 0, 128),    # Violet - Fear  
5     3: (0, 255, 255),    # Jaune - Happy  
6     4: (255, 0, 0),      # Bleu - Sad  
7     5: (0, 165, 255),    # Orange - Surprise  
8     6: (128, 128, 128)   # Gris - Neutral  
9 }
```

Listing 14 – Couleurs par émotion

6 Amélioration du Dataset

6.1 Avantages du Balanced AffectNet

Par rapport à FER2013, le dataset Balanced AffectNet présente plusieurs avantages majeurs :

1. **Résolution supérieure** : Images 75×75 pixels (vs 48×48)
2. **Images RGB** : 3 canaux de couleur offrant plus d'informations

Émotion	Nombre d'images	Pourcentage
Angry	5,126	12.5%
Contempt	5,126	12.5%
Disgust	5,126	12.5%
Fear	5,126	12.5%
Happy	5,126	12.5%
Neutral	5,126	12.5%
Sad	5,126	12.5%
Surprise	5,126	12.5%

TABLE 5 – Distribution parfaitement équilibrée des classes dans Balanced AffectNet.

3. **Dataset équilibré** : Pas de biais de classe, contrairement à FER2013
4. **8 classes natives** : Inclut Contempt (mépris) dès le départ
5. **Meilleure qualité** : Annotations plus fiables

6.2 Chargement du Dataset

Le dataset Balanced AffectNet est organisé en dossiers par émotion :

```
1 class BalancedAffectNetDataset(Dataset):
2     EMOTION_CLASSES = {
3         'Anger': 0, 'Disgust': 1, 'Fear': 2, 'Happy': 3,
4         'Sad': 4, 'Surprise': 5, 'Neutral': 6, 'Contempt': 7,
5     }
6
7     def __init__(self, root_dir, split='train', transform=None):
8         self.images = []
9         self.labels = []
10
11        for emotion_name, emotion_idx in self.EMOTION_CLASSES.items():
12            emotion_dir = os.path.join(root_dir, split, emotion_name)
13            for img_name in os.listdir(emotion_dir):
14                self.images.append(os.path.join(emotion_dir, img_name))
15                self.labels.append(emotion_idx)
16
17    def __getitem__(self, idx):
18        image = Image.open(self.images[idx]).convert('RGB')
19        image = np.array(image) # 75x75x3
20        label = self.labels[idx]
21
22        if self.transform:
```

```

23     image = self.transform(image=image) ['image']
24
25     return image, label

```

Listing 15 – Utilisation du dataset AffectNet

6.3 Équilibrage des Classes

Pour compenser le déséquilibre (notamment le manque de *Disgust*), deux techniques sont implémentées :

6.3.1 Poids de classe dans la loss

```

1 def get_class_weights(dataset):
2     class_counts = np.bincount(labels, minlength=7)
3     weights = 1.0 / class_counts
4     weights = weights / weights.sum() * len(weights)
5     return torch.FloatTensor(weights)
6
7 class_weights = get_class_weights(train_dataset)
8 criterion = nn.CrossEntropyLoss(weight=class_weights)

```

Listing 16 – CrossEntropyLoss avec poids de classe

Le poids de chaque classe est inversement proportionnel à sa fréquence :

$$w_c = \frac{N}{N_c \times C} \quad (4)$$

où N est le nombre total d'échantillons, N_c le nombre d'échantillons de la classe c , et C le nombre de classes.

6.3.2 WeightedRandomSampler

