

PROJET THÉMATIQUE

2ÈME ANNÉE - ELECTRONIQUE

Compétition Kaggle IA/ML en Vision par ordinateur

Élèves :

Badr LAMBARKI EL ALLIOUI
Amer SAIDI

Encadrant :

Mr. DONIAS

Table des matières

1	Introduction	3
1.1	Objectifs du projet	3
1.2	Présentation Kaggle	3
1.2.1	Compétitions	3
1.2.2	Datasets	3
1.2.3	Kernels/Notebooks	3
1.2.4	Formations	3
2	Choix du projet	4
3	Formalités de Kaggle sur notre projet	4
3.1	Soumission attendue	5
3.2	Métrique d'évaluation	5
3.3	Caractéristiques de la Dataset	5
4	Approches des Participants	5
4.1	Modèles Couramment Utilisés	5
4.2	Observations	6
5	Nos propres approches	6
6	CNN classique	7
6.1	Définition d'un resnet50	7
6.2	Explication du Code	8
6.3	Choix du Modèle	8
6.4	Prétraitement et Augmentation des Données	8
6.5	Entraînement	8
6.6	Planification du Taux d'Apprentissage (Learning Rate)	9
6.7	Fonction de Décision par Sigmoides	10
6.8	analyse des résultat	11
6.9	test sur autre dataset	12
7	Fusion entre deux modèles	14
7.1	Fusion ConvNeXt + Swin Transformer	14
7.2	Autocast (Mixed Precision)	15
7.3	Rôle du Global Average Pooling (GAP) en Mixed Precision	16
7.4	Présentation des résultats	16
8	VAE et analyse fréquentielle (FIRE)	18
8.1	Principe du Variational Autoencoder (VAE)	18
8.2	Application à la détection IA/réel : méthode FIRE	18
8.3	Détails de l'implémentation	19
8.4	Synthèse des simplifications mémoire	19
8.5	Analyse des courbes d'apprentissage	20
9	Faillle dans la dataset de la compétition	20
10	Conclusion	21

11 Annexes	22
11.1 Fire méthode :	22
11.2 Code resnet :	30
11.3 Convnext+swin	41

1 Introduction

1.1 Objectifs du projet

Le projet a pour but d'appliquer les concepts théoriques de l'IA étudiée en S7 et son application sur traitement des images étudiée tout au long de notre cursus. Les principales objectifs du projet :

- Bilan des compétitions Kaggle portant sur la vision par ordinateur.
- Identification d'une problématique d'intérêt et abordable.
- Inscription dans la compétition, téléchargement des données et de codes disponibles proposés par des compétiteurs.
- Confrontation des résultats obtenus avec ceux visibles dans le tableau des scores.
- Réflexion sur une amélioration de la stratégie engagée et, éventuellement, soumission de nouveaux résultats.

1.2 Présentation Kaggle

Kaggle est une plateforme en ligne dédiée à la data science, au machine learning et à l'analyse de données. Elle est très populaire auprès des data scientists, des chercheurs et des passionnés d'IA. Les Principales Fonctionnalités :

1.2.1 Compétitions

Les **compétitions** Kaggle sont des défis sponsorisés par des entreprises ou des institutions de recherche, proposant souvent des prix aux meilleures solutions. Ces compétitions couvrent divers domaines tels que : Prédiction de fraudes (finance), reconnaissance d'images (computer vision), traitement du langage (NLP), optimisation de modèles.

Elles sont classées par niveau de difficulté, ce qui les rend accessibles aussi bien aux débutants qu'aux experts. L'aspect compétitif motive les participants à innover, tandis que les problèmes abordés sont directement utiles aux organisations qui les proposent.

1.2.2 Datasets

Kaggle héberge une **vaste collection de datasets publics**, gratuits de type variés (Texte , image ...) et couvrant des domaines variés.

1.2.3 Kernels/Notebooks

La plateforme offre un environnement de développement intégré avec :

- Accès à des ressources de calcul (GPU/TPU) . Chaque utilisateur vérifié a une quota de 30h par semaines sur CPU T4 x2 ou CPU P100 et 20h sur TPU.
- Partage et collaboration facilités, très utile pour les groupes

1.2.4 Formations

Kaggle propose des parcours d'apprentissage en : Machine Learning et Deep Learning, visualisation de données ,python et R pour la data science.

2 Choix du projet

La recherche du sujet constitue l'étape fondamentale pour développer un projet solide et pertinent. Nous avons exploré trois axes principaux : les thèmes d'actualité, les défis en machine learning et les applications en traitement d'image.

Parmi les nombreuses pistes étudiées, la détection de langage des signes présentait un intérêt particulier grâce à son utilité sociale pour les personnes non verbales.

Cependant, notre choix final s'est porté sur un projet émergent : la détection d'images générées par IA versus images humaines. Ce sujet, particulièrement actuel avec l'essor des générateurs d'images, offre à la fois un défi technique stimulant et des applications concrètes dans la lutte contre la désinformation.

Aujourd'hui, on a un vrai problème : les images créées par des intelligences artificielles deviennent de plus en plus difficiles à repérer. Pourquoi ? Deux explications possibles :

- Soit les sites internet bloquent ces images artificielles
- Soit les IA sont devenues si fortes que leurs images sont maintenant parfaites, comme si un humain les avait faites

C'est exactement pour ça que notre projet est utile : on veut créer un outil capable de faire la différence entre les vraies images et celles fabriquées par les IA, même les plus récentes.



FIGURE 1 – AI vs Human

3 Formalités de Kaggle sur notre projet

En premier temps, la compétition a commencé le 15 Janvier 2025 et est terminée le 8 Mars 2025. Ce qui est aligné avec la durée de notre projet.



FIGURE 2 – Durée du projet

Kaggle met à disposition deux ensembles de données :

- `train.csv` - Ensemble d'entraînement contenant 79 950 images
- `test.csv` - Ensemble de test contenant 19 986 images

Format des fichiers CSV :

```
id - train_data/a6dcb93f596a43249135678dfcfc17ea.jpg
label - 1
```

3.1 Soumission attendue

Un fichier CSV contenant les prédictions avec :

- `image_id` : Identifiant unique de l'image
- `label` : 0 (image humaine) ou 1 (image IA)

3.2 Métrique d'évaluation

La métrique principale est le F1-score, moyenne harmonique entre précision et rappel :

$$F1 = 2 \times \frac{\text{Précision} \times \text{Rappel}}{\text{Précision} + \text{Rappel}}$$
$$\text{Précision} = \frac{\text{Vrais positifs}}{\text{Vrais positifs} + \text{Faux positifs}}$$
$$\text{Rappel} = \frac{\text{Vrais positifs}}{\text{Vrais positifs} + \text{Faux négatifs}}$$

3.3 Caractéristiques de la Dataset

- La soumission répétée du même code donne les mêmes résultats (données de test constantes)
- Deux ensembles de test existent :
 - Ensemble visible pour tester les codes
 - Ensemble caché pour le calcul final

4 Approches des Participants

Les soumissions publiques des autres compétiteurs nous ont été précieuses pour comprendre les attentes des organisateurs. La majorité des solutions proposées reposent sur des architectures CNN classiques, avec des variations au niveau des classifieurs.

4.1 Modèles Couramment Utilisés

- **ResNet-18/50 (2015)**
Architecture CNN profonde avec connexions résiduelles. Pré-entraînée sur ImageNet.
But : Extraction de caractéristiques visuelles de base.
- **RegNet-18/50 (2020, Facebook AI)**
CNN optimisé via une recherche architecturale automatisée. Entraîné sur ImageNet.
But : Bons compromis vitesse/précision.

- **EfficientNet** (2019, Google)
CNN scalable utilisant un coefficient de mise à l'échelle. Pré-entraînement ImageNet.
But : Performance optimisée pour ressources limitées.
- **ViT (Vision Transformer)** (2020)
Adaptation des transformers au traitement d'images. Entraîné sur JFT-300M/ImageNet.
But : Capturer des relations globales dans l'image.
- **Swin Transformer** (2021, Microsoft)
Transformer hiérarchique avec fenêtres glissantes. Pré-entraînement ImageNet-22K.
But : Réduire la complexité computationnelle des ViT.
- **ConvNeXt** (2022)
Modernisation des CNN inspirée par les transformers. Dataset ImageNet.
But : Combiner forces des CNN et transformers.
- **YOLOv11** (Variante de YOLO)
Détecteur one-shot adapté à la classification. COCO/ImageNet.
But : Efficacité temps réel (sur-optimisé pour ce projet).
- **VAE Classifiers**
Variational Autoencoder avec couche de classification. Entraînement end-to-end.
But : Détection d'anomalies dans les images générées.

4.2 Observations

- Les modèles vision transformers (ViT, Swin) montrent des résultats prometteurs pour cette tâche
- L'approche VAE est intéressante mais peu représentée dans les soumissions
- EfficientNet reste populaire pour son efficacité computationnelle

5 Nos propres approches

De ce qu'on a tiré des approches précédentes, c'est qu'il existe au moins trois grandes familles de réseaux de neurones :

- CNN (Convolutional Neural Networks) : sont efficaces pour extraire des caractéristiques locales comme les bords, les textures ou les formes.
- Transformers : à l'origine conçus pour le traitement du langage naturel, ils sont aujourd'hui largement utilisés aussi en vision par ordinateur . Leur force réside dans l'utilisation du mécanisme d'attention, qui permet de modéliser des relations globales au sein des données.
- VAE (Variational Autoencoders) : ce sont des autoencodeurs probabilistes reposant sur une architecture encodeur-décodeur. L'encodeur projette les données dans un espace latent, et le décodeur les reconstruit à partir de cet espace. Ils sont particulièrement utiles pour : la compression, la génération d'images réalistes, la détection d'anomalies.

6 CNN classique

6.1 Définition d'un resnet50

ResNet50 (Residual Network with 50 layers) est un réseau de neurones convolutifs (CNN) profond introduit par Microsoft Research en 2015 dans le cadre du paper "**Deep Residual Learning for Image Recognition**".

Principe Fondamental : Apprentissage Résiduel Pour une couche donnée dans ResNet50 :

$$y = F(x, \{W_i\}) + x$$

Où :

- x : entrée du bloc résiduel
- $\{W_i\}$: poids des convolutions dans le bloc
- $F(x, \{W_i\})$: transformation apprise par les couches (le "résidu")
- $+x$: connexion directe (skip connection)
- y : sortie du bloc

Fonctionnement d'apprentissage Résiduel

1. Si $F(x) \approx 0$ (les couches n'apprennent rien d'utile) :

$$y = 0 + x = x$$

⇒ Le bloc se comporte comme une simple connexion directe.

2. Le gradient peut circuler directement via x :

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial y} \cdot \left(1 + \frac{\partial F}{\partial x}\right)$$

⇒ Le terme "1" garantit que le gradient ne s'annule jamais complètement.

l'apprentissage résiduel permet d'empiler des centaines de couches dans un réseau sans que le gradient ne disparaisse (problème du vanishing gradient), grâce à des connexions directes (skip connections) qui préservent l'information d'origine.

En pratique dans ResNet50 :

- x : une image (caractéristiques extraites par les couches précédentes)
- $F(x)$: apprend à corriger x pour mieux détecter les artefacts des images générées par IA
- $y = F(x) + x$: combine l'information originale et les corrections

Pourquoi ResNet50 pour la classification IA vs Réel ?

- Transfer Learning efficace : Pré-entraîné sur ImageNet, il a appris des caractéristiques génériques utiles pour la détection d'artefacts dans les images IA.

- Robustesse aux variations : Grâce aux résidus, il conserve une bonne sensibilité même sur des différences subtiles (ex : textures non naturelles dans les images générées).
- Équilibre complexité/performance : Plus léger que ResNet101/152, mais suffisamment profond pour capturer des motifs discriminants.

6.2 Explication du Code

Le code développé repose sur l'utilisation de PyTorch pour la classification des images générées par IA contre les images réelles. Implémentation du Modèle

6.3 Choix du Modèle

Nous avons utilisé l'architecture **ResNet50**, un réseau convolutionnel profond basé sur l'apprentissage résiduel. Ce choix se justifie par sa robustesse et sa capacité à éviter le problème du *vanishing gradient* grâce aux *skip connections*. Le modèle pré-entraîné sur ImageNet (`ResNet50_Weights.IMAGENET1K_V2`) est réutilisé en tant que base.

Seule la couche de classification finale a été modifiée pour l'adapter à notre tâche binaire :

```
1 model.fc = nn.Sequential(  
2     nn.Linear(num, 512),  
3     nn.ReLU(),  
4     nn.Dropout(0.3),  
5     nn.Linear(512, 1)  
6 )
```

Listing 1 – Modification de la couche finale

L'entraînement a été effectué uniquement sur la dernière couche ajoutée au ResNet. Nous avons également testé l'impact de l'entraînement sur plusieurs couches (comme le layer4), mais avons observé que cela n'apportait pas d'amélioration significative des performances.

6.4 Prétraitement et Augmentation des Données

Les images sont normalisées et soumises à plusieurs transformations aléatoires pour améliorer la robustesse du modèle :

- **Train Transform** : redimensionnement aléatoire, rotation, jitter de couleur, inversion horizontale, etc.
- **Validation Transform** : recadrage central et normalisation.

