

ДОДАТОК А

Перелік джерел посилання за науковими напрямками керівника та науковців
кафедри програмної інженерії

1. Кіценко Ю. О., Смеляков К. С. Розробка моделі виявлення фейкового контенту на основі архітектури EfficientNet. *Матеріали 28-го міжнародного молодіжного форуму «Радіоелектроніка та молодь у XXI столітті»*: Зб. матеріалів форуму. Т. 6. Харків, 2024. С. 435–436.

2. Smelyakov K., Kitsenko Y., Chupryna A. Deepfake detection models based on machine learning technologies. *2024 IEEE open conference of electrical, electronic and information sciences (estream)*. 2024. P. 1–6.
URL: <https://doi.org/10.1109/eStream61684.2024.10542582>.

43. The neural network technologies effectiveness for face detection / K. Smelyakov et al. *2020 IEEE third international conference on data stream mining & processing (DSMP)*. 2020. P. 201–205.

44. The neural network models effectiveness for face detection and face recognition / K. Smelyakov et al. *2021 IEEE open conference of electrical, electronic and information sciences (estream)*. 2021. P. 1–7.

ДОДАТОК Б

Слайди презентації



МІНІСТЕРСТВО
ОСВІТИ І НАУКИ
УКРАЇНИ



ХАРКІВСЬКИЙ
НАЦІОНАЛЬНИЙ
УНІВЕРСИТЕТ
РАДІОЕЛЕКТРОНИКИ

ДОСЛІДЖЕННЯ МОДЕЛЕЙ ВІЯВЛЕННЯ ФЕЙКОВОГО КОНТЕНТУ НА ОСНОВІ ТЕХНОЛОГІЙ МАШИННОГО НАВЧАННЯ

Кіценко Юрій Олександрович, ІПЗм-22-4

Науковий керівник: д.т.н., проф.каф. ПІ
Смеляков Кирило Сергійович



18 червня 2024

Зміст

- 1. Мета дослідження та актуальність проблеми
- 2. Еволюція наборів даних та моделей виявлення фейкового контенту (аналітичний огляд)
- 3. Обробка та підготовка даних
- 4. Вибір архітектури згорткової нейронної мережі
- 5. Аугментації
- 6. Результати навчання моделей (експериментальне дослідження)
- 7. Висновки

Мета дослідження

Об'єкт дослідження – моделі виявлення фейкового контенту з заміною ідентичності на базі згорткових нейронних мереж.

Мета дослідження – аналітичний огляд існуючих моделей виявлення фейкового контенту з заміною ідентичності та наборів даних, які використовуються для виявлення маніпуляцій з заміною ідентичності; побудова та навчання декількох моделей виявлення фейкового контенту на основі нейронних мереж.

Результат дослідження – проведено огляд наборів даних маніпуляцій з заміни ідентичності та моделей виявлення фейкового контенту. Виконано навчання декількох моделей виявлення фейкового контенту на основі згорткових нейронних мереж. Навчання проводилося за допомогою сучасного набору даних Deepfake Detection Challenge Dataset.

Термін “DeepFake” (заміна ідентичності) відноситься до техніки на основі глибокого навчання, здатної створювати фальшиві відео шляхом заміни обличчя однієї людини на обличчя іншої. Цей термін з’явився після того, як користувач Reddit під ім’ям “deepfakes” заявив у кінці 2017 року, що розробив алгоритм машинного навчання, який допоміг йому перенести обличчя знаменитостей у порно відео.