```

1 def get_balanced_sampler(dataset):
2     class_counts = np.bincount(labels, minlength=7)
3     weights = 1.0 / class_counts
4     sample_weights = weights[labels]
5
6     sampler = WeightedRandomSampler(
7         weights=sample_weights,
8         num_samples=len(sample_weights),
9         replacement=True
10    )
11    return sampler

```

Listing 17 – Sampler équilibré

Le sampler sur-échantillonne les classes minoritaires pendant l'entraînement.

7 Techniques d'Entraînement Avancées

Pour améliorer significativement les performances, notamment sur les classes difficiles (Disgust, Contempt, Fear, Angry), nous avons développé un script d'entraînement avancé (`train_advanced.py`) intégrant les techniques les plus récentes du deep learning.

7.1 Focal Loss pour le Déséquilibre des Classes

La Cross-Entropy standard traite tous les exemples de manière égale. La **Focal Loss** [?] ajoute un facteur de modulation qui réduit la contribution des exemples faciles et se concentre sur les exemples difficiles :

$$\text{FL}(p_t) = -\alpha_t(1 - p_t)^\gamma \log(p_t) \quad (5)$$

où p_t est la probabilité prédictive pour la vraie classe, α_t est le poids de la classe, et γ (gamma) contrôle le focus sur les exemples difficiles.

```

1 class FocalLoss(nn.Module):
2     def __init__(self, alpha=None, gamma=2.0, reduction='mean'):
3         super().__init__()
4         self.alpha = alpha # Poids par classe
5         self.gamma = gamma # Focus parameter (2.0 recommandé)
6         self.reduction = reduction
7
8     def forward(self, inputs, targets):
9         ce_loss = F.cross_entropy(inputs, targets,
10                             weight=self.alpha,
11                             reduction='none')
12         pt = torch.exp(-ce_loss) # p_t
13         focal_loss = ((1 - pt) ** self.gamma) * ce_loss
14
15         if self.reduction == 'mean':
16             return focal_loss.mean()
17         return focal_loss

```

Listing 18 – Implémentation de la Focal Loss

Avantage : Avec $\gamma = 2$, un exemple bien classifié ($p_t = 0.9$) contribue $100\times$ moins qu'un exemple difficile ($p_t = 0.1$), permettant au modèle de se concentrer sur les émotions sous-représentées comme Disgust.

7.2 Mixup : Régularisation par Interpolation

Mixup [?] crée de nouveaux exemples d'entraînement en interpolant linéairement des paires d'images et leurs labels :

$$\tilde{x} = \lambda x_i + (1 - \lambda)x_j \quad (6)$$

$$\tilde{y} = \lambda y_i + (1 - \lambda)y_j \quad (7)$$

où $\lambda \sim \text{Beta}(\alpha, \alpha)$ avec $\alpha = 0.2$.

```

1 def mixup_data(x, y, alpha=0.2):
2     if alpha > 0:
3         lam = np.random.beta(alpha, alpha)
4     else:
5         lam = 1
6
7     batch_size = x.size(0)
8     index = torch.randperm(batch_size).to(x.device)
9
10    mixed_x = lam * x + (1 - lam) * x[index, :]
11    y_a, y_b = y, y[index]
12    return mixed_x, y_a, y_b, lam

```

```

13
14 def mixup_criterion(criterion, pred, y_a, y_b, lam):
15     return lam * criterion(pred, y_a) + (1 - lam) * criterion(pred, y_b)

```

Listing 19 – Implémentation de Mixup

7.3 CutMix : Augmentation par Découpage

CutMix [?] est une variante de Mixup qui découpe et colle des régions rectangulaires entre images :

```

1 def cutmix_data(x, y, alpha=1.0):
2     lam = np.random.beta(alpha, alpha)
3     batch_size = x.size(0)
4     index = torch.randperm(batch_size).to(x.device)
5
6     # Calcul de la boîte de découpe
7     W, H = x.size(2), x.size(3)
8     cut_rat = np.sqrt(1. - lam)
9     cut_w = int(W * cut_rat)
10    cut_h = int(H * cut_rat)
11
12    cx = np.random.randint(W)
13    cy = np.random.randint(H)
14
15    bbx1 = np.clip(cx - cut_w // 2, 0, W)
16    bby1 = np.clip(cy - cut_h // 2, 0, H)
17    bbx2 = np.clip(cx + cut_w // 2, 0, W)
18    bby2 = np.clip(cy + cut_h // 2, 0, H)
19
20    x[:, :, bbx1:bbx2, bby1:bby2] = x[index, :, bbx1:bbx2, bby1:bby2]
21
22    # Ajuster lambda selon la surface découpée
23    lam = 1 - ((bbx2-bbx1)*(bby2-bby1) / (W*H))
24    return x, y, y[index], lam

```

Listing 20 – Implémentation de CutMix

Avantage : CutMix force le modèle à utiliser des indices partiels du visage (un œil, la bouche seule) plutôt que de dépendre de l'image entière, améliorant la robustesse.

7.4 Label Smoothing

Le label smoothing régularise en remplaçant les labels hard (one-hot) par des labels soft :

$$y_{\text{smooth}} = (1 - \epsilon) \cdot y_{\text{hard}} + \frac{\epsilon}{K} \quad (8)$$

avec $\epsilon = 0.1$ et $K = 8$ classes.