6.5 Entraînement

Le modèle est entraîné avec une perte de type `BCEWithLogitsLoss` et un optimiseur Adam. L'entraînement est effectué sur plusieurs époques, en validant la performance à chaque époque :

Fonction de perte et optimisation :

Le modèle est entraîné à l'aide de la fonction de perte `BCEWithLogitsLoss`, adaptée aux tâches de classification binaire. Cette fonction combine la sigmoïde et l'entropie croisée binaire, ce qui la rend numériquement plus stable que l'application séparée d'une fonction sigmoïde suivie d'une `BCELoss`. La formule mathématique est la suivante :

$$\mathcal{L}(x, y) = -[y \cdot \log(\sigma(x)) + (1 - y) \cdot \log(1 - \sigma(x))]$$

où :

- x est la sortie (logit) du modèle,
- $y \in \{0, 1\}$ est le label cible,
- $\sigma(x) = \frac{1}{1 + e^{-x}}$ est la fonction sigmoïde.

Cette formulation est particulièrement adaptée à la classification binaire car :

- elle permet d'interpréter la sortie comme une probabilité appartenant à la classe 1,
- elle pénalise fortement les erreurs de classification,
- elle est compatible avec l'entraînement direct à partir de logits.

L'optimisation du modèle est réalisée à l'aide de l'algorithme `Adam`, qui adapte dynamiquement le taux d'apprentissage pour chaque paramètre à l'aide des moments d'ordre 1 et 2 des gradients :

$$\theta_{t+1} = \theta_t - \alpha \cdot \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon}$$

où :

- \hat{m}_t et \hat{v}_t sont les estimations corrigées du moment et du carré du moment,
- α est le taux d'apprentissage,
- ϵ est un terme de stabilisation.

L'entraînement est effectué sur plusieurs époques, avec validation des performances (accuracy, F1-score, etc.) à la fin de chaque époque pour suivre la progression du modèle.

6.6 Planification du Taux d'Apprentissage (Learning Rate)

Nous utilisons un *scheduler* cosinusoidal (`CosineAnnealingLR`) pour adapter dynamiquement le taux d'apprentissage. Cette approche varie le *learning rate* η selon :

$$\eta_t = \eta_{min} + \frac{1}{2}(\eta_{max} - \eta_{min}) \left(1 + \cos \left(\frac{t}{T_{max}} \pi \right) \right)$$

Avantages :

- Réduction progressive des oscillations en fin d'entraînement
- Évite les minima locaux grâce aux variations cycliques
- Pas de paramètres supplémentaires à optimiser

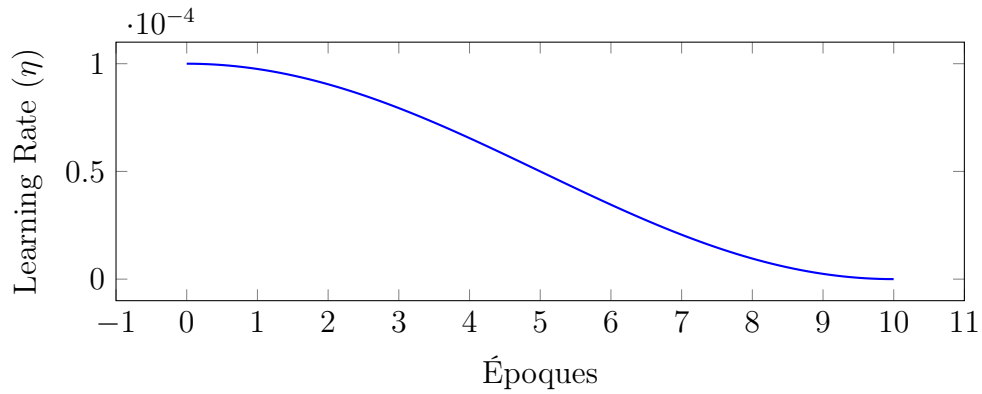


FIGURE 3 – Évolution cosinusoidale du taux d'apprentissage ($\eta_{max} = 10^{-4}$, $T_{max} = 10$).

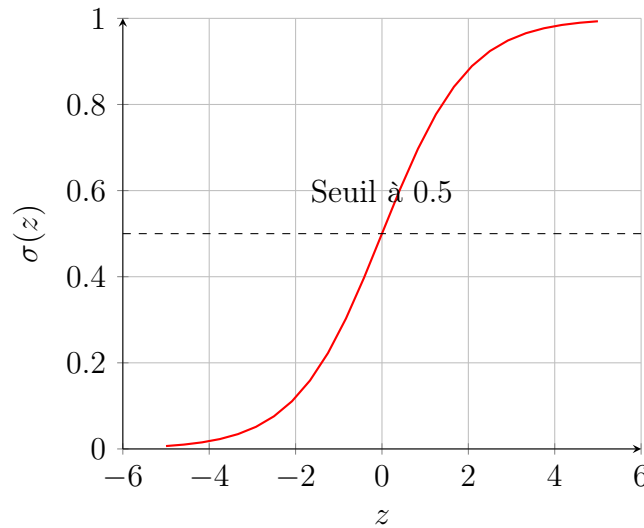


FIGURE 4 – Fonction sigmoïde et seuil de décision.

6.7 Fonction de Décision par Sigmoides

Pour la classification binaire, nous utilisons une sigmoïde avec seuil à 0.5 :

$$\sigma(z) = \frac{1}{1 + e^{-z}} \quad \text{où } z = \text{logits du modèle}$$

Processus de décision :

$$\text{Classe} = \begin{cases} 0 & (\text{R  el}) & \text{si } \sigma(z) \leq 0.5 \\ 1 & (\text{IA}) & \text{si } \sigma(z) > 0.5 \end{cases}$$

Justification :

- Interpr  tation probabiliste naturelle ($P(y = 1|x) = \sigma(z)$)
- Compatible avec la fonction de perte `BCEWithLogitsLoss` (plus stable num  riquement que `BCELoss` + sigmo  de manuelle)
- Seuil 0.5 optimal pour des classes   quilibr  es (50% IA / 50% R  el dans notre cas)

6.8 analyse des résultat

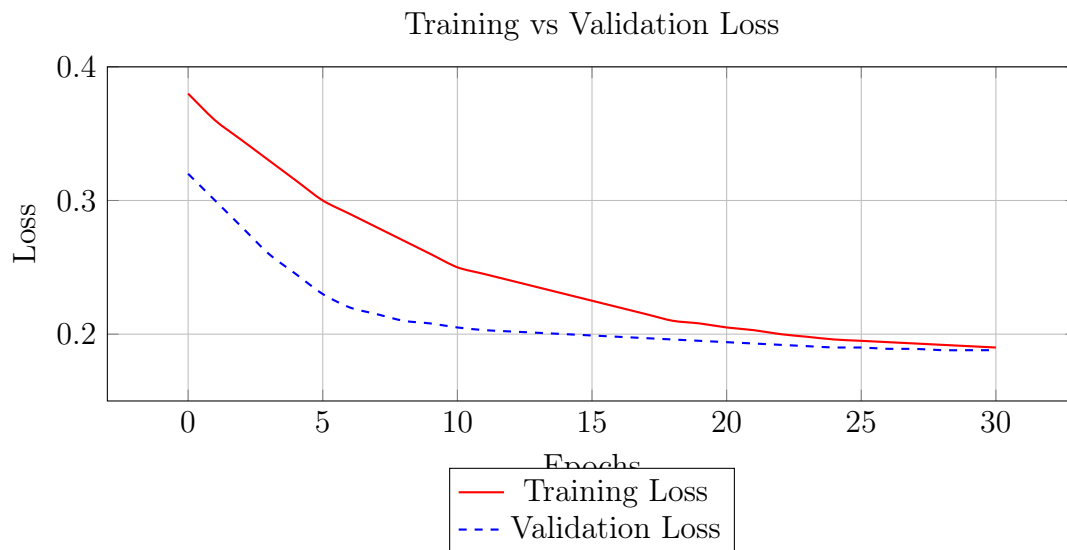


FIGURE 5 – Comparaison entre la **loss d'entraînement** (rouge) et la **loss de validation** (bleu). Les deux courbes convergent vers 0,19 à l'epoch 30.

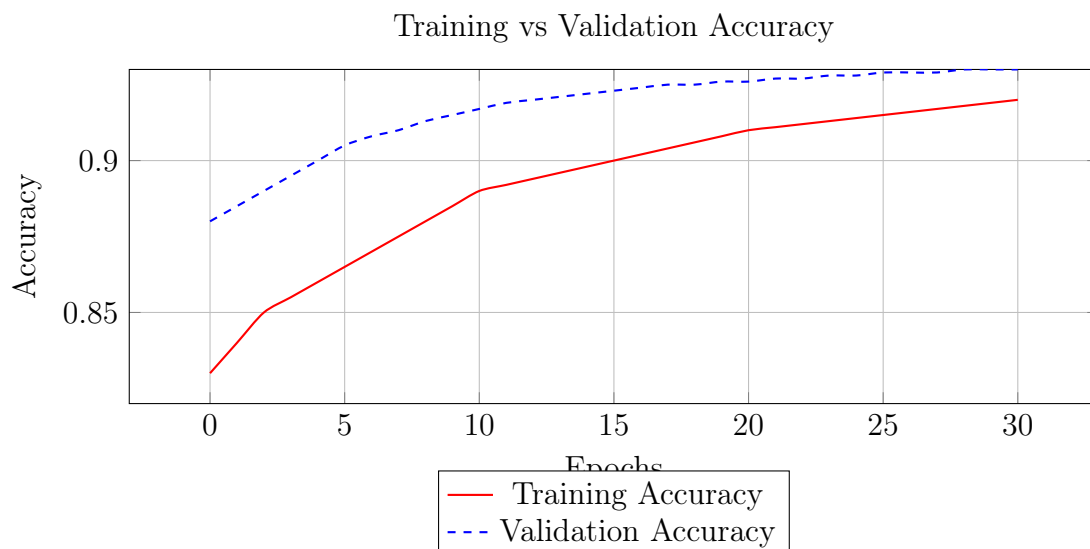


FIGURE 6 – Comparaison entre la **précision d'entraînement** (rouge) et la **précision de validation** (bleu). Les deux courbes convergent vers 93 % à l'epoch 30.

La courbe ROC (Receiver Operating Characteristic) permet d'évaluer la capacité du modèle à distinguer entre les deux classes (images réelles et images générées par IA), en traçant le taux de vrais positifs (*True Positive Rate*) contre le taux de faux positifs (*False Positive Rate*) à différents seuils de décision.

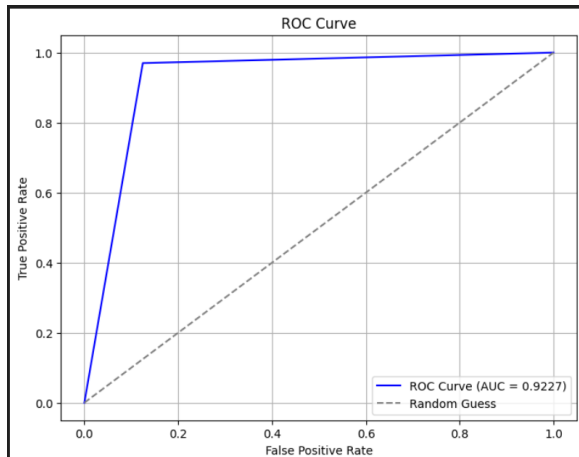


FIGURE 7 – Courbe ROC

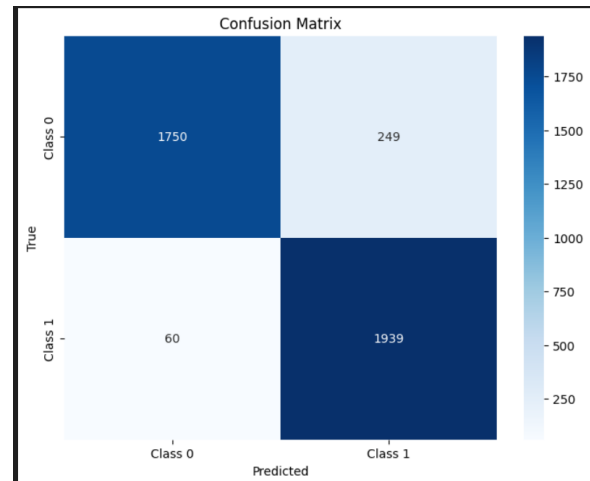


FIGURE 8 – Matrice de confusion

FIGURE 9 – Résultats du modèle ResNet

L'aire sous la courbe (AUC) mesurée est de **0.92**, ce qui indique une excellente capacité de séparation. Plus l'AUC est proche de 1, meilleure est la performance du modèle. Une AUC de 0.5 correspondrait à un modèle aléatoire, tandis qu'une AUC de 1.0 indique une séparation parfaite.

Les performances du modèle sur l'ensemble de validation sont résumées ci-dessous à l'aide du `classification report` de `scikit-learn` :

Classe	Précision	Rappel	F1-score	Support
0 (Image Réelle)	0.97	0.88	0.92	1999
1 (Image IA)	0.89	0.97	0.93	1999
Accuracy	0.92 (3998 images)			
Macro Avg	0.93	0.92	0.92	3998
Weighted Avg	0.93	0.92	0.92	3998

TABLE 1 – Rapport de classification du modèle sur l'ensemble de validation

On observe que le modèle présente une bonne précision globale (92%), avec une légère asymétrie entre les deux classes : les images réelles sont parfois mal classées en IA, ce qui se reflète dans leur rappel légèrement inférieur (88%). Toutefois, les scores F1 équilibrés (0.92 et 0.93) montrent une performance stable.

le score de notre proposition est 0.71653

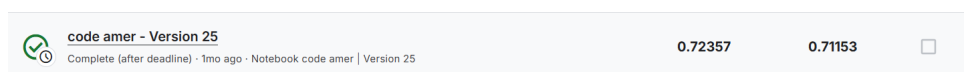


FIGURE 10

6.9 test sur autre dataset

Exemple des images qui sont détecté par le resnet :



FIGURE 11



FIGURE 12

FIGURE 13 – Images IA détectée par Resnet50

Détection des images générées par ResNet50 Les images générées, bien qu'esthétiquement réalistes à première vue, présentent plusieurs caractéristiques qui les rendent détectables par le réseau de neurones convolutif ResNet50. Premièrement, les textures hyperréalistes du pelage, avec des motifs de fourrure anormalement réguliers, lisses et symétriques, diffèrent des irrégularités naturelles présentes dans des photos réelles. Deuxièmement, les couleurs pastel irréalistes, la lumière ambiante irréprochable et les arrière-plans flous artificiellement uniformes génèrent des artefacts statistiques dans les couches profondes du réseau. ResNet50, grâce à ses nombreux blocs résiduels profonds, capte ces subtilités dans les motifs de bas niveau (bords, gradients) et de haut niveau (textures, cohérence contextuelle), ce qui permet au modèle de distinguer efficacement les images IA des vraies images.