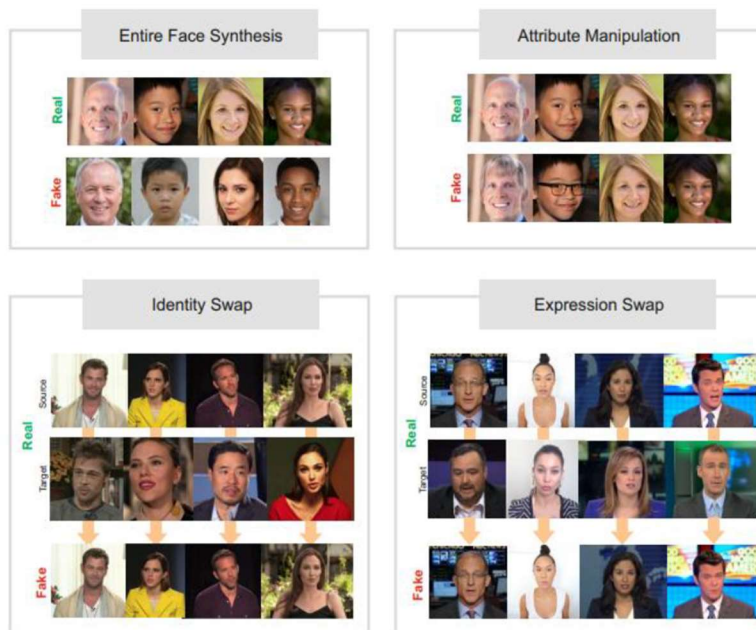
Актуальність дослідження



Приклад заміни ідентичності: однак зображення ліворуч є фейковим.

[M.F. Sohan et al., “A survey on deepfake video detection datasets,” Indonesian Journal of Electrical Engineering and Computer Science, 2023.

Інші приклади маніпуляцій з обличчям



Реальні та фейкові приклади основних типів маніпуляцій з обличчям

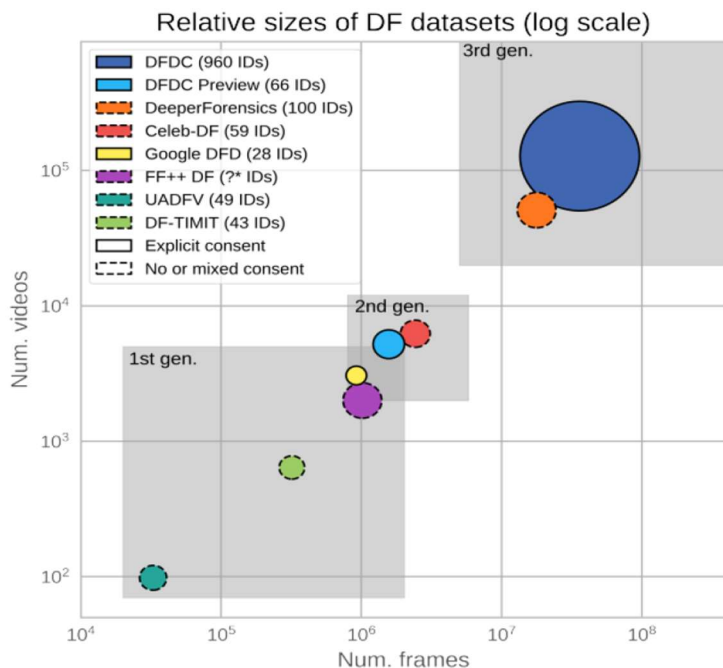
[R. Tolosana et al., 2020]

Постановка задачі

Задачею даного дослідження є ознайомлення з існуючими методами, моделями та наборами даних, які використовуються для виявлення маніпуляцій з заміною ідентичності, їх всебічний аналіз та оцінка та створення та навчання власних моделей виявлення фейкового контенту. Ця мета включає кілька ключових аспектів:

1. Огляд існуючих наборів даних, які використовуються при створенні моделей виявлення маніпуляцій із заміною ідентичності;
2. Огляд існуючих моделей виявлення маніпуляцій із заміною ідентичності
3. Вибір набору даних та архітектури нейронної мережі для навчання
4. Навчання декількох моделей згорткових мереж виявлення фейкового контенту.
5. Аналіз отриманих результатів та їх порівняння з результатами інших досліджень.