```

1 class LabelSmoothingLoss(nn.Module):
2     def __init__(self, classes, smoothing=0.1):
3         super().__init__()
4         self.smoothing = smoothing
5         self.cls = classes
6
7     def forward(self, pred, target):

```

```

8     confidence = 1.0 - self.smoothing
9     smooth_val = self.smoothing / (self.cls - 1)
10
11    one_hot = torch.zeros_like(pred)
12    one_hot.fill_(smooth_val)
13    one_hot.scatter_(1, target.unsqueeze(1), confidence)
14
15    log_prob = F.log_softmax(pred, dim=1)
16    return -(one_hot * log_prob).sum(dim=1).mean()

```

Listing 21 – Application du Label Smoothing

7.5 Augmentation Avancée avec Albumentations

Nous utilisons la bibliothèque Albumentations pour des transformations plus sophistiquées :

```

1 import albumentations as A
2
3 train_transform = A.Compose([
4     A.HorizontalFlip(p=0.5),
5     A.ShiftScaleRotate(
6         shift_limit=0.1,
7         scale_limit=0.15,
8         rotate_limit=15,
9         p=0.5
10    ),
11    A.OneOf([
12        A.MotionBlur(blur_limit=3, p=1.0),
13        A.GaussianBlur(blur_limit=3, p=1.0),
14        A.GaussNoise(var_limit=(10, 50), p=1.0),
15    ], p=0.3),
16    A.OneOf([
17        A.RandomBrightnessContrast(
18            brightness_limit=0.2,
19            contrast_limit=0.2,
20            p=1.0
21        ),
22        A.CLAHE(clip_limit=2.0, p=1.0),
23    ], p=0.5),
24    A.CoarseDropout(
25        max_holes=1,
26        max_height=12,
27        max_width=12,
28        fill_value=0,
29        p=0.25
30    ),
31 ])

```

Listing 22 – Transformations Albumentations

TABLE 6 – Nouvelles transformations Albumentations

Transformation	Justification
MotionBlur / GaussianBlur	Simule le flou de mouvement en conditions réelles
GaussNoise	Robustesse au bruit de capteur
CLAHE	Améliore le contraste local, meilleur que l'égalisation standard
CoarseDropout	Similaire à Cutout, force la redondance des features

7.6 Gradient Accumulation

Pour simuler des batch sizes plus grands sur GPU avec mémoire limitée :

```

1 ACCUMULATION_STEPS = 4 # Effective batch = 32 * 4 = 128
2
3 for i, (inputs, targets) in enumerate(train_loader):
4     outputs = model(inputs)
5     loss = criterion(outputs, targets)
6     loss = loss / ACCUMULATION_STEPS
7     loss.backward()
8
9     if (i + 1) % ACCUMULATION_STEPS == 0:
10         torch.nn.utils.clip_grad_norm_(
11             model.parameters(), max_norm=1.0
12         )
13         optimizer.step()
14         optimizer.zero_grad()
```

Listing 23 – Gradient Accumulation

Avantage : Un batch effectif de 128 images stabilise les gradients et améliore la convergence.

7.7 Cosine Annealing avec Warm Restarts

Scheduler qui suit une courbe cosinus avec redémarrages périodiques :