Description du Dataset Le dataset utilisé contient 60 000 images équilibrées :

- **30 000 images IA :**
 - 10 000 de *Stable Diffusion*
 - 10 000 de *MidJourney*
 - 10 000 de *DALL-E*
- **30 000 images réelles :**
 - 22 500 de *Pexels/Unsplash* (photos)
 - 7 500 de *WikiArt* (art numérique)

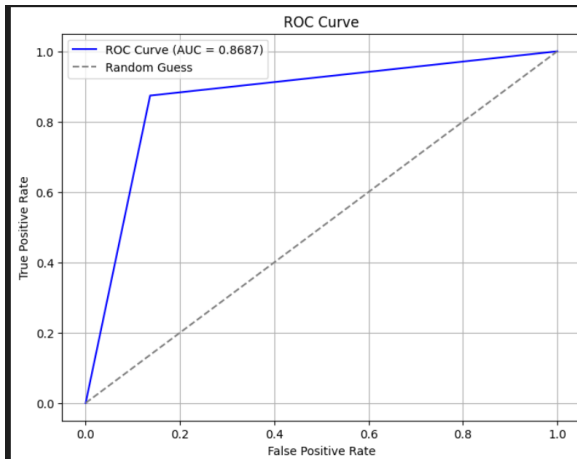


FIGURE 14 – Courbe ROC

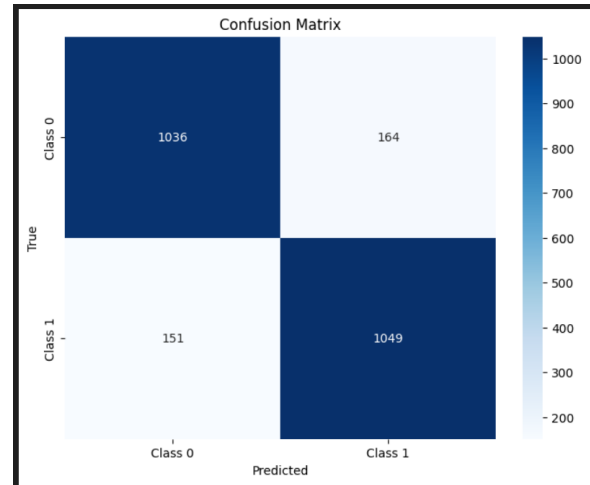


FIGURE 15 – Matrice de confusion

FIGURE 16 – Résultats du modèle ResNet

Performance et Analyse du Modèle Notre modèle ResNet50 démontre des performances solides sur le dataset de test, avec une précision globale de **87%** et une AUC de **0.8687**, indiquant une bonne capacité à distinguer images réelles et générées par IA. La matrice de confusion révèle un équilibre entre les classes (F1-score de 0.87 pour chacune), bien que **15% d'erreurs** persistent, notamment sur des cas ambigus. Ces résultats suggèrent que le modèle capture efficacement les artéfacts typiques des IA (textures, incohérences), mais pourrait être amélioré par un fine-tuning ciblé sur les images "borderline" (ex : IA photoréalistes) et l'ajout de données difficiles. Cette performance valide l'approche pour des applications de détection de contenu synthétique, avec une marge de progression vers les 90%+.

Classe	Précision	Rappel	F1-score	Support
0 (Image Réelle)	0.87	0.86	0.87	1200
1 (Image IA)	0.86	0.87	0.87	1200
Accuracy	0.87 (2400 images)			
Macro Avg	0.87	0.87	0.87	2400
Weighted Avg	0.87	0.87	0.87	2400

TABLE 2 – Rapport de classification sur l'ensemble de test

7 Fusion entre deux modèles

7.1 Fusion ConvNeXt + Swin Transformer

Nous combinons les forces de **ConvNeXt Large**, expert en textures locales, et du **Swin Transformer Base**, performant sur les relations globales. ConvNeXt exploite la structure spatiale avec l'efficacité des CNNs, tandis que Swin utilise l'attention à fenêtres glissantes pour capter les dépendances à longue portée. Leur complémentarité améliore la robustesse du modèle, en capturant à la fois les détails fins et la cohérence contextuelle, essentielle pour distinguer images réelles et générées.

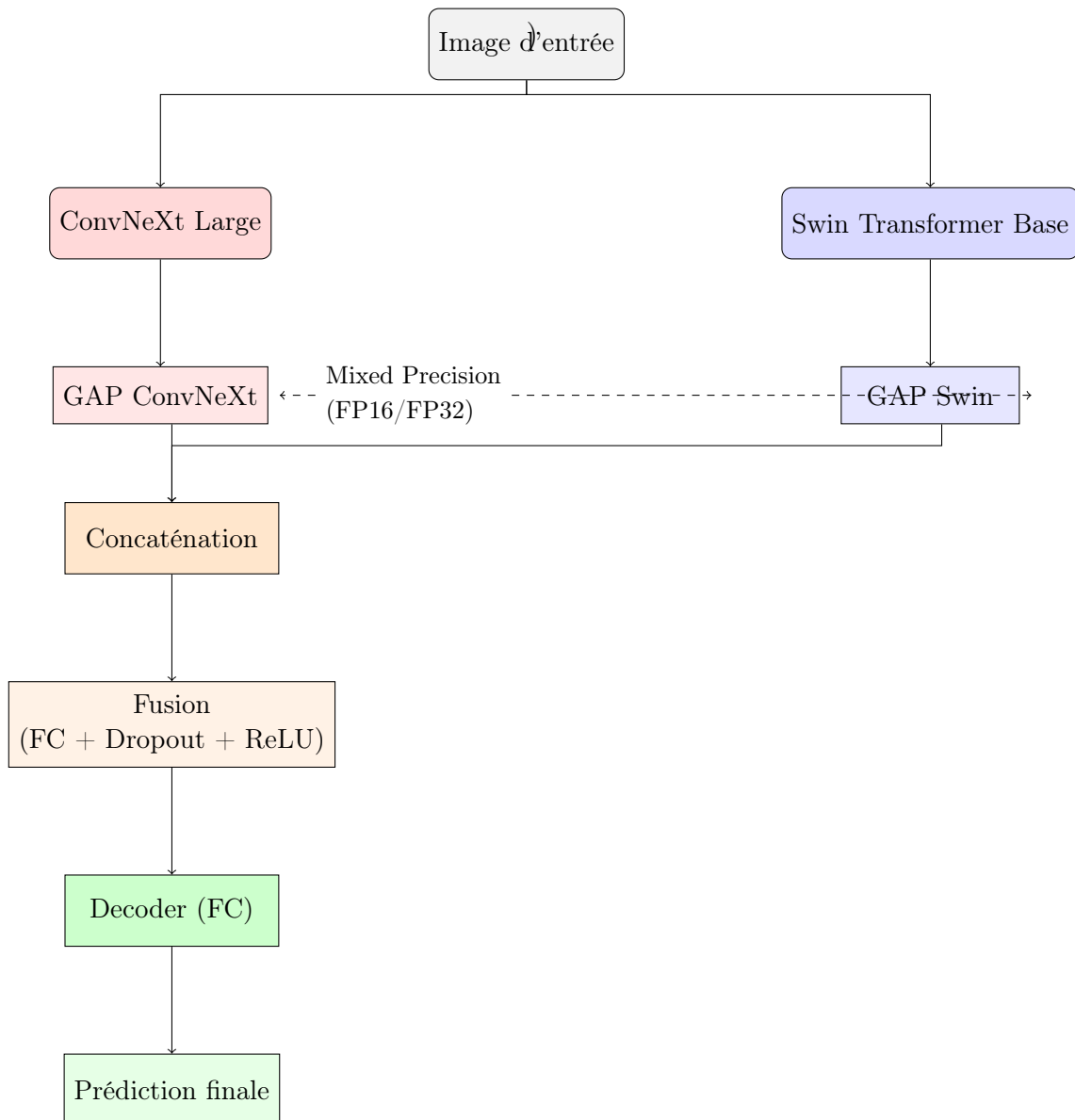


FIGURE 17 – Fusion des caractéristiques extraites par ConvNeXt Large et Swin Transformer Base, avec apprentissage en précision mixte.

7.2 Autocast (Mixed Precision)

L'autocast (mixed precision) a pour objectif d'accélérer l'entraînement des réseaux de neurones sans perte de précision en combinant judicieusement le `float16` pour les calculs rapides (notamment les opérations matricielles) et le `float32` pour les opérations sensibles (comme les sommes et normalisations). Cette technique repose sur plusieurs principes clés : le mixed precision utilise le `float16` pour accélérer l'entraînement tout en conservant le `float32` pour la stabilité numérique, particulièrement pour le calcul de la loss et les opérations de réduction qui sont systématiquement sauvegardées en `float32` pour éviter les erreurs numériques. Le `GradScaler` joue un rôle crucial en amplifiant les gradients trop petits avant la rétropropagation pour prévenir les *underflows*. À noter que les CPUs ne bénéficient pas de cette accélération, traitant le `float16` aussi lentement que le `float32` par manque de circuits dédiés. Enfin, sans `GradScaler`, les gradients en `float16` risquent

de devenir nuls, ce qui paralyserait complètement l'apprentissage.

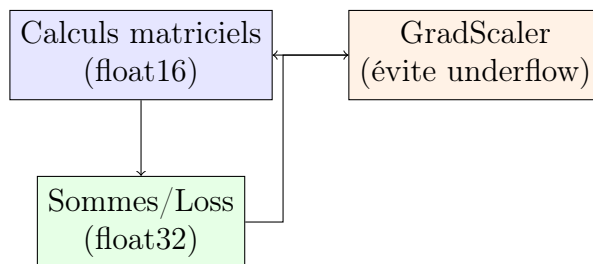


FIGURE 18 – Workflow du mixed precision. Les gradients sont rescaling par le GradScaler avant la rétropropagation.

7.3 Rôle du Global Average Pooling (GAP) en Mixed Precision

Le Global Average Pooling (GAP) est essentiel pour stabiliser l'entraînement en mixed precision. En réduisant chaque carte de caractéristiques à sa valeur moyenne par canal (transformant par exemple un tenseur de forme $[1, 1536, 7, 7]$ en $[1, 1536]$), le GAP élimine les opérations conflictuelles entre `float16` et `float32` qui surviennent dans les couches convolutives. Cette réduction dimensionnelle diminue radicalement le risque d'*overflow/underflow* en limitant le nombre d'opérations sensibles (de 12544 valeurs à 256 moyennes par image), tout en optimisant l'utilisation de la mémoire GPU. Concrètement, l'opération `x_convnext.unsqueeze(2).view(x_convnext.size(0), -1)` garantit la compatibilité avec `BCEWithLogitsLoss` en maintenant une sortie de dimension $[batch, 1]$ pour la classification binaire, tout en préservant la stabilité numérique grâce au mélange contrôlé des précisions.

7.4 Présentation des résultats



FIGURE 19

Amélioration par fusion L'architecture hybride (ConvNeXt + Swin Transformer) atteint un score de **0.8** contre **0.7** pour ResNet seul grâce à sa capacité à capturer simultanément :

- **Textures locales** via ConvNeXt (optimisé pour les motifs bas-niveau)
- **Relations globales** via Swin Transformer (attention aux artefacts contextuels)

Cette complémentarité permet de détecter à la fois les anomalies fines (ex : artefacts de génération dans les cheveux) et les incohérences structurelles (ex : symétrie parfaite anormale), là où ResNet seul ne percevait que les défauts locaux. Le gain de 10% reflète ainsi la synergie entre approches CNN et Transformer.