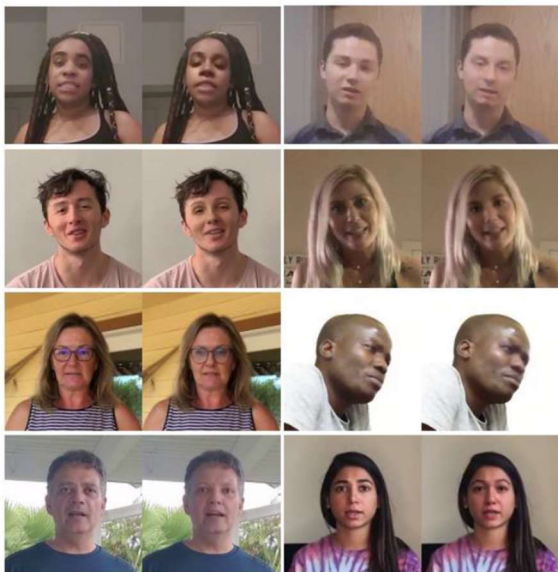
Еволюція датасетів та моделей виявлення фейкового контенту



Генерації датасетів з маніпуляціями типу заміни ідентичності

[B. Dolhansky et al., 2020]

Еволюція датасетів та моделей виявлення фейкового контенту - DFDC



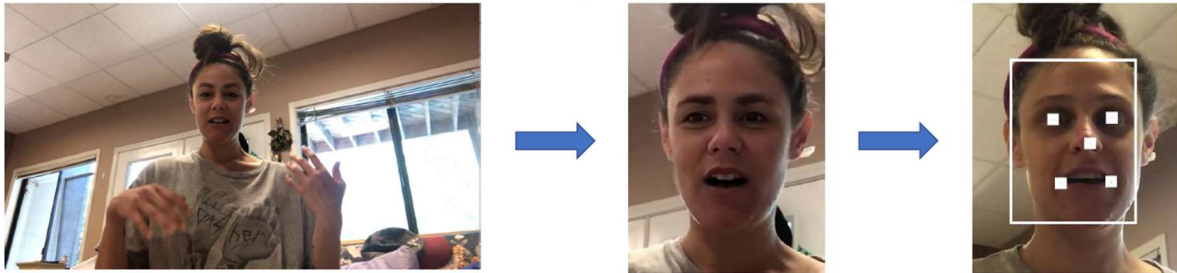
Приклади заміни обличчя з датасету DFDF [B. Dolhansky et al., 2019]

Facebook у співпраці з Microsoft, Amazon, MIT та іншими, запустили наприкінці 2019 року новий виклик під назвою Deepfake Detection Challenge (DFDC).

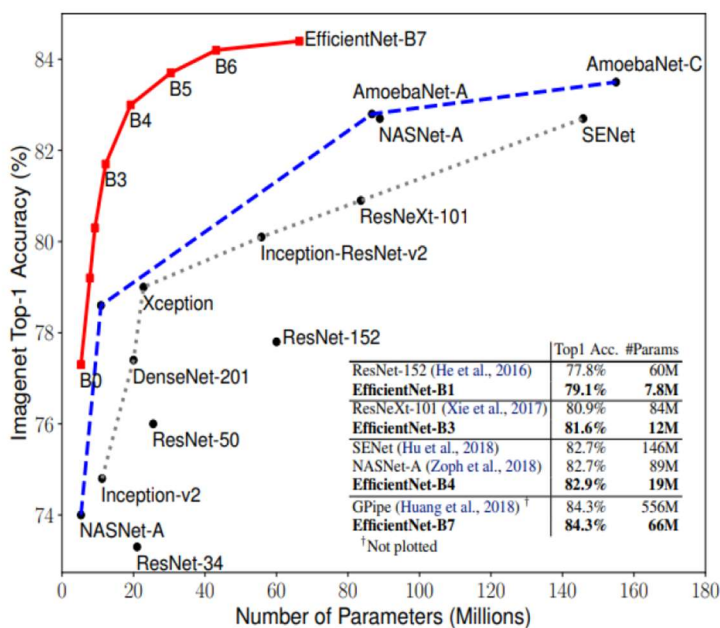
Спочатку вони випустили демонстраційний набір даних, що складається з 1,131 реальних відео (66 акторів) та 4,119 фальшивих відео. Пізніше був випущений повний набір даних DFDC, який містить понад 100 000 відеокліпів з участю 3 426 оплачених акторів (470 ГБ контенту). Ці відеокліпи були створені за допомогою кількох методів, зокрема Deepfake, GAN-заснованих та нелінійних методів навчання.