```

1 scheduler = optim.lr_scheduler.CosineAnnealingWarmRestarts(
2     optimizer,
3     T_0=10,      # Première période de 10 epochs
4     T_mult=2,    # Périodes suivantes x2 (10, 20, 40...)
5     eta_min=1e-6
6 )
```

Listing 24 – CosineAnnealingWarmRestarts

Learning Rate suit une courbe cosinus :

$$\eta_t = \eta_{\min} + \frac{1}{2}(\eta_{\max} - \eta_{\min})(1 + \cos(\frac{T_{\text{cur}}}{T_i}\pi))$$

FIGURE 2 – Le learning rate diminue selon une courbe cosinus puis remonte à chaque restart

8 Application Étendue avec MediaPipe

Pour améliorer la détection des émotions difficiles et enrichir le panel d'emojis affichables, nous avons intégré **MediaPipe**, la bibliothèque de Google pour l'analyse faciale et gestuelle en temps réel.

8.1 MediaPipe Face Mesh : Analyse des Landmarks Faciaux

Face Mesh détecte 468 points de repère 3D sur le visage, permettant une analyse fine des caractéristiques faciales :

```
1 import mediapipe as mp
2
3 mp_face_mesh = mp.solutions.face_mesh
4 face_mesh = mp_face_mesh.FaceMesh(
5     max_num_faces=1,
6     refine_landmarks=True, # Points supplémentaires iris
7     min_detection_confidence=0.5,
8     min_tracking_confidence=0.5
9 )
```

Listing 25 – Initialisation de Face Mesh

8.1.1 Features Extraites

```
1 class FacialAnalyzer:
2     # Indices des landmarks clés
3     LEFT_EYE = [33, 160, 158, 133, 153, 144]
4     RIGHT_EYE = [362, 385, 387, 263, 373, 380]
5     MOUTH = [61, 291, 0, 17, 78, 308]
6     LEFT_EYEBROW = [70, 63, 105, 66, 107]
7     RIGHT_EYEBROW = [336, 296, 334, 293, 300]
8
9     def analyze(self, landmarks):
10         features = {}
11
12         # Ouverture des yeux (Eye Aspect Ratio)
13         features['left_eye_open'] = self._eye_aspect_ratio(
14             [landmarks[i] for i in self.LEFT_EYE]
15         )
16         features['right_eye_open'] = self._eye_aspect_ratio(
17             [landmarks[i] for i in self.RIGHT_EYE]
18         )
19
20         # Ouverture de la bouche
21         features['mouth_open'] = self._mouth_aspect_ratio(landmarks)
22
23         # Position des sourcils
24         features['brow_raise'] = self._brow_raise(landmarks)
25         features['brow_squeeze'] = self._brow_squeeze(landmarks)
26
27         # Sourire (ratio largeur/hauteur bouche)
28         features['smile'] = self._smile_ratio(landmarks)
29
30         return features
31
```

```

32     def _eye_aspect_ratio(self, eye_points):
33         # Ratio vertical/horizontal pour détecter clignement
34         vertical = np.linalg.norm(
35             np.array(eye_points[1]) - np.array(eye_points[5])
36         )
37         horizontal = np.linalg.norm(
38             np.array(eye_points[0]) - np.array(eye_points[3])
39         )
40         return vertical / (horizontal + 1e-6)

```

Listing 26 – Classe FacialAnalyzer pour l'extraction de features

8.1.2 Boost des Émotions Difficiles

Les features faciales sont utilisées pour renforcer la détection des émotions que le CNN a du mal à identifier :