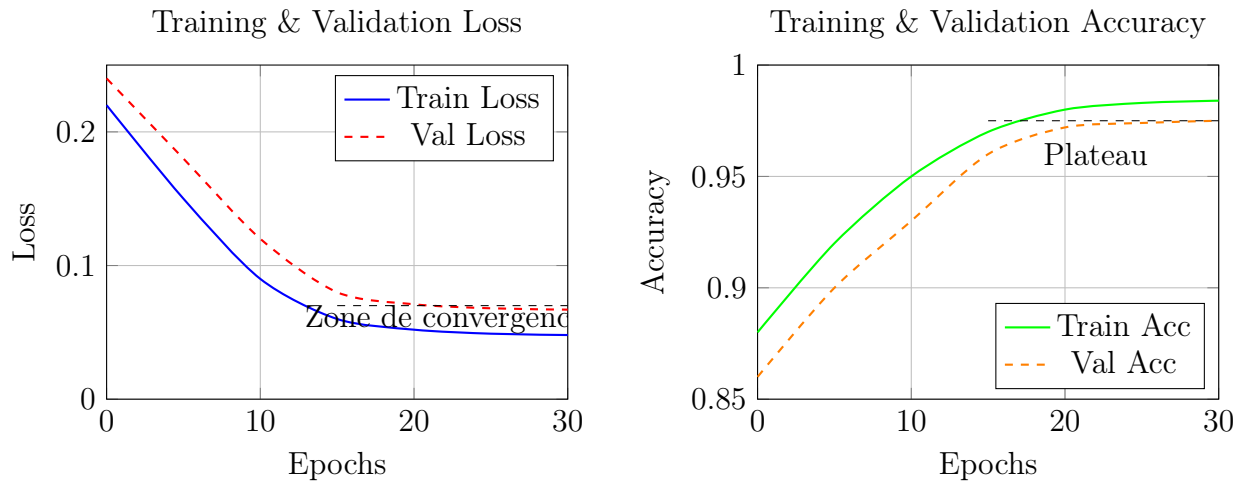


FIGURE 20 – Courbes d'apprentissage sur 30 epochs. Convergence stable atteinte vers l'epoch 15, avec un écart train/val inférieur à 2% en phase finale, indiquant un bon équilibre biais-variance. La loss finale se stabilise autour de 0.05 (train) et 0.07 (validation).

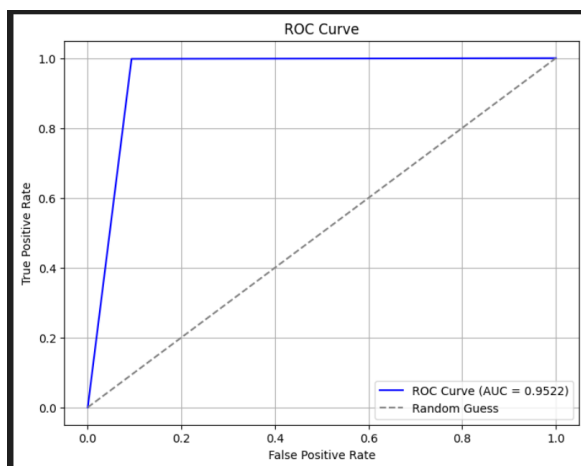


FIGURE 21 – Courbe AUC-ROC

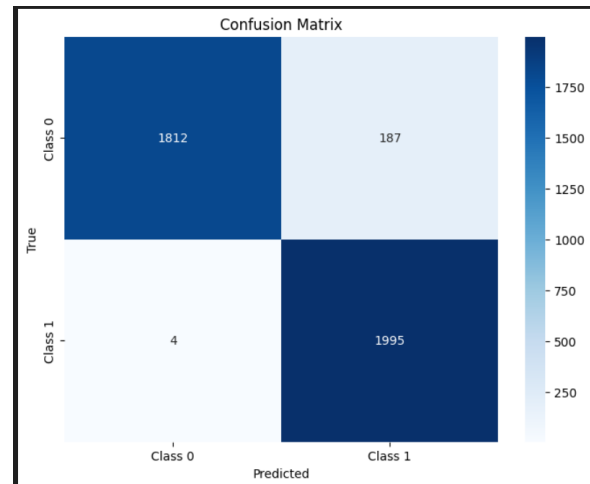


FIGURE 22 – Matrice de confusion

Performance du modèle Le modèle démontre d'excellentes performances avec une précision globale de **95%** (accuracy), confirmée par une **AUC de 0.952** sur la courbe ROC, indiquant une forte capacité à distinguer les classes. La matrice de confusion révèle un équilibre entre les deux classes :

- **Classe 0** (Réelle) : Précision parfaite (1.00) mais rappel légèrement inférieur (0.91)
- **Classe 1** (IA) : Rappel parfait (1.00) avec une précision de 0.91

Les F1-scores symétriques (**0.95** pour chaque classe) montrent une harmonie entre précision et rappel. Le faible taux de faux négatifs (4 cas seulement pour la Classe 1) souligne la robustesse du modèle pour la détection d'images générées par IA. Ces résultats suggèrent que le modèle maîtrise aussi bien les caractéristiques locales que globales discriminantes.

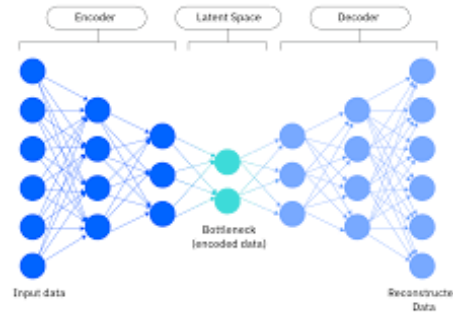


FIGURE 23 – VAE vu d'intérieur

8 VAE et analyse fréquentielle (FIRE)

8.1 Principe du Variational Autoencoder (VAE)

Le **Variational Autoencoder** (VAE) est un modèle génératif probabiliste basé sur une architecture encodeur-décodeur. L'encodeur transforme une image d'entrée \mathbf{x} en une distribution latente $\mathcal{N}(\mu, \sigma^2)$, tandis que le décodeur tente de reconstruire l'image à partir d'un échantillon \mathbf{z} tiré de cette distribution. Contrairement à un autoencodeur classique, le VAE encode chaque entrée comme une distribution (et non un point unique), ce qui permet de générer de nouvelles images réalistes par simple échantillonnage dans l'espace latent.

Mathématiquement, le VAE optimise la *fonction de perte variationnelle* suivante :

$$\mathcal{L}_{\text{VAE}} = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})] - D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))$$

où q_{ϕ} est l'encodeur, p_{θ} le décodeur, et D_{KL} la divergence de Kullback-Leibler.

8.2 Application à la détection IA/réel : méthode FIRE

La méthode **FIRE** (*Frequency-guided Reconstruction Error*) exploite une propriété fondamentale des modèles de diffusion : leur incapacité à reconstruire fidèlement les informations de fréquence moyenne dans les images réelles. L'algorithme procède ainsi :

1. **Décomposition fréquentielle** : l'image est transformée dans le domaine fréquentiel (FFT), puis un masque prédéfini isole la bande de fréquences intermédiaires (rayon 40 à 120 pixels autour du centre du spectre), initialisé sur la bande et régularisé par entraînement (notre valeur ajoutée) pour rester proche de ce masque.
2. **Création d'une image pseudo-générée** : on supprime ces fréquences intermédiaires pour obtenir une version de l'image ressemblant à une image générée par IA.
3. **Reconstruction par VAE** : le VAE reconstruit à la fois l'image originale et l'image pseudo-générée.
4. **Erreur de reconstruction** : on calcule l'erreur de reconstruction pour chaque image.
5. **Classification** : la différence entre les erreurs de reconstruction sert de descripteur pour un classifieur binaire. Dans notre cas, un ResNet18 modifié.

L'hypothèse clé est que les images réelles présentent une différence d'erreur de reconstruction bien plus marquée que les images générées par IA, car la suppression de la bande intermédiaire "efface" une information difficile à synthétiser pour les modèles de diffusion.

8.3 Détails de l'implémentation

L'implémentation suit fidèlement le schéma du papier, tout en étant adaptée aux contraintes de mémoire de l'environnement Kaggle :

- **Préparation des données :**
 - Sélection de 5 000 images réelles et 5 000 images générées, soit un total de 10 000 images, pour limiter la charge mémoire.
 - Split train/validation équilibré.
- **Augmentations de données :**
 - Utilisation d'augmentations fortes (flip, crop, jitter, blur, dropout) pour améliorer la robustesse sans augmenter la taille effective du dataset.
- **Architecture du module FMRE :**
 - Encodeur à 4 couches convolutionnelles.
 - Deux décodeurs indépendants (pour m_{mid} et m_{mid_c}), chacun suivi d'un `PixelShuffle`.
- **Utilisation du VAE de Stable Diffusion :**
 - Chargement du VAE pré-entraîné (`sd-vae-ft-mse`), gel des poids, utilisation en mode inference uniquement.
- **Classifieur :**
 - CNN léger (1 conv, pooling adaptatif, FC).

8.4 Synthèse des simplifications mémoire

Pour rendre l'entraînement réalisable sur Kaggle, plusieurs simplifications ont été nécessaires :

- **Réduction drastique du dataset :** 10 000 images au lieu de 80 000+.
- **Batch size minimal :** 4 images.
- **Formation du modèle en *mixed precision* :** calculs en float16.
- **Accumulation des gradients :** permet de simuler un batch effectif plus important sans surcharge mémoire.
- **VAE seul, pipeline de diffusion complet évité :** on ne réalise que la reconstruction par autoencodeur, ce qui réduit considérablement la mémoire utilisée.
- **Classifieur compact :** architecture CNN minimale pour la dernière étape.



FIGURE 24 – Score de l'approche

8.5 Analyse des courbes d'apprentissage

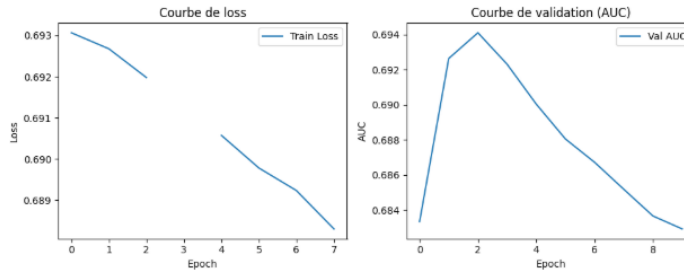


FIGURE 25 – Courbes d'apprentissage du modèle

Les courbes d'apprentissage montrent que :

- **La loss d'entraînement** diminue lentement et reste proche de 0.693, ce qui correspond à la loss d'un classifieur binaire aléatoire. Cela indique que le modèle peine à extraire des motifs discriminants entre images IA et réelles.
- **L'AUC de validation** atteint un maximum à 0.694, puis redescend rapidement pour stagner autour de 0.684.

En conclusion, l'utilisation du VAE dans la méthode FIRE permet de révéler une signature fréquentielle propre aux images réelles, mais la performance dépend fortement de la bonne extraction de cette information par le masque fréquentiel et de la qualité des données d'entraînement.

9 Faille dans la dataset de la compétition

Cependant, il est essentiel de souligner une faille méthodologique majeure dans la plupart des benchmarks de détection d'images générées par IA : le biais de taille et de compression des images. Comme l'ont montré plusieurs études récentes, de nombreux datasets pour la détection d'images IA comportent des différences systématiques entre images réelles et images générées, notamment au niveau de la résolution (les images IA étant souvent générées à des tailles fixes, tandis que les images réelles présentent une grande variété de dimensions) et du format de compression (JPEG pour les vraies, PNG pour les générées). Ce biais technique conduit les modèles à apprendre à distinguer des propriétés superficielles (dimensions, artefacts de compression) plutôt que de véritables signatures de génération. Dans notre cas, la métrique basée sur la moyenne de l'histogramme couleur 3D s'est révélée fortement corrélée à la taille de l'image, ce qui a permis d'atteindre une accuracy artificiellement élevée (jusqu'à 96%) simplement en détectant la résolution, et non l'origine réelle ou synthétique de l'image. Ce phénomène met en évidence la nécessité de construire des jeux de données équilibrés et débarrassés de ces biais pour garantir la robustesse et la généralisabilité des détecteurs d'images IA. Pour montrer que ces solutions ne sont pas consistantes, on l'a testé sur notre propre dataset et on a eu des résultats d'accuracy très bas. Ce qui est logique car d'un seul histogramme on peut créer plusieurs images.

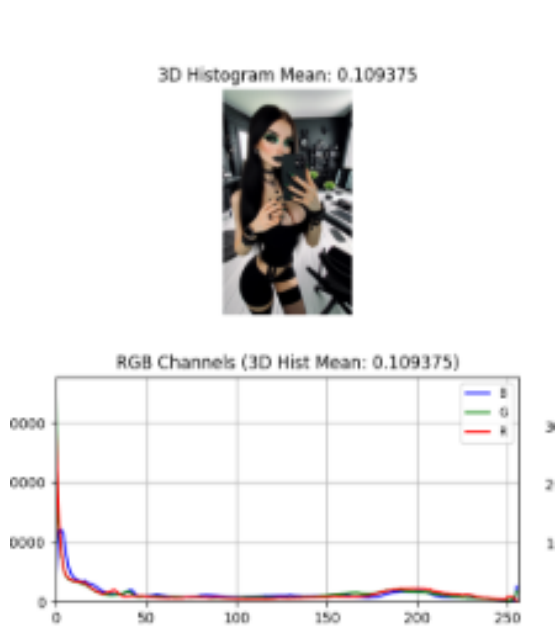


FIGURE 26 – Image IA



FIGURE 27 – Image réelle

FIGURE 28 – Comparaison entre une image générée par IA et une image réelle.



FIGURE 29 – Score de la solution biaisée

10 Conclusion

Ce projet a exploré plusieurs approches pour la détection d'images générées par IA, en comparant les forces et limites des architectures modernes. Les résultats montrent que :

- Les CNN classiques (ResNet50) offrent une solution robuste et généralisable, avec une précision de 92% sur des datasets équilibrés.
- Les architectures hybrides (ConvNeXt + Swin Transformer) améliorent les performances à 95% en combinant extraction locale et analyse contextuelle.
- L'approche fréquentielle FIRE avec VAE montre un potentiel pour les modèles de diffusion, mais nécessite de larges ressources en mémoire vu les optimisation qu'on a fait , pour surpasser les autres méthodes.
- La faille identifiée dans les datasets (biais de taille/résolution) souligne l'importance cruciale d'une curation rigoureuse des données pour éviter les artefacts techniques.