Обробка та підготовка даних

1. В якості основного набору використовуємо DFDC
2. Застосовуємо підхід на основі аналізу окремих кадрів з відеофайлів
3. За допомогою детектора обличчя (MTCCN) знаходимо рамки обличчя та ключові точки
4. Вирізаємо обличчя
5. Доповнюємо тренувальний набір за допомогою аугментації (розділ 5)
6. Виконуємо тренування на основі попередньо навчених моделей (розділ 4)



Вибір архітектури CNN



Наш вибір:

ResNet-50

Xception

EfficientNet-B4

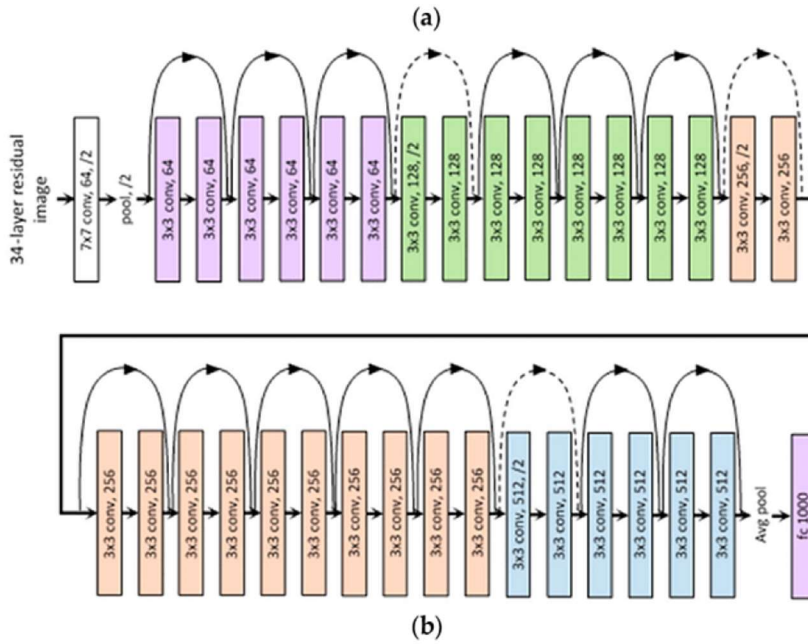
EfficientNet-B5 } + Noisy Student (NS)

EfficientNet-B7

Залежність точності моделі від кількості параметрів

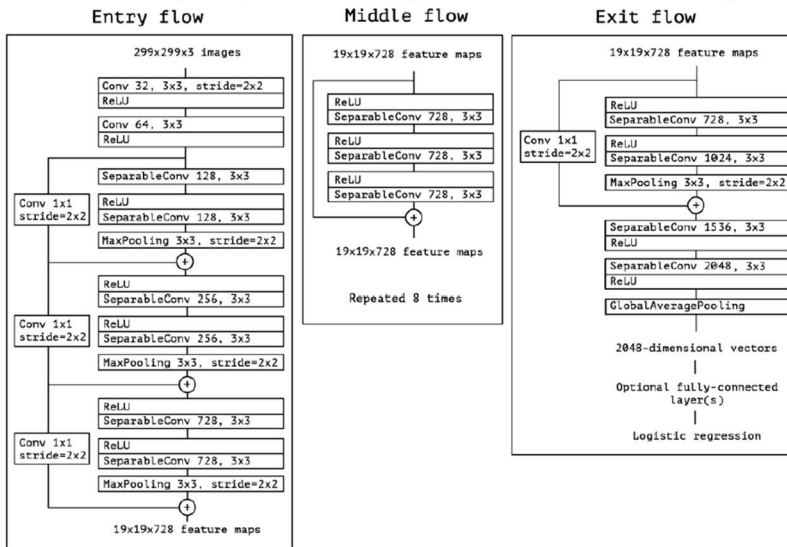
[Mingxing Tan and Quoc V. Le, EfficientNet: Rethinking model scaling for convolutional neural networks, 2019]