```

1 EMOTION_BOOST = {
2     'Angry': 1.4,      # Boost de 40%
3     'Disgust': 1.5,    # Boost de 50%
4     'Fear': 1.3,       # Boost de 30%
5     'Contempt': 1.4,   # Boost de 40%
6     'Neutral': 0.85    # Reduction de 15%
7 }
8
9 def boost_emotions(probs, features):
10    boosted = probs.copy()
11
12    # Boost Angry si sourcils froncés
13    if features.get('brow_squeeze', 0) > 0.6:
14        boosted[ANGRY_IDX] *= 1.3
15
16    # Boost Disgust si nez froncé
17    if features.get('nose_wrinkle', 0) > 0.5:
18        boosted[DISGUST_IDX] *= 1.4
19
20    # Boost Fear si yeux grands ouverts
21    avg_eye = (features['left_eye_open'] +
22                features['right_eye_open']) / 2
23    if avg_eye > 0.4:
24        boosted[FEAR_IDX] *= 1.2
25
26    return boosted / boosted.sum()    # Re-normaliser

```

Listing 27 – Système de boost des émotions

8.2 Extension du Mapping Emoji

Au-delà des 8 émotions de base, les features faciales permettent d'afficher des emojis plus nuancés :

```

1 EXTENDED_EMOJI_MAP = {
2     # Emotions de base
3     'Happy': '😊',
4     'Sad': '😢',
5     'Angry': '😠',
6     'Surprise': '😲',

```

```

7   'Fear': '😱',
8   'Disgust': '🤮',
9   'Neutral': '😐',
10  'Contempt': '😎',
11
12  # Combinaisons basees sur features
13  'very_happy': '😊',           # Sourire large
14  'laughing': '😂',           # Bouche ouverte + sourire
15  'thinking': '🤔',           # Un sourcil leve
16  'sleepy': '😴',             # Yeux fermes
17  'wink': '😉',               # Un oeil ferme
18  'kiss': '😘',                # Bouche en O
19  'love': '😍',                # Coeurs si smile fort
20  'cool': '😎',               # Neutre + confiant
21  'skeptical': '🤨',          # Sourcil asymetrique
22 }
23
24 def get_extended_emoji(emotion, features):
25     # Vérifier conditions spéciales
26     if features['left_eye_open'] < 0.15:
27         if features['right_eye_open'] > 0.25:
28             return EXTENDED_EMOJI_MAP['wink']
29
30     if features['smile'] > 0.7 and features['mouth_open'] > 0.4:
31         return EXTENDED_EMOJI_MAP['laughing']
32
33     if emotion == 'Happy' and features['smile'] > 0.8:
34         return EXTENDED_EMOJI_MAP['very_happy']
35
36     return EXTENDED_EMOJI_MAP.get(emotion, ' ')

```

Listing 28 – Mapping emoji étendu

8.3 MediaPipe Hands : Détection des Gestes

Pour enrichir l’interaction, nous avons ajouté la détection des gestes de la main :

```

1 mp_hands = mp.solutions.hands
2 hands = mp_hands.Hands(
3     static_image_mode=False,
4     max_num_hands=2,
5     min_detection_confidence=0.7,
6     min_tracking_confidence=0.5
7 )

```

Listing 29 – Initialisation de MediaPipe Hands

8.3.1 Gestes Reconnus

```

1 class HandRecognizer:
2     def recognize_gesture(self, landmarks):
3         # Extraire les etats des doigts (leve/baisse)
4         fingers = self._get_finger_states(landmarks)
5         thumb, index, middle, ring, pinky = fingers
6
7         # Pouce leve
8         if thumb and not any([index, middle, ring, pinky]):

```