11 Annexes

11.1 Fire méthode :

```
1
2 import os
3 import random
4 import numpy as np
5 import pandas as pd
6 from PIL import Image
7 from tqdm import tqdm
8 import albumentations as A
9 from sklearn.metrics import roc_auc_score
10 import matplotlib.pyplot as plt
11
12 import torch
13 import torch.nn as nn
14 import torch.nn.functional as F
15 from torch.utils.data import Dataset, DataLoader
16 from diffusers import AutoencoderKL
17
18 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
19 print(f"Using device: {device}")
20
21 base_dir = '/kaggle/input/ai-vs-human-generated-dataset'
22 train_csv_path = os.path.join(base_dir, 'train.csv')
23 test_csv_path = os.path.join(base_dir, 'test.csv')
24
25 df_train = pd.read_csv(train_csv_path)
26 df_train['file_name'] = df_train['file_name'].apply(lambda x: os.path.
    join(base_dir, x))
27 df_test = pd.read_csv(test_csv_path)
28 df_test['id'] = df_test['id'].apply(lambda x: os.path.join(base_dir, x))
29
30 # Reduction 50000 images BALANCE (25000 reelles + 25000
    generated)
31 real_samples = df_train[df_train['label'] == 0].sample(n=25000,
    random_state=42)
32 fake_samples = df_train[df_train['label'] == 1].sample(n=25000,
    random_state=42)
33 df_train = pd.concat([real_samples, fake_samples]).sample(frac=1,
    random_state=42)
34
35 train_data = df_train.sample(frac=0.9, random_state=42)
36 val_data = df_train.drop(train_data.index)
37
38
39 class FireDataset(Dataset):
40     def __init__(self, paths, labels, img_size=256, train=True):
41         self.paths = paths
42         self.labels = labels
43         self.train = train
44         self.img_size = img_size
45
46         # CORRECTION: Param tres valides pour CoarseDropout
47         self.strong_aug = A.Compose([
```

```

48         A.HorizontalFlip(p=0.5),
49         A.RandomResizedCrop(size=(img_size, img_size), scale=(0.8,
1.0)),
50         A.ColorJitter(brightness=0.1, contrast=0.1, saturation=0.1,
hue=0.05, p=0.8),
51         A.GaussianBlur(blur_limit=(3, 7), p=0.5),
52         A.CoarseDropout(max_holes=1, max_height=32, max_width=32,
fill_value=0, p=0.3)
53     ])
54
55     def __len__(self):
56         return len(self.paths)
57
58     def __getitem__(self, idx):
59         img = Image.open(self.paths[idx]).convert('RGB')
60         if self.train and random.random() < 0.7:
61             img = np.array(img)
62             img = self.strong_aug(image=img)['image']
63             img = Image.fromarray(img)
64         img = img.resize((self.img_size, self.img_size))
65         img = torch.from_numpy(np.array(img)).permute(2, 0, 1).float() /
255.0
66         return img, torch.tensor(self.labels[idx], dtype=torch.float32)
67
68
69     def create_fixed_mask(size=256):
70
71         y, x = torch.meshgrid(torch.arange(size), torch.arange(size),
indexing='ij')
72         center = size // 2
73         dist = torch.sqrt((x - center)**2 + (y - center)**2)
74         mask = ((dist >= 40) & (dist <= 120)).float()
75         mask_c = 1.0 - mask
76         return mask.unsqueeze(0).unsqueeze(0), mask_c.unsqueeze(0).unsqueeze
(0) # [1,1,H,W]
77
78
79     class FMRE(nn.Module):
80         def __init__(self, init_mask, init_mask_c):
81             super().__init__()
82             # Masques pr d finis comme buffers (device-aware)
83             self.register_buffer('M_mid', init_mask)
84             self.register_buffer('M_mid_c', init_mask_c)
85
86             # Encodeur
87             self.encoder = nn.Sequential(
88                 nn.Conv2d(1, 64, 3, stride=2, padding=1),
89                 nn.ReLU(),
90                 nn.Conv2d(64, 128, 3, stride=2, padding=1),
91                 nn.ReLU(),
92                 nn.Conv2d(128, 256, 3, stride=2, padding=1),
93                 nn.ReLU()
94             )
95
96             # D codeurs avec initialisation guid e
97             self.decoder_mid = self._build_decoder(self.M_mid)
98             self.decoder_mid_c = self._build_decoder(self.M_mid_c)

```



```

99
100     def _build_decoder(self, target_mask):
101         decoder = nn.Sequential(
102             nn.Conv2d(256, 128, 3, padding=1),
103             nn.PixelShuffle(2),
104             nn.Conv2d(32, 1, 1)
105         )
106         # Initialisation proche du masque cible
107         with torch.no_grad():
108             decoder[-1].weight.data = target_mask.mean() * torch.
109             ones_like(decoder[-1].weight)
110             decoder[-1].bias.data.zero_()
111         return decoder
112
113     def forward(self, x_fft):
114         encoded = self.encoder(x_fft)
115         m_mid = torch.sigmoid(self.decoder_mid(encoded))
116         m_mid_c = torch.sigmoid(self.decoder_mid_c(encoded))
117
118         # Redimensionnement si n cessaire
119         if m_mid.shape[-1] != x_fft.shape[-1]:
120             m_mid = F.interpolate(m_mid, size=x_fft.shape[-2:], mode='
121             bilinear', align_corners=False)
122             m_mid_c = F.interpolate(m_mid_c, size=x_fft.shape[-2:], mode
123             ='bilinear', align_corners=False)
124
125         return m_mid, m_mid_c
126
127 class FIRE(nn.Module):
128     def __init__(self, img_size=256):
129         super().__init__()
130         # Initialisation des masques fixes
131         M_mid, M_mid_c = create_fixed_mask(img_size)
132         self.fmre = FMRE(M_mid, M_mid_c)
133
134         # VAE de Stable Diffusion
135         self.vae = AutoencoderKL.from_pretrained(
136             "stabilityai/sd-vae-ft-mse",
137             torch_dtype=torch.float16,
138             use_safetensors=True
139         ).to(device)
140         self.vae.enable_tiling(False)
141         self.vae.requires_grad_(False)
142
143         # Classifieur l ger
144         self.classifier = nn.Sequential(
145             nn.Conv2d(6, 64, 3, stride=2),
146             nn.ReLU(),
147             nn.AdaptiveAvgPool2d(1),
148             nn.Flatten(),
149             nn.Linear(64, 1)
150         )
151
152     def apply_frequency_mask(self, x, mask):
153         # Adapter le masque la taille de l'image
154         if mask.shape[-2:] != x.shape[-2:]:

```

```

153         mask = F.interpolate(mask, size=x.shape[-2:], mode='bilinear
154     ', align_corners=False)
155
156     # Adapter les canaux si n cessaire
157     if mask.size(1) == 1 and x.size(1) == 3:
158         mask = mask.repeat(1, 3, 1, 1)
159
160     # FFT -> Masquage -> IFFT
161     freq = torch.fft.fftshift(torch.fft.fft2(x, dim=(-2, -1)), dim
162 =(-2, -1))
163     masked_freq = freq * mask
164     return torch.fft.ifft2(torch.fft.ifftshift(masked_freq, dim=(-2,
165 -1)), dim=(-2, -1)).real
166
167 def forward(self, x):
168     gray_x = torch.mean(x, dim=1, keepdim=True)
169     x_fft = torch.fft.fftshift(torch.fft.fft2(gray_x, dim=(-2, -1)),
170 dim=(-2, -1)).abs().log()
171
172     # G n ration des masques via FMRE
173     m_mid, m_mid_c = self.fmre(x_fft)
174
175     # Cr ation de l'image pseudo-g n r e
176     x_pseudo = self.apply_frequency_mask(x, m_mid_c)
177
178     # Reconstruction via VAE
179     with torch.amp.autocast(device_type='cuda', dtype=torch.float16)
180 :
181         latent_x = self.vae.encode(x).latent_dist.sample()
182         latent_pseudo = self.vae.encode(x_pseudo).latent_dist.sample
183 ()
184         recon_x = self.vae.decode(latent_x).sample.float()
185         recon_pseudo = self.vae.decode(latent_pseudo).sample.float()
186
187     # Calcul des erreurs de reconstruction
188     delta_x = (recon_x - x).abs()
189     delta_pseudo = (recon_pseudo - x_pseudo).abs()
190
191     return self.classifier(torch.cat([delta_x, delta_pseudo], dim=1)
192 )
193
194 def compute_loss(self, x, y):
195     batch_size = x.size(0)
196
197     # Pr paration pour FMRE
198     gray_x = torch.mean(x, dim=1, keepdim=True)
199     x_fft = torch.fft.fftshift(torch.fft.fft2(gray_x, dim=(-2, -1)),
200 dim=(-2, -1)).abs().log()
201
202     # Obtention des masques
203     m_mid, m_mid_c = self.fmre(x_fft)
204
205     # Cr ation des images filtr es
206     x_mid = self.apply_frequency_mask(x, m_mid)
207     x_pseudo = self.apply_frequency_mask(x, m_mid_c)
208
209     # Reconstruction

```

```

202     with torch.amp.autocast(device_type='cuda', dtype=torch.float16)
203     :
204         latent_x = self.vae.encode(x).latent_dist.sample()
205         latent_pseudo = self.vae.encode(x_pseudo).latent_dist.sample
206     ()
207         recon_x = self.vae.decode(latent_x).sample.float()
208         recon_pseudo = self.vae.decode(latent_pseudo).sample.float()
209
210     # Erreurs de reconstruction
211     delta_x = (recon_x - x).abs()
212     delta_pseudo = (recon_pseudo - x_pseudo).abs()
213
214     # 1. L_mid_rec: alignement mid-freq avec erreur
215     L_mid_rec = F.mse_loss(x_mid, delta_x)
216
217     # 2. L_mask: guidage vers masques pr d finis (batch-aware)
218     M_mid = self.fmre.M_mid.expand(batch_size, -1, -1, -1)
219     M_mid_c = self.fmre.M_mid_c.expand(batch_size, -1, -1, -1)
220
221     L_mask = (
222         F.mse_loss(m_mid, M_mid) +
223         F.mse_loss(m_mid_c, M_mid_c) +
224         F.mse_loss(1.0 - m_mid - m_mid_c, torch.zeros_like(m_mid))
225     )
226
227     # 3. L_ce: perte de classification
228     output = self.classifier(torch.cat([delta_x, delta_pseudo], dim
229 =1)).squeeze()
230     L_ce = F.binary_cross_entropy_with_logits(output, y)
231
232     # Coefficients du papier (section 3.4)
233     total_loss = 0.2 * L_mid_rec + 0.2 * L_mask + 0.6 * L_ce
234
235     # Gestion des NaN
236     if torch.isnan(total_loss):
237         print("NaN detected! Applying corrective measures...")
238         total_loss = torch.tensor(0.0, device=device, requires_grad=
239 True)
240
241     return total_loss, {
242         "L_mid_rec": L_mid_rec.item(),
243         "L_mask": L_mask.item(),
244         "L_ce": L_ce.item()
245     }
246
247 def train_kaggle():
248     train_dataset = FireDataset(train_data['file_name'].values,
249 train_data['label'].values)
250     val_dataset = FireDataset(val_data['file_name'].values, val_data['
251 label'].values, train=False)
252     train_loader = DataLoader(train_dataset, batch_size=4, shuffle=True)
253     val_loader = DataLoader(val_dataset, batch_size=4)
254
255     model = FIRE().to(device)
256     optimizer = torch.optim.AdamW(model.parameters(), lr=1e-5) #
257     Learning rate r duit

```

```

252     scaler = torch.amp.GradScaler(device_type='cuda')
253
254     train_losses = []
255     val_aucs = []
256     component_losses = {"L_mid_rec": [], "L_mask": [], "L_ce": []}
257
258     for epoch in range(15):
259         model.train()
260         running_loss = 0.0
261         epoch_components = {"L_mid_rec": 0.0, "L_mask": 0.0, "L_ce":
0.0}
262
263         for batch_idx, (imgs, labels) in enumerate(tqdm(train_loader)):
264             # Verification NaN
265             if torch.isnan(imgs).any():
266                 print(f"NaN in batch {batch_idx}, skipping...")
267                 continue
268
269             imgs, labels = imgs.to(device), labels.to(device)
270             optimizer.zero_grad()
271
272             with torch.amp.autocast(device_type='cuda', dtype=torch.
float16):
273                 loss, components = model.compute_loss(imgs, labels)
274
275                 if not torch.isnan(loss):
276                     scaler.scale(loss).backward()
277
278                     # Clip gradient pour stabilit 
279                     torch.nn.utils.clip_grad_norm_(model.parameters(),
max_norm=1.0)
280
281                     if (batch_idx + 1) % 4 == 0:
282                         scaler.step(optimizer)
283                         scaler.update()
284                         optimizer.zero_grad()
285
286                     running_loss += loss.item()
287                     for k in components:
288                         epoch_components[k] += components[k]
289                 else:
290                     print(f"Skipping batch {batch_idx} due to NaN loss")
291
292             # Calcul des moyennes
293             avg_train_loss = running_loss / len(train_loader)
294             train_losses.append(avg_train_loss)
295
296             for k in epoch_components:
297                 component_losses[k].append(epoch_components[k] / len(
train_loader))
298             print(f"Epoch {epoch+1} | {k}: {component_losses[k][-1]:.4f}
")
299
300             # Validation
301             model.eval()
302             val_preds, val_labels = [], []
303             with torch.no_grad():