Вибір архітектури CNN - ResNet



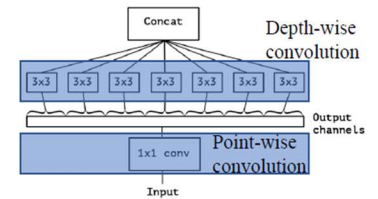
Архітектура моделі ResNet-34
 [Kaiming He et al, “Deep residual learning for image recognition”, 2016]

Вибір архітектури CNN - XceptionNet



Зверніть увагу, що після всіх шарів Conv і SeparableConv виконується пакетна нормалізація (ці шари не включені до діаграми).

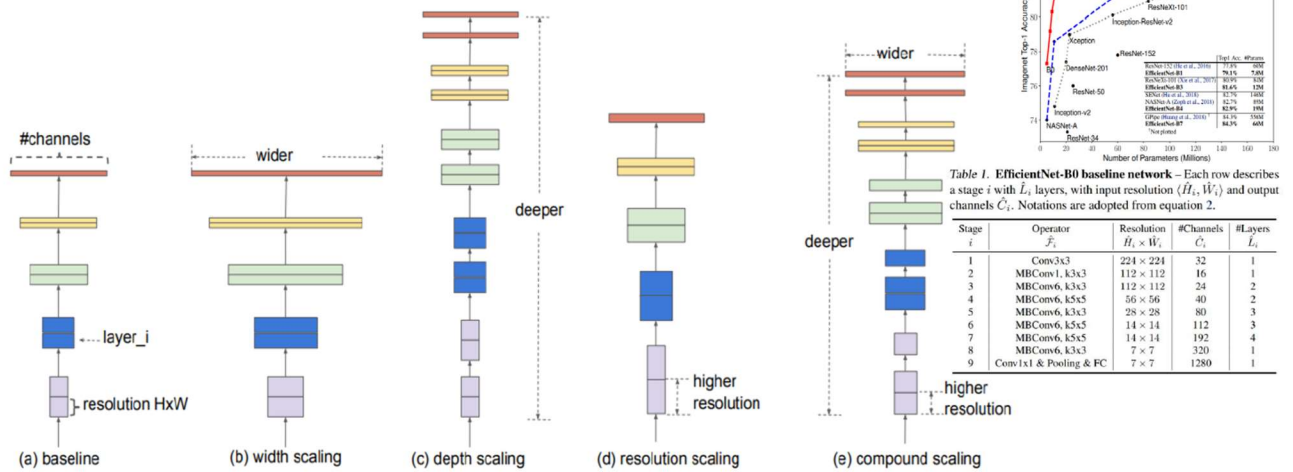
«Екстремальна» версія модуля Inception (SeparableConv)



Архітектура моделі XceptionNet

[Francois Chollet, “Xception: Deep Learning with Depthwise Separable Convolutions”, 2017]

Вибір архітектури CNN - EfficientNet



Масштабування архітектури EfficientNet: (а) базова мережа; (b)-(d) звичайне масштабування, яке збільшує ширину, глибину або роздільну здатність відповідно; (е) – комплексний метод масштабування, який рівномірно масштабує всі три виміри з фіксованим співвідношенням

[M. Tan and Q.V. Le, “EfficientNet: Rethinking model scaling for convolutional neural networks”, 2019]