```

9         return 'thumbs_up', ' '
10
11     # Signe de paix (V)
12     if index and middle and not ring and not pinky:
13         return 'peace', ' '
14
15     # Signe OK
16     if self._is_ok_gesture(landmarks):
17         return 'ok', ' '
18
19     # Rock (cornes)
20     if index and pinky and not middle and not ring:
21         return 'rock', ' '
22
23     # Shaka (pouce + auriculaire)
24     if thumb and pinky and not index and not middle:
25         return 'shaka', ' '
26
27     # Poing ferme
28     if not any(fingers):
29         return 'fist', ' '
30
31     # Main ouverte (tous les doigts)
32     if all(fingers):
33         return 'wave', ' '
34
35     # Pointer
36     if index and not middle and not ring and not pinky:
37         return 'point', ' '
38
39     return None, None
40
41 def _get_finger_states(self, landmarks):
42     # Comparer position des bouts vs articulations
43     thumb = landmarks[4].y < landmarks[3].y
44     index = landmarks[8].y < landmarks[6].y
45     middle = landmarks[12].y < landmarks[10].y
46     ring = landmarks[16].y < landmarks[14].y
47     pinky = landmarks[20].y < landmarks[18].y
48     return [thumb, index, middle, ring, pinky]

```

Listing 30 – Classe HandRecognizer

TABLE 7 – Gestes reconnus et emojis correspondants

Geste	Description	Emoji
Thumbs Up	Pouce levé seul	
Peace	Index et majeur levés	
OK	Pouce et index formant un cercle	
Rock	Index et auriculaire levés	
Shaka	Pouce et auriculaire	
Wave	Tous les doigts levés	
Fist	Poing fermé	
Point	Index seul levé	

9 Optimisations de l'Inférence (app_v3.py)

La version 3 de l'application intègre plusieurs optimisations pour améliorer la précision et la stabilité des prédictions.

9.1 Test-Time Augmentation (TTA)

Le TTA applique plusieurs transformations à l'image d'entrée et moyenne les prédictions :

```
1 class EmotionClassifier:
2     def predict_with_tta(self, face_roi, num_augmentations=3):
3         predictions = []
4
5         # Prediction originale
6         predictions.append(self._predict_single(face_roi))
7
8         # Flip horizontal
9         flipped = cv2.flip(face_roi, 1)
10        predictions.append(self._predict_single(flipped))
11
12        # Légères variations de luminosité
13        for _ in range(num_augmentations - 2):
14            factor = np.random.uniform(0.9, 1.1)
15            adjusted = np.clip(face_roi * factor, 0, 255)
16            predictions.append(
17                self._predict_single(adjusted.astype(np.uint8)))
18
19        # Moyenne des predictions
20        return np.mean(predictions, axis=0)
21
```

Listing 31 – Implémentation du TTA

Avantage : Le TTA réduit la variance des prédictions et améliore la robustesse aux petites perturbations.

9.2 Prétraitement CLAHE

Contrairement à l'égalisation d'histogramme globale, CLAHE (Contrast Limited Adaptive Histogram Equalization) travaille sur des régions locales :

```
1 def preprocess_face(self, face_roi):
2     # Redimensionner
3     face = cv2.resize(face_roi, (48, 48))
4
5     # CLAHE pour le contraste local
6     clahe = cv2.createCLAHE(clipLimit=2.0, tileSize=(4, 4))
7     face = clahe.apply(face)
8
9     # Normalisation
10    face = face.astype(np.float32) / 255.0
11
12    return face
```

Listing 32 – Application de CLAHE

TABLE 8 – Comparaison des méthodes d'égalisation

Méthode	Avantages	Inconvénients
Histogramme global	Simple, rapide	Peut sur-amplifier le bruit
CLAHE	Préserve les détails locaux, contrôle le contraste	Légèrement plus lent

9.3 Lissage Temporel Amélioré

La version 3 utilise un lissage exponentiel plus sophistiqué :

```

1 SMOOTHING_WINDOW = 7
2
3 def get_smoothed_prediction(self, current_probs):
4     self.prediction_history.append(current_probs)
5
6     if len(self.prediction_history) < 2:
7         return current_probs
8
9     # Poids exponentiels (frames récentes plus importantes)
10    n = len(self.prediction_history)
11    weights = np.exp(np.linspace(-1, 0, n))
12    weights = weights / weights.sum()
13
14    smoothed = np.zeros(8)
15    for i, probs in enumerate(self.prediction_history):
16        smoothed += weights[i] * probs
17
18    return smoothed

```

Listing 33 – Lissage temporel avec pondération exponentielle

9.4 Architecture Modulaire

L'application V3 sépare clairement les responsabilités :