```

```

304         for imgs, labels in val_loader:
305             outputs = model(imgs.to(device)).squeeze()
306             outputs = torch.nan_to_num(outputs, nan=0.5) # Gestion
NaN
307             val_preds.extend(torch.sigmoid(outputs).cpu().numpy())
308             val_labels.extend(labels.numpy())
309
310         try:
311             auc = roc_auc_score(val_labels, val_preds)
312         except:
313             auc = 0.5
314         val_aucs.append(auc)
315         print(f"Epoch {epoch+1} | Train Loss: {avg_train_loss:.4f} | Val
AUC: {auc:.4f}")
316
317     # Visualisation
318     plt.figure(figsize=(12,5))
319     plt.subplot(1,2,1)
320     plt.plot(train_losses, label='Train Loss')
321     plt.xlabel('Epoch')
322     plt.ylabel('Loss')
323     plt.legend()
324
325     plt.subplot(1,2,2)
326     plt.plot(val_aucs, label='Val AUC')
327     plt.xlabel('Epoch')
328     plt.ylabel('AUC')
329     plt.legend()
330     plt.tight_layout()
331     plt.show()
332
333     return model
334
335
336 def generate_submission(model):
337     test_dataset = FireDataset(df_test['id'].values, [0]*len(df_test),
train=False)
338     test_loader = DataLoader(test_dataset, batch_size=4)
339     model.eval()
340     predictions = []
341
342     with torch.no_grad():
343         for imgs, _ in tqdm(test_loader):
344             outputs = model(imgs.to(device)).squeeze()
345             outputs = torch.nan_to_num(outputs, nan=0.5) # Remplacer
NaN
346             predictions.extend(torch.sigmoid(outputs).cpu().numpy())
347
348     binary_labels = (np.array(predictions) > 0.5).astype(int)
349     submission = pd.DataFrame({
350         'id': ['test_data_v2/' + os.path.basename(p) for p in df_test['
id'].values],
351         'label': binary_labels
352     })
353     submission.to_csv('submission.csv', index=False)
354     return submission
355

```

```
356
357 if __name__ == "__main__":
358     trained_model = train_kaggle()
359     submission_df = generate_submission(trained_model)
360     print(submission_df.head())
361     torch.save(trained_model.state_dict(), 'model.pth')
```

11.2 Code resnet :

```
1
2 import pandas as pd
3 import numpy as np
4 import os
5 from PIL import Image
6 import matplotlib.pyplot as plt
7 from tqdm import tqdm
8 import torch
9 import torch.nn as nn
10 import torch.nn.functional as F
11 from torchvision.models import regnet_y_32gf, RegNet_Y_32GF_Weights,
    resnet50, ResNet50_Weights
12 import torchvision.transforms as transforms
13
14 def load_data():
15     train_path="/kaggle/input/detect-ai-vs-human-generated-images/train.
    csv"
16     test_path="/kaggle/input/detect-ai-vs-human-generated-images/test.
    csv"
17
18     dataset_path="/kaggle/input/ai-vs-human-generated-dataset"
19
20     # on cherche      lire les fichier csv, donc on cherche      appliquer
    pd.read_csv
21
22     train_df=pd.read_csv(train_path)
23     test_df=pd.read_csv(test_path)
24
25     #img=train_df.iloc[0,1]
26     #image_path=os.path.join(dataset_path, img)
27
28
29     #pic=Image.open(image_path)
30     #plt.imshow(pic)
31     #plt.axis("off")
32     #plt.show()
33     return train_df, test_df, dataset_path
34
35
36
37
38 train_df, test_df, dataset_path=load_data()
39 print(test_df.shape)
40 img=test_df.iloc[1,0]
41 img=os.path.join(dataset_path,img)
42
43 pic=Image.open(img)
44 plt.imshow(pic)
45 plt.axis("off")
46 plt.show()
47
48 display(train_df.head())
49 print("-----")
50 print(test_df.head())
51 print("-----")
```

```

52
53
54
55 submission_path="/kaggle/input/solution/submission.csv"
56 solution_df=pd.read_csv(submission_path)
57
58 img=solution_df.iloc[1,0]
59 img=os.path.join(dataset_path,img)
60
61 pic=Image.open(img)
62 plt.imshow(pic)
63 plt.axis("off")
64 plt.show()
65 display(solution_df.head)
66
67 d finir le filtre passe bande
68
69 model=resnet50(weights=ResNet50_Weights.IMAGENET1K_V2)
70 for param in model.parameters():
71     param.requires_grad=False
72 num=model.fc.in_features
73 model.fc=nn.Sequential(
74     nn.Linear(num,512),
75     nn.ReLU(),
76     nn.Dropout(0.4),
77     nn.Linear(512,1)
78 )
79 for param in model.layer4.parameters():
80     param.required_grad=True
81 for param in model.fc.parameters():
82     param.requires_grad=True
83
84
85 device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
86 model = model.to(device)
87
88
89
90 from torchvision.transforms import InterpolationMode
91 import albumentations.core.composition
92 import albumentations.augmentations as A
93 from albumentations.pytorch import ToTensorV2
94 import cv2
95
96 train_transform=transforms.Compose([
97     transforms.Resize(224, interpolation=InterpolationMode.BICUBIC),
98     transforms.RandomResizedCrop(224),
99     transforms.RandomHorizontalFlip(),
100    transforms.RandomVerticalFlip(),
101    transforms.RandomRotation(20),
102    #transforms.GaussianBlur(kernel_size=(7, 13), sigma=(0.1, 1.0)),
103    transforms.ColorJitter(brightness=0.2, contrast=0.2, saturation=0.2,
104    hue=0.1),
105    transforms.ToTensor(),
106    transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224,
107    0.225])

```



```

107 ])
108
109
110
111 train_transform2 = transforms.Compose([
112     transforms.Resize(256), # Augmentation l g re avant crop
113     transforms.RandomHorizontalFlip(p=0.5),
114     transforms.RandomResizedCrop(224, scale=(0.08, 1.0), ratio=(0.75,
115     1.33), antialias=True),
116     transforms.ColorJitter(brightness=0.04, contrast=0.04, saturation
117     =0.04, hue=0.1),
118     transforms.RandomGrayscale(p=0.2),
119     transforms.RandomErasing(scale=(96/224, 96/224), ratio=(1, 1), value
120     =128, p=0.2),
121     transforms.GaussianBlur(kernel_size=(5, 5), sigma=(0.1, 2.0)), #
122     Ajust pour rester raisonnable
123     transforms.ToTensor(),
124     transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224,
125     0.225])
126 ])
127
128 #strong_transform = albumentations.core.composition.Compose([
129 #    A.geometric.resize.SmallestMaxSize(max_size=512),
130 #    A.geometric.transforms.HorizontalFlip(p=0.5),
131 #    A.crops.transforms.RandomResizedCrop(height=512, width=512, scale
132 #    =(0.08, 1.0), ratio=(0.75, 1.0/0.75), p=0.2),
133 #    A.crops.RandomCrop(height=512, width=512),
134 #    A.transforms.ColorJitter(brightness=0.04, contrast=0.04, saturation
135 #    =0.04, hue=0.1, p=0.8),
136 #    A.transforms.ToGray(p=0.2),
137 #    A.dropout.CoarseDropout(max_holes=1, min_holes=1, hole_height_range
138 #    =(96, 96), hole_width_range=(96, 96), fill_value=128, p=0.2),
139 #    A.transforms.GaussNoise(var_limit=(10.0, 50.0), p=0.2),
140 #    A.GaussianBlur(),
141 #
142 #    A.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225]),
143 #    ToTensorV2()
144 #])
145
146
147 validation_transform=transforms.Compose([
148     transforms.Resize(224,interpolation=InterpolationMode.BICUBIC),
149     transforms.CenterCrop(224),
150     transforms.ToTensor(),
151     transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224,
152     0.225])
153 ])
154
155 test_transform=transforms.Compose([
156     transforms.Resize(224,interpolation=InterpolationMode.BICUBIC),
157     transforms.CenterCrop(224),
158     transforms.ToTensor(),
159     transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224,
160     0.225])
161 ])

```

```

154 ]
155 ])
156 #print(img)
157 #image=cv2.imread(img)
158 #image = cv2.cvtColor(image, cv2.COLOR_BGR2RGB)
159 #augmented = strong_transform(image=image)
160 #transformed_image = augmented["image"]
161 #transformed_image = transformed_image.permute(1, 2, 0).numpy()
162
163 # Afficher l'image originale et transformée
164 #fig, axes = plt.subplots(1, 2, figsize=(10, 5))
165 #axes[0].imshow(image)
166 #axes[0].set_title("Image originale")
167 #axes[1].imshow(transformed_image)
168 #axes[1].set_title("Image transformée")
169 #plt.show()
170
171 from torch.utils.data import Dataset, DataLoader
172 class CustomDataset(Dataset):
173     def __init__(self, csv_file, transform=None):
174         self.data = pd.read_csv(csv_file) # Charger le CSV
175         self.transform = transform
176
177     def __len__(self):
178         return len(self.data)
179
180     def __getitem__(self, idx):
181         img_path = self.data.iloc[idx, 1] # Supposons que la première
182         # colonne est le chemin de l'image
183         dataset_path="/kaggle/input/ai-vs-human-generated-dataset"
184         img_path=os.path.join(dataset_path,img_path)
185         label = self.data.iloc[idx, 2] # Supposons que la deuxième
186         # colonne est le label (0 ou 1)
187
188         # Charger l'image
189         image = Image.open(img_path).convert("RGB")
190
191         # Appliquer les transformations
192         if self.transform:
193             image = self.transform(image)
194
195         return image, label
196 class CustomDataset2(Dataset):
197     def __init__(self, file_list, labels=None, transform=None):
198         self.file_list = file_list
199         self.labels = labels
200         self.transform = transform
201
202     def __len__(self):
203         return len(self.file_list)
204
205     def __getitem__(self, idx):
206         img_path = self.file_list[idx]
207         img = Image.open(img_path).convert("RGB")
208         if self.transform:
209             img = self.transform(img)

```

```

209         if self.labels is not None:
210             label = self.labels[idx]
211             return img, label
212         else:
213             # Pas de label en inference
214             return img
215
216 class CustomDataset3(Dataset):
217     def __init__(self, file_list, labels=None, transform=None,
218                 strong_transform=None):
219         self.file_list = file_list
220         self.labels = labels
221         self.transform = transform
222         self.strong_transform=strong_transform
223
224     def __len__(self):
225         return len(self.file_list)
226
227     def __getitem__(self, idx):
228         img_path = self.file_list[idx]
229         image = cv2.imread(img_path)
230         image=cv2.cvtColor(image,cv2.COLOR_BGR2RGB)
231         if self.transform:
232             image = self.transform(img)
233         if strong_transform:
234             augmented=self.strong_transform(image=image)
235             image=augmented["image"]
236         if self.labels is not None:
237             label = self.labels[idx]
238             return image, label
239         else:
240             # Pas de label en inference
241             return image
242
243 # Charge le dataset
244 train_dataset = CustomDataset(csv_file="/kaggle/input/detect-ai-vs-human
245                               -generated-images/train.csv", transform=train_transform)
246 # Cr er un DataLoader pour l'entra nement
247 print(len(train_dataset))
248 # V rifier si a fonctionne
249
250 #from torch.utils.data import random_split
251 #from torch.utils.data import random_split, DataLoader
252
253 # D finir la taille des ensembles
254 #train_size = int(0.9 * len(train_dataset)) # 90% pour l'entra nement
255 #val_size = len(train_dataset) - train_size # 10% pour la validation
256
257 # D couper train_dataset en train et validation
258 #train_dataset, val_dataset = random_split(train_dataset, [train_size,
259                                             val_size])
260 #train_loader = DataLoader(train_dataset, batch_size=64, shuffle=True)
261 #val_loader=DataLoader(val_dataset,batch_size=64,shuffle=False)
262 #for image, label in train_loader:

```

```

263 # print(image.shape) # Devrait tre torch.Size([32, 3, 224, 224])
264 # print(label.shape) # Devrait tre torch.Size([32])
265 # break
266
267 #print(len(train_dataset))
268 #print(len(val_dataset))
269
270
271 # Paths to the dataset
272 from sklearn.model_selection import train_test_split
273 base_dir = '/kaggle/input/ai-vs-human-generated-dataset'
274 train_csv_path = os.path.join(base_dir, 'train.csv')
275 test_csv_path = os.path.join(base_dir, 'test.csv')
276
277 # Reading the training CSV file
278 df_train = pd.read_csv(train_csv_path)
279 # Example of a row: file_name="train_data/041be3153810...", label=0 or 1
280
281 # Reading the testing CSV file
282 df_test = pd.read_csv(os.path.join(base_dir, 'test.csv'))
283 # Exemple: df_test['id'] = "test_data/e25323c62af644fba97afb846261b05b.
    jpg", etc.
284
285 # Adding the full path to the file_name instead of just "
    trainORtest_data/xxx.jpg"
286 df_test['id'] = df_test['id'].apply(lambda x: os.path.join(base_dir, x))
287 df_train['file_name'] = df_train['file_name'].apply(lambda x: os.path.
    join(base_dir, x))
288
289 all_image_paths = df_train['file_name'].values
290 all_labels = df_train['label'].values
291
292 # Splitting train/validation (95% / 5%)
293 train_paths, val_paths, train_labels, val_labels = train_test_split(
294     all_image_paths,
295     all_labels,
296     test_size=0.05,
297     stratify=all_labels,
298     random_state=42
299 )
300
301
302 #train_paths = train_paths[:1500]
303 #val_paths = val_paths[:150]
304 print(f"Train Data: {len(train_paths)}")
305 print(f"Validation Data: {len(val_paths)}")
306 print(train_paths)
307
308 train_data = CustomDataset2(train_paths, train_labels, transform=
    train_transform)
309 val_data = CustomDataset2(val_paths, val_labels, transform=
    validation_transform)
310
311 train_loader = DataLoader(dataset=train_data, batch_size=64, shuffle=
    True, num_workers=4)
312 val_loader = DataLoader(dataset=val_data, batch_size=64, shuffle=
    False, num_workers=4)