Вибір архітектури CNN – Noisy Student

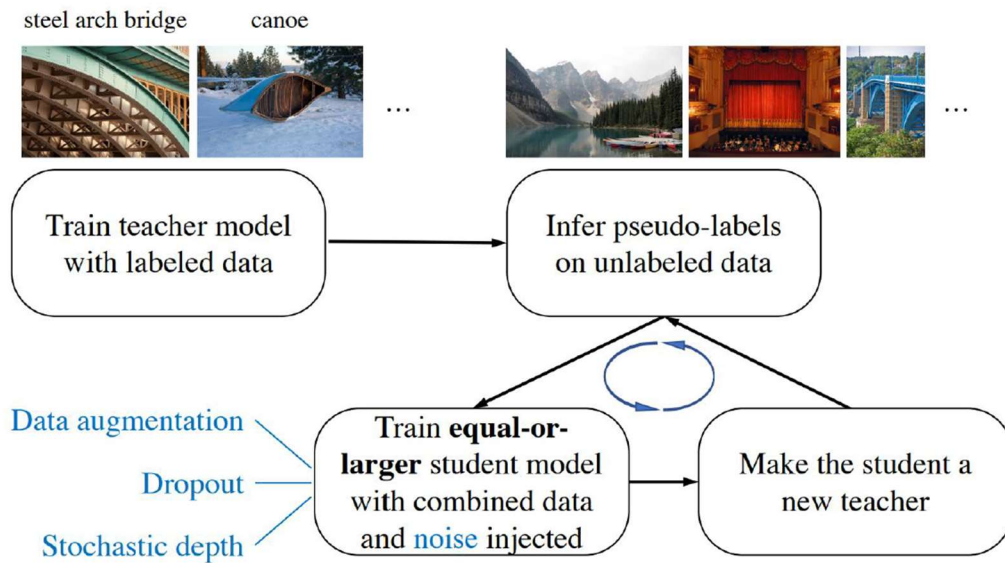
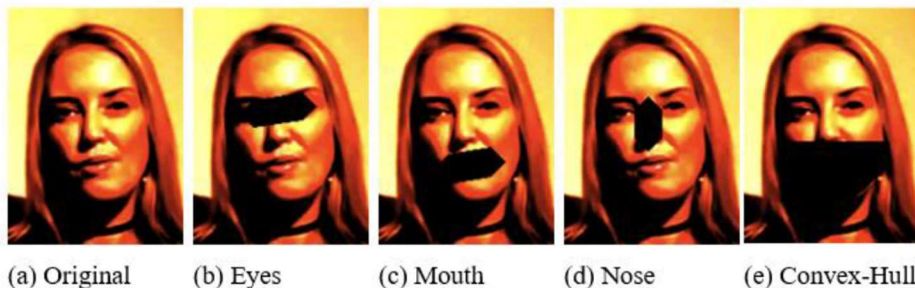


Illustration of the Noisy Student Training (images are from ImageNet)

[Q. Xie, M.T. Luong, E. Hovy, & Q.V. Le, “Self-training with noisy student improves imagenet classification”, 2020]

Аугментації

В якості аугментацій при тренуванні моделі використовувались переважно стандартні аугментації з бібліотеки Albumentations, такі як компресія зображення, накладання гаусового шуму та розмиття, горизонтальне відображення, змінення розміру, змінення яскравості, контрастності та насиченості, Fancy PCA, перетворення кольорового зображення в градації сірого та застосування афінних перетворень (переміщення, масштабування та обертання). Зазначенні перетворення здійснювались з різними імовірностями від 0.05 до 0.7. Також використовувались Face-Cutout аугментації.



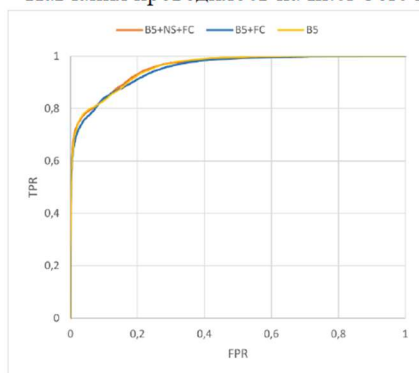
Face-Cutout аугментації: (a) оригінальне обличчя без накладання аугментацій; (b), (c), (d) вирізання очей, рота, та носа відповідно; (e) випадкові вирізання опуклих оболонок з заповненням вирізання випадковими значеннями [Das, Sowmen et al., “Towards solving the deepfake problem: An analysis on improving deepfake detection using dynamic face augmentation”, 2021]

Результати навчання моделей

Результати навчання моделей виявлення фейкового контенту з використанням набору даних DFDC