```

1 class EmotionClassifier:
2     """Gestion du modèle CNN et predictions"""
3     pass
4
5 class FacialAnalyzer:
6     """Analyse des landmarks faciaux avec MediaPipe"""
7     pass
8
9 class HandRecognizer:
10    """Détection des gestes de la main"""
11    pass
12
13 class EmotionApp:
14    """Application principale orchestrant les composants"""
15    def __init__(self):
16        self.classifier = EmotionClassifier()
17        self.facial_analyzer = FacialAnalyzer()
18        self.hand_recognizer = HandRecognizer()
19
20    def process_frame(self, frame):

```

```

21     # 1. Detecter visage
22     # 2. Analyser features faciales
23     # 3. Predire emotion avec CNN + TTA
24     # 4. Boost avec features
25     # 5. Detecter gestes
26     # 6. Afficher resultats
27     pass

```

Listing 34 – Architecture modulaire de app_v3.py

9.5 Datasets Alternatifs

Pour des performances encore meilleures, d'autres datasets sont supportés :

TABLE 9 – Comparaison des datasets d'émotions faciales

Dataset	Images	Résolution	Classes	Accès
FER2013	35,887	48×48 (Gray)	7	Gratuit (Kaggle)
FER+	35,887	48×48 (Gray)	8	Gratuit (GitHub)
Balanced AffectNet	41,008	75×75 (RGB)	8	Gratuit (Kaggle)
AffectNet (original)	450,000	Variable	8	Demande requise
RAF-DB	30,000	100×100	7	Demande requise

```

1 class AffectNetDataset(Dataset):
2     # Mapping AffectNet vers les 7 classes FER
3     AFFECTNET_TO_FER = {
4         0: 6,    # Neutral
5         1: 3,    # Happy
6         2: 4,    # Sad
7         3: 5,    # Surprise
8         4: 2,    # Fear
9         5: 1,    # Disgust
10        6: 0,   # Anger
11        7: 6,   # Contempt -> Neutral
12    }
13
14    def __getitem__(self, idx):
15        image = Image.open(self.images[idx]).convert('L')
16        image = image.resize((48, 48))
17        label = self.AFFECTNET_TO_FER[self.labels[idx]]
18        return image, label

```

Listing 35 – Support multi-datasets

10 Résultats Attendus

10.1 Amélioration de la Précision

TABLE 10 – Amélioration attendue de la précision avec Balanced AffectNet

Configuration	Précision estimée	Gain
FER2013 baseline (grayscale 48×48)	60-65%	-
Balanced AffectNet (RGB 75×75)	70-75%	+10%
+ Data augmentation RGB	73-78%	+3%
+ Mixup/CutMix	76-80%	+3%
+ Label Smoothing	78-82%	+2%
+ TTA à l’inférence	80-84%	+2%
Total avec optimisations	80-84%	+20-24%

10.2 Amélioration des Classes Difficiles

TABLE 11 – Impact attendu sur les classes (avec dataset équilibré)

Émotion	FER2013 Baseline	Balanced AffectNet	Gain
Angry	55%	75%	+20%
Disgust	35%	72%	+37%
Fear	45%	70%	+25%
Contempt	40%	68%	+28%

Note : L’amélioration majeure sur Disgust et Contempt est due au dataset parfaitement équilibré et aux images RGB de meilleure qualité.

10.3 Amélioration de l’Expérience Utilisateur

- **Stabilité** : Le lissage temporel amélioré élimine le *flickering*
- **Robustesse** : CLAHE gère mieux les variations d’éclairage que l’égalisation standard
- **Précision TTA** : Le test-time augmentation réduit la variance des prédictions
- **Emojis étendus** : Plus de 15 emojis possibles grâce à l’analyse des features faciales
- **Gestes de la main** : 8 gestes reconnus avec emojis correspondants
- **Lisibilité** : Les couleurs et barres de progression améliorent la compréhension
- **Réactivité** : Les paramètres de détection optimisés améliorent la fluidité

11 Conclusion

Ce projet a permis de développer un système complet et avancé de reconnaissance d’émotions faciales en temps réel. Les principales contributions sont :

1. **Architecture CNN optimisée** avec double convolution, dropout progressif, et Global Average Pooling

2. **Pipeline d'entraînement avancé** incluant :
 - Focal Loss pour le déséquilibre des classes
 - Mixup et CutMix pour la régularisation
 - Label Smoothing pour la généralisation
 - Augmentations avancées avec Albumentations
3. **Intégration MediaPipe** pour :
 - Analyse des 468 landmarks faciaux (Face Mesh)
 - Détection de 8 gestes de la main
 - Boost contextuel des émotions difficiles
4. **Application temps réel optimisée (V3)** avec :
 - Test-Time Augmentation (TTA)
 - Prétraitement CLAHE
 - Lissage temporel exponentiel
 - Mapping emoji étendu (15+ emojis)
5. **Support multi-datasets** permettant d'utiliser FER+, AffectNet, ou RAF-DB

11.1 Perspectives

Pour aller plus loin, les améliorations suivantes pourraient être envisagées :

- Utilisation de transfer learning (VGGFace, ResNet pré-entraîné sur visages)
- Architecture avec attention mechanism (Transformer, CBAM)
- Détection multi-tâches (émotion + âge + genre)
- Analyse de la valence et de l'arousal (modèle dimensionnel)
- Déploiement sur edge devices (quantification INT8, pruning)
- Intégration d'un modèle audio pour l'analyse multimodale

A Structure du Projet