```

```

313
314 print(f"Train Dataset size: {len(train_data)}")
315 print(f"Validation Dataset size: {len(val_data)}")
316
317 for image, label in train_loader:
318     print(image.shape) # Devrait tre torch.Size([32, 3, 224, 224])
319     print(label.shape) # Devrait tre torch.Size([32])
320     break
321
322 import torch.optim as optim
323 from sklearn.metrics import roc_curve, auc
324 from sklearn.metrics import f1_score, roc_auc_score
325
326 # Initialisation
327 training_param = [p for p in model.parameters() if p.requires_grad]
328 print(len(training_param))
329 optimizer = torch.optim.AdamW(training_param, lr=1e-4, weight_decay=1e-2)
330 scheduler = torch.optim.lr_scheduler.CosineAnnealingLR(optimizer, T_max
    =10)
331 criterion = nn.BCEWithLogitsLoss()
332 torch.cuda.empty_cache()
333
334
335 train_losses, train_accuracies = [], []
336 val_losses, val_accuracies, val_f1_scores, val_roc_aucs = [], [], [], []
337
338 epochs = 2
339 best_val_loss = float('inf')
340 early_stopping_patience = 5
341 early_stopping_counter = 0
342
343 for epoch in range(epochs):
344     epoch_loss = 0.0
345     epoch_accuracy = 0.0
346     model.train()
347
348     for data, labels in tqdm(train_loader, desc=f"Training epoch {epoch
+1}"):
349         data = data.to(device)
350         labels = labels.to(device).float()
351
352         optimizer.zero_grad()
353
354         outputs = model(data).squeeze(1)
355         loss = criterion(outputs, labels)
356
357         loss.backward()
358         optimizer.step()
359         preds = (torch.sigmoid(outputs) > 0.5).float()
360         epoch_accuracy += (preds == labels).float().mean().item()
361         epoch_loss += loss.item()
362
363     train_losses.append(epoch_loss/len(train_loader))
364     train_accuracies.append(epoch_accuracy/len(train_loader))
365
366     # Validation
367     model.eval()

```

```

368     epoch_val_loss = 0.0
369     epoch_val_accuracy = 0.0
370     val_pred_class = []
371     val_label_class = []
372
373     with torch.no_grad():
374         for data, labels in tqdm(val_loader, desc=f"validation_epoch {
epoch+1}"):
375             data, labels = data.to(device), labels.to(device).float()
376
377
378             outputs = model(data).squeeze(1)
379             loss = criterion(outputs, labels)
380
381             preds = (torch.sigmoid(outputs) > 0.5).float()
382             epoch_val_accuracy += (preds == labels).float().mean().item
()
383             epoch_val_loss += loss.item()
384             val_pred_class.extend(preds.cpu().numpy())
385             val_label_class.extend(labels.cpu().numpy())
386
387     val_loss = epoch_val_loss/len(val_loader)
388     val_acc = epoch_val_accuracy/len(val_loader)
389     val_f1 = f1_score(np.array(val_label_class), np.array(val_pred_class
))
390     val_roc_auc = roc_auc_score(np.array(val_label_class), np.array(
val_pred_class))
391
392     val_losses.append(val_loss)
393     val_accuracies.append(val_acc)
394     val_f1_scores.append(val_f1)
395     val_roc_aucs.append(val_roc_auc)
396
397     print(f"Epoch [{epoch+1}/{epochs}] "
398           f"Train Loss: {epoch_loss/len(train_loader):.4f} | "
399           f"Train Acc: {epoch_accuracy/len(train_loader):.4f} | "
400           f"Val Loss: {val_loss:.4f} | "
401           f"Val Acc: {val_acc:.4f} | "
402           f"Val F1: {val_f1:.4f} | "
403           f"Val ROC AUC: {val_roc_auc:.4f}")
404
405     scheduler.step()
406
407     # Early stopping
408     if epoch_val_loss < best_val_loss:
409         best_val_loss = epoch_val_loss
410         torch.save(model.state_dict(), 'best_model.pth')
411         early_stopping_counter = 0
412     else:
413         early_stopping_counter += 1
414         if early_stopping_counter >= early_stopping_patience:
415             print("Early stopping triggered")
416             break
417
418     torch.save(model.state_dict(), "model_weights.pth")
419
420     #torch.save(model, "model.pth")

```

```

421
422 fig, axs = plt.subplots(1, 2, figsize=(14, 6))
423
424 # Plot Loss
425 axs[0].plot(range(1, len(train_losses) + 1), train_losses, label='Train
    Loss', color='blue')
426 axs[0].plot(range(1, len(val_losses) + 1), val_losses, label='Validation
    Loss', color='red')
427 axs[0].set_title('Loss per Epoch')
428 axs[0].set_xlabel('Epochs')
429 axs[0].set_ylabel('Loss')
430 axs[0].legend()
431
432 # Plot Accuracy
433 axs[1].plot(range(1, len(train_accuracies) + 1), train_accuracies, label
    ='Train Accuracy', color='blue')
434 axs[1].plot(range(1, len(val_accuracies) + 1), val_accuracies, label='
    Validation Accuracy', color='red')
435 axs[1].set_title('Accuracy per Epoch')
436
437 axs[1].set_xlabel('Epochs')
438 axs[1].set_ylabel('Accuracy')
439 axs[1].legend()
440
441 # Show the plots
442 plt.tight_layout()
443 plt.savefig('Accuracy per Epoch.png')
444 plt.show()
445
446 from sklearn.metrics import roc_curve
447
448 # Compute ROC curve
449 fpr, tpr, _ = roc_curve(np.array(val_label_class , dtype=int), np.array(
    val_pred_class, dtype=int))
450
451 # Plot ROC Curve
452 plt.figure(figsize=(8, 6))
453 plt.plot(fpr, tpr, color='blue', label=f'ROC Curve (AUC = {val_roc_auc
    :.4f})')
454 plt.plot([0, 1], [0, 1], linestyle='--', color='gray', label='Random
    Guess')
455 plt.xlabel("False Positive Rate")
456 plt.ylabel("True Positive Rate")
457 plt.title("ROC Curve")
458 plt.legend()
459 plt.grid()
460 plt.show()
461
462 from sklearn.metrics import confusion_matrix, classification_report
463 import seaborn as sns
464 # Generate and plot Confusion Matrix
465 conf_matrix = confusion_matrix(val_label_class , val_pred_class)
466
467 # Plot confusion matrix
468 plt.figure(figsize=(8, 6))
469 sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues', xticklabels
    =['Class 0', 'Class 1'],

```

```

470         yticklabels=['Class 0', 'Class 1'])
471 plt.title('Confusion Matrix')
472 plt.xlabel('Predicted')
473 plt.ylabel('True')
474 plt.show()
475
476 # Print Classification Report
477 print("\nClassification Report:")
478 print(classification_report(val_label_class , val_pred_class))
479
480 base_dir = '/kaggle/input/ai-vs-human-generated-dataset'
481 test_csv_path = os.path.join(base_dir, 'test.csv')
482 df_test = pd.read_csv(os.path.join(base_dir, 'test.csv'))
483 df_test['id'] = df_test['id'].apply(lambda x: os.path.join(base_dir, x))
484 class TestImageDataset(Dataset):
485     def __init__(self, file_list, transform=None):
486         self.file_list = file_list
487         self.transform = transform
488
489     def __len__(self):
490         return len(self.file_list)
491
492     def __getitem__(self, idx):
493         img_path = self.file_list[idx]
494         img = Image.open(img_path).convert("RGB")
495         if self.transform:
496             img = self.transform(img)
497         return img, os.path.basename(img_path)
498
499 test_dataset = TestImageDataset(df_test['id'].values, transform=
    test_transform)
500 test_loader = DataLoader(test_dataset, batch_size=64, shuffle=False,
    num_workers=4)
501
502
503
504 model.eval()
505 predictions = []
506 image_names = []
507
508 with torch.no_grad():
509     for data, names in tqdm(test_loader, desc="Predicting"):
510         data = data.to(device)
511         output = model(data).squeeze(1) # Ensure correct dimensions
512
513         # Apply sigmoid activation to convert logits to probabilities
514         probs = torch.sigmoid(output)
515
516         # Convert probabilities to binary class (0 or 1) using threshold
517         0.5
518         preds = (probs > 0.5).int()
519
520         predictions.extend(preds.cpu().numpy())
521         image_names.extend([f"test_data_v2/{name}" for name in names])
522
523 # Cr er le DataFrame au format "id,label"
524 submission_df = pd.DataFrame({

```



```
524     'id': image_names,
525     'label': predictions
526 })
527
528 submission_df.head()
529
530 submission_df.to_csv("submission.csv", index=False)
531 print("Submission file generated: submission.csv")
532
533 # Sauvegarder les poids
534 torch.save(model.state_dict(), 'model_weights.pth')
535
536 # D placer le fichier vers le dossier d'output
537 import shutil
538 shutil.move('model_weights.pth', '/kaggle/working/model_weights.pth')
```

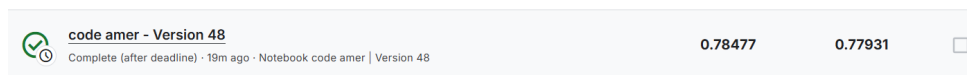


FIGURE 30

11.3 Convnext+swin

```
1 import pandas as pd
2 import numpy as np
3 import os
4 from PIL import Image
5 import matplotlib.pyplot as plt
6 from tqdm import tqdm
7 import torch
8 import torch.nn as nn
9 import torch.nn.functional as F
10 from torchvision.models import regnet_y_32gf, RegNet_Y_32GF_Weights,
    resnet50, ResNet50_Weights
11 import torchvision.transforms as transforms
12 from timm import create_model
13
14 # Data avec augmentation de Dataset
15
16 train_path="/kaggle/input/detect-ai-vs-human-generated-images/train.csv"
17 test_path="/kaggle/input/detect-ai-vs-human-generated-images/test.csv"
18
19 dataset_path="/kaggle/input/ai-vs-human-generated-dataset"
20
21 train_df=pd.read_csv(train_path)
22 test_df=pd.read_csv(test_path)
23
24 image_path=os.path.join(dataset_path,train_df.iloc[0,1])
25 image=Image.open(image_path)
26 plt.imshow(image)
27
28 train_transforms=transforms.Compose([
29     transforms.Resize((232)),
30     transforms.RandomResizedCrop(224),
31     transforms.RandomHorizontalFlip(),
32     transforms.RandomRotation(10),
33     transforms.ColorJitter(brightness=0.2, contrast=0.2, saturation=0.2,
    hue=0.1),
34     transforms.ToTensor(),
35     transforms.Normalize([0.485, 0.456, 0.406], [0.229, 0.224, 0.225])
36 ])
37
38 image2=train_transforms(image)
39 plt.imshow(image2.permute(1,2,0))
40
41 test_transforms=transforms.Compose([
42     transforms.Resize(232),
43     transforms.CenterCrop(224),
44     transforms.ToTensor(),
45     transforms.Normalize([0.485, 0.456, 0.406], [0.229, 0.224, 0.225])
46 ])
47 validation_transform=transforms.Compose([
48     transforms.Resize(224,interpolation=transforms.InterpolationMode.
    BICUBIC),
49     transforms.CenterCrop(224),
50     transforms.ToTensor(),
51     transforms.Normalize(mean=[0.485, 0.456, 0.406], std=[0.229, 0.224,
    0.225])
```

```

52
53 ])
54
55 import gc
56 import torch
57
58 def clear_gpu_memory():
59     # tape 1 : Supprimer toutes les r f r ences
60     gc.collect()
61
62     # tape 2 : Vider le cache PyTorch
63     torch.cuda.empty_cache()
64
65     # tape 3 : Reset CUDA (nettoyage profond)
66     if torch.cuda.is_available():
67         torch.cuda.reset_peak_memory_stats()
68         torch.cuda.reset_max_memory_allocated()
69         torch.cuda.reset_max_memory_cached()
70
71 # Model
72
73 class convnextandswin(nn.Module):
74     def __init__(self, num_classes=1):
75         super(convnextandswin, self).__init__()
76
77         # Load ConvNeXt Large
78         self.convnext = create_model("convnext_large", pretrained=True,
num_classes=0)
79         convnext_out = self.convnext.num_features
80
81         # Load Swin Transformer
82         self.swin = create_model("swin_base_patch4_window7_224",
pretrained=True, num_classes=0)
83         swin_out = self.swin.num_features
84
85         # Global Average Pooling for each model
86         self.global_avg_pooling_convnext = nn.AdaptiveAvgPool1d(1)
87         self.global_avg_pooling_swin = nn.AdaptiveAvgPool1d(1)
88
89         # Fully Connected Layers for feature fusion
90         self.feature_fusion = nn.Sequential(
91             nn.BatchNorm1d(convnext_out + swin_out),
92             nn.Linear(convnext_out + swin_out, 1024),
93             nn.ReLU(),
94             nn.Dropout(0.4),
95             nn.Linear(1024, 512),
96             nn.ReLU(),
97             nn.Dropout(0.4),
98             nn.Linear(512, 256),
99             nn.ReLU(),
100            nn.Dropout(0.4),
101        )
102
103         # Decoder: Additional layers to output classification results
104         self.decoder = nn.Sequential(
105             nn.Linear(256, 128),
106             nn.ReLU(),

```

```

107         nn.Linear(128, num_classes)
108     )
109
110     def forward(self, x):
111         # Pass through ConvNeXt and Swin Transformer
112         x_convnext = self.convnext(x)
113         x_swin = self.swin(x)
114
115         # Debug print statements to inspect tensor shapes
116         #print("Shape of x_convnext before pooling:", x_convnext.shape)
117         #print("Shape of x_swin before pooling:", x_swin.shape)
118
119         # Apply global average pooling
120         x_convnext = self.global_avg_pooling_convnext(x_convnext.
unsqueeze(2)).view(x_convnext.size(0), -1)
121         x_swin = self.global_avg_pooling_swin(x_swin.unsqueeze(2)).view(
x_swin.size(0), -1)
122
123         # Debug print statements to inspect tensor shapes after pooling
124         #print("Shape of x_convnext after pooling:", x_convnext.shape)
125         #print("Shape of x_swin after pooling:", x_swin.shape)
126
127         # Concatenate both feature vectors
128         x_combined = torch.cat((x_convnext, x_swin), dim=1)
129         x_fused = self.feature_fusion(x_combined)
130
131         # Pass through the decoder to output the final classification
result
132         decoded_output = self.decoder(x_fused)
133
134         return decoded_output
135
136     # Initialize the model with ConvNeXt Large and Swin Transformer
137     device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
138     model = convnextandswin(num_classes=1).to(device)
139
140     # Freeze All Layers
141     for param in model.convnext.parameters():
142         param.requires_grad = False
143
144     for param in model.swin.parameters():
145         param.requires_grad = False
146
147     # Unfreeze Last 10 Layers
148     for param in list(model.convnext.parameters())[-20:]:
149         param.requires_grad = True
150
151     for param in list(model.swin.parameters())[-20:]:
152         param.requires_grad=True
153     model_path="/kaggle/input/fusion-model/best_model.pth"
154     model.load_state_dict(torch.load(model_path, map_location=device,
weights_only=True))
155
156
157
158     # Data
159

```

```

160 from torch.utils.data import Dataset, DataLoader
161 class CustomDataset2(Dataset):
162     def __init__(self, file_list, labels=None, transform=None):
163         self.file_list = file_list
164         self.labels = labels
165         self.transform = transform
166
167     def __len__(self):
168         return len(self.file_list)
169
170     def __getitem__(self, idx):
171         img_path = self.file_list[idx]
172         img = Image.open(img_path).convert("RGB")
173         if self.transform:
174             img = self.transform(img)
175         if self.labels is not None:
176             label = self.labels[idx]
177             return img, label
178         else:
179             # Pas de label en inference
180             return img
181
182 # Paths to the dataset
183 from sklearn.model_selection import train_test_split
184 base_dir = '/kaggle/input/ai-vs-human-generated-dataset'
185 train_csv_path = os.path.join(base_dir, 'train.csv')
186 test_csv_path = os.path.join(base_dir, 'test.csv')
187
188 # Reading the training CSV file
189 df_train = pd.read_csv(train_csv_path)
190 # Example of a row: file_name="train_data/041be3153810...", label=0 or 1
191
192 # Reading the testing CSV file
193 df_test = pd.read_csv(os.path.join(base_dir, 'test.csv'))
194 # Exemple: df_test['id'] = "test_data/e25323c62af644fba97afb846261b05b.
195         jpg", etc.
196
197 # Adding the full path to the file_name instead of just "
198         trainORtest_data/xxx.jpg"
199 df_test['id'] = df_test['id'].apply(lambda x: os.path.join(base_dir, x))
200 df_train['file_name'] = df_train['file_name'].apply(lambda x: os.path.
201         join(base_dir, x))
202
203 all_image_paths = df_train['file_name'].values
204 all_labels = df_train['label'].values
205
206 # Splitting train/validation (95% / 5%)
207 train_paths, val_paths, train_labels, val_labels = train_test_split(
208     all_image_paths,
209     all_labels,
210     test_size=0.05,
211     stratify=all_labels,
212     random_state=42
213 )
214
215 #train_paths = train_paths[:1500]

```

```

214 #val_paths = val_paths[:150]
215 print(f"Train Data: {len(train_paths)}")
216 print(f"Validation Data: {len(val_paths)}")
217 print(train_paths)
218
219 train_data = CustomDataset2(train_paths,train_labels, transform=
    train_transforms)
220 val_data    = CustomDataset2(val_paths,val_labels ,transform=
    validation_transform)
221
222 train_loader = DataLoader(dataset=train_data, batch_size=8, shuffle=True
    , num_workers=4,pin_memory=True)
223 val_loader   = DataLoader(dataset=val_data, batch_size=8, shuffle=
    False, num_workers=4, pin_memory=True)
224
225 print(f"Train Dataset size: {len(train_data)}")
226 print(f"Validation Dataset size: {len(val_data)}")
227
228 for image, label in train_loader:
229     print(image.shape) # Devrait tre torch.Size([32, 3, 224, 224])
230     print(label.shape) # Devrait tre torch.Size([32])
231     break
232
233 device='cuda' if torch.cuda.is_available() else 'cpu'
234 print(device)
235
236 import torch.optim as optim
237 from sklearn.metrics import roc_curve, auc
238 from sklearn.metrics import f1_score, roc_auc_score
239
240 # Initialisation
241 training_param = [p for p in model.parameters() if p.requires_grad]
242 print(len(training_param))
243 optimizer = torch.optim.AdamW(training_param, lr=1e-4,weight_decay=1e-2)
244 scheduler = torch.optim.lr_scheduler.CosineAnnealingLR(optimizer, T_max
    =10)
245 criterion = nn.BCEWithLogitsLoss()
246 torch.cuda.empty_cache()
247 scaler = torch.GradScaler(device='cuda')
248 print(scaler)
249
250 train_losses, train_accuracies = [], []
251 val_losses, val_accuracies, val_f1_scores, val_roc_aucs = [], [], [], []
252
253 epochs = 5
254 best_val_loss = float('inf')
255 early_stopping_patience = 7
256 early_stopping_counter = 0
257
258 for epoch in range(epochs):
259     epoch_loss = 0.0
260     epoch_accuracy = 0.0
261     model.train()
262
263     for data, labels in tqdm(train_loader, desc=f"Training epoch {epoch
    +1}"):
264         data = data.to(device)

```

```

265         labels = labels.to(device).float()
266
267         optimizer.zero_grad()
268         with torch.autocast(device_type='cuda'):
269             outputs = model(data).squeeze(1)
270             loss = criterion(outputs, labels)
271
272         scaler.scale(loss).backward()
273         scaler.step(optimizer)
274         scaler.update()
275
276         preds = (torch.sigmoid(outputs) > 0.5).float()
277         epoch_accuracy += (preds == labels).float().mean().item()
278         epoch_loss += loss.item()
279
280     train_losses.append(epoch_loss/len(train_loader))
281     train_accuracies.append(epoch_accuracy/len(train_loader))
282
283     # Validation
284     model.eval()
285     epoch_val_loss = 0.0
286     epoch_val_accuracy = 0.0
287     val_pred_class = []
288     val_label_class = []
289
290     with torch.no_grad():
291         for data, labels in tqdm(val_loader, desc=f"validation_epoch {
epoch+1}"):
292             data, labels = data.to(device), labels.to(device).float()
293
294             with torch.autocast(device_type='cuda'):
295                 outputs = model(data).squeeze(1)
296                 loss = criterion(outputs, labels)
297
298             preds = (torch.sigmoid(outputs) > 0.5).float()
299             epoch_val_accuracy += (preds == labels).float().mean().item
()
300
301             epoch_val_loss += loss.item()
302             val_pred_class.extend(preds.cpu().numpy())
303             val_label_class.extend(labels.cpu().numpy())
304
305     val_loss = epoch_val_loss/len(val_loader)
306     val_acc = epoch_val_accuracy/len(val_loader)
307     val_f1 = f1_score(np.array(val_label_class), np.array(val_pred_class
))
308
309     val_roc_auc = roc_auc_score(np.array(val_label_class), np.array(
val_pred_class))
310
311     val_losses.append(val_loss)
312     val_accuracies.append(val_acc)
313     val_f1_scores.append(val_f1)
314     val_roc_aucs.append(val_roc_auc)
315
316     print(f"Epoch [{epoch+1}/{epochs}] "
317           f"Train Loss: {epoch_loss/len(train_loader):.4f} | "
318           f"Train Acc: {epoch_accuracy/len(train_loader):.4f} | "
319           f"Val Loss: {val_loss:.4f} | ")

```

```

318         f"Val Acc: {val_acc:.4f} | "
319         f"Val F1: {val_f1:.4f} | "
320         f"Val ROC AUC: {val_roc_auc:.4f}")
321
322     scheduler.step()
323
324     # Early stopping
325     if epoch_val_loss < best_val_loss:
326         best_val_loss = epoch_val_loss
327         torch.save(model.state_dict(), 'best_model.pth')
328         early_stopping_counter = 0
329     else:
330         early_stopping_counter += 1
331         if early_stopping_counter >= early_stopping_patience:
332             print("Early stopping triggered")
333             break
334
335     # Plot Loss and Accuracy
336
337     fig, axs = plt.subplots(1, 2, figsize=(14, 6))
338
339     # Plot Loss
340     axs[0].plot(range(1, len(train_losses) + 1), train_losses, label='Train
341                 Loss', color='blue')
342     axs[0].plot(range(1, len(val_losses) + 1), val_losses, label='Validation
343                 Loss', color='red')
344     axs[0].set_title('Loss per Epoch')
345     axs[0].set_xlabel('Epochs')
346     axs[0].set_ylabel('Loss')
347     axs[0].legend()
348
349     # Plot Accuracy
350     axs[1].plot(range(1, len(train_accuracies) + 1), train_accuracies, label
351                 = 'Train Accuracy', color='blue')
352     axs[1].plot(range(1, len(val_accuracies) + 1), val_accuracies, label='
353                 Validation Accuracy', color='red')
354     axs[1].set_title('Accuracy per Epoch')
355
356     axs[1].set_xlabel('Epochs')
357     axs[1].set_ylabel('Accuracy')
358     axs[1].legend()
359
360     # Show the plots
361     plt.tight_layout()
362     plt.savefig('Accuracy per Epoch.png')
363     plt.show()
364
365     import matplotlib.pyplot as plt
366     from sklearn.metrics import roc_curve
367
368     # Compute ROC curve
369     fpr, tpr, _ = roc_curve(np.array(val_label_class , dtype=int), np.array(
370         val_pred_class, dtype=int))
371
372     # Plot ROC Curve
373     plt.figure(figsize=(8, 6))

```



```

369 plt.plot(fpr, tpr, color='blue', label=f'ROC Curve (AUC = {val_roc_auc
    :.4f})')
370 plt.plot([0, 1], [0, 1], linestyle='--', color='gray', label='Random
    Guess')
371 plt.xlabel("False Positive Rate")
372 plt.ylabel("True Positive Rate")
373 plt.title("ROC Curve")
374 plt.legend()
375 plt.grid()
376 plt.show()
377
378
379 from sklearn.metrics import confusion_matrix, classification_report
380 import seaborn as sns
381 # Generate and plot Confusion Matrix
382 conf_matrix = confusion_matrix(val_label_class , val_pred_class)
383
384 # Plot confusion matrix
385 plt.figure(figsize=(8, 6))
386 sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues', xticklabels
    =['Class 0', 'Class 1'],
387             yticklabels=['Class 0', 'Class 1'])
388 plt.title('Confusion Matrix')
389 plt.xlabel('Predicted')
390 plt.ylabel('True')
391 plt.show()
392
393 # Print Classification Report
394 print("\nClassification Report:")
395 print(classification_report(val_label_class , val_pred_class))
396
397
398 base_dir = '/kaggle/input/ai-vs-human-generated-dataset'
399 test_csv_path = os.path.join(base_dir, 'test.csv')
400 df_test = pd.read_csv(os.path.join(base_dir, 'test.csv'))
401 df_test['id'] = df_test['id'].apply(lambda x: os.path.join(base_dir, x))
402 class TestImageDataset(Dataset):
403     def __init__(self, file_list, transform=None):
404         self.file_list = file_list
405         self.transform = transform
406
407     def __len__(self):
408         return len(self.file_list)
409
410     def __getitem__(self, idx):
411         img_path = self.file_list[idx]
412         img = Image.open(img_path).convert("RGB")
413         if self.transform:
414             img = self.transform(img)
415         return img, os.path.basename(img_path)
416
417 test_dataset = TestImageDataset(df_test['id'].values, transform=
    test_transforms)
418 test_loader = DataLoader(test_dataset, batch_size=8, shuffle=False,
    num_workers=4)
419
420 model.eval()

```

```
421 predictions = []
422 image_names = []
423
424 with torch.no_grad():
425     for data, names in tqdm(test_loader, desc="Predicting"):
426         data = data.to(device)
427
428         # On récupère la classe prédite
429         output = model(data).squeeze(1) # Ensure correct dimensions
430
431         # Apply sigmoid activation to convert logits to probabilities
432         probs = torch.sigmoid(output)
433
434         # Convert probabilities to binary class (0 or 1) using threshold
435         0.5
436         preds = (probs > 0.5).int()
437
438         predictions.extend(preds.cpu().numpy())
439         image_names.extend([f"test_data_v2/{name}" for name in names])
440
441 # Créer le DataFrame au format "id,label"
442 submission_df = pd.DataFrame({
443     'id': image_names,
444     'label': predictions
445 })
446
447 submission_df.head()
448 submission_df.to_csv("submission.csv", index=False)
449 print("Submission file generated: submission.csv")
```