```
Final_project/
|-- app.py                      # Application temps reel (version initiale)
|-- app_extended.py              # Version avec MediaPipe Face Mesh
|-- app_extended_v2.py           # Version avec gestes de la main
|-- app_v3.py                    # Version optimisee (TTA, CLAHE, modulaire)
|-- model.py                     # Architecture CNN (RGB 75x75)
|-- train_affectnet.py           # Entrainement sur Balanced AffectNet
|-- train.py                      # Script d'entraînement (legacy FER2013)
|-- dataset_affectnet.py         # Dataset Balanced AffectNet
|-- dataset.py                   # Dataset FER2013 (legacy)
|-- download_datasets.py          # Script de telechargement des datasets
|-- emotion_model.pth            # Poids du modele entraine
|-- emotion_model_best.pth       # Meilleur modele (early stopping)
|-- data/
|   +- affectnet/                # Dataset Balanced AffectNet
|       |-- train/
|           |-- Anger/
|           |-- Contempt/
|           |-- Disgust/
```

```

|       |   |-- Fear/
|       |   |-- Happy/
|       |   |-- Neutral/
|       |   |-- Sad/
|       |   +-+ Surprise/
|       |   |-- val/
|       |   +-+ test/
|-- report/
|   +-+ report.tex          # Ce rapport
+-+ README.md

```

B Références

1. Goodfellow, I. J., et al. (2013). "Challenges in representation learning : A report on three machine learning contests." *ICML Workshop*.
2. Barsoum, E., et al. (2016). "Training Deep Networks for Facial Expression Recognition with Crowd-Sourced Label Distribution." *ACM ICMI*.
3. Mollahosseini, A., et al. (2017). "AffectNet : A Database for Facial Expression, Valence, and Arousal Computing in the Wild." *IEEE Trans. Affective Computing*.
4. Li, S., et al. (2017). "Reliable Crowdsourcing and Deep Locality-Preserving Learning for Expression Recognition in the Wild." *CVPR*.
5. He, K., et al. (2015). "Delving Deep into Rectifiers : Surpassing Human-Level Performance on ImageNet Classification." *ICCV*.
6. Lin, T. Y., et al. (2017). "Focal Loss for Dense Object Detection." *ICCV*.
7. Zhang, H., et al. (2018). "mixup : Beyond Empirical Risk Minimization." *ICLR*.
8. Yun, S., et al. (2019). "CutMix : Regularization Strategy to Train Strong Classifiers with Localizable Features." *ICCV*.
9. Lugaresi, C., et al. (2019). "MediaPipe : A Framework for Building Perception Pipelines." *arXiv preprint arXiv :1906.08172*.
10. Buslaev, A., et al. (2020). "Albumentations : Fast and Flexible Image Augmentations." *Information*, 11(2), 125.