Model	Parameters (millions)	Train time (hours)	BCE	AUC (%)	AP (%)
B4+NS+FC	19	4,4	0,374	95,2	99,49
B5	30	5,5	0,347	95,7	99,53
B5+FC	30	5,7	0,337	95,2	99,49
B5+NS+FC	30	5,8	0,302	95,6	99,53
B7+NS+FC	43	9,4	0,338	95,2	99,47
Xception	22,9	3,5	0,351	93,2	99,21
ResNet50	25,6	3,3	0,322	93,6	99,31

Навчання проводилось на Intel Core i7-13700K 16C/24T з 64 Gb RAM, 2+2 Tb SSD та GeForce RTX 3090 24 Gb.



ROC-криві для моделей B5, B5+FC, B5+NS+FC

Висновки

Під час виконання кваліфікаційної роботи було:

- дано визначення поняття фейкового контенту та розглянуто його види;
- проведено аналітичний огляд наборів даних із маніпуляціями з заміни ідентичності та моделей їх виявлення;
- обрано набір даних для проведення досліджень;
- обрано архітектури згорткових нейронних мереж для проведення досліджень;
- проведено навчання моделей згорткових нейронних мереж на основі архітектур EfficientNet, XceptionNet та ResNet та отримано метрики якості проведеного навчання;
- проведено аналіз отриманих результатів та їх порівняння з даними інших досліджень.

Апробація результатів

Результати кваліфікаційної роботи було опубліковано у 2 виданнях:

1. Кіщенко Ю. О., Смеляков К. С. Розробка моделі виявлення фейкового контенту на основі архітектури EfficientNet. Матеріали 28-го міжнародного молодіжного форуму «Радіоелектроніка та молодь у XXI столітті»: Зб. матеріалів форуму. Т. 6. Харків, 2024. С. 435–436.

2. Smelyakov K., Kitsenko Y., Chupryna A. Deepfake detection models based on machine learning technologies. 2024 IEEE open conference of electrical, electronic and information sciences (estream). 2024. P. 1–6. URL: <https://doi.org/10.1109/eStream61684.2024.10542582>.



ДОДАТОК В

Апробація результатів кваліфікаційної роботи

1. Диплом за результатами доповіді на 28-му міжнародному молодіжному форумі «Радіоелектроніка та молодь у ХХІ столітті»



2. Тези доповіді на 28-му міжнародному молодіжному форумі
«Радіоелектроніка та молодь у ХХІ столітті»

МІНІСТЕРСТВО ОСВІТИ І НАУКИ УКРАЇНИ
ХАРКІВСЬКИЙ НАЦІОНАЛЬНИЙ УНІВЕРСИТЕТ
РАДІОЕЛЕКТРОНІКИ

МАТЕРІАЛИ ХХVІІІ МІЖНАРОДНОГО МОЛОДІЖНОГО
ФОРУМУ

**«РАДІОЕЛЕКТРОНІКА ТА МОЛОДЬ
У ХХІ СТОЛІТТІ»**

16 – 18 квітня 2024 р.

Том 6

**КОНФЕРЕНЦІЯ
«ІНФОРМАЦІЙНІ ІНТЕЛЕКТУАЛЬНІ СИСТЕМИ»
INFORMATION INTELLIGENT SYSTEMS**

Харків 2024

РОЗРОБКА МОДЕЛІ ВИЯВЛЕННЯ ФЕЙКОВОГО КОНТЕНТУ НА ОСНОВІ АРХІТЕКТУРИ EFFICIENTNET

Кіценко Ю. О.

Науковий керівник – д.т.н., проф. Смеляков К. С.

Харківський національний університет радіоелектроніки, каф. ПІ

м. Харків, Україна

e-mail: yurii.kitsenko@nure.ua

The research is devoted to efficiency evaluation of modern deepfake detection models based on convolutional neural networks (CNN). In today's world, with the growing influence of digital technology and increasing volume of information on the internet, detection of fake images and videos has become increasingly important. Fake content spread through social media and other platforms can cause serious damage, from personal attacks to manipulation of public opinion on a global level. During the study, we trained a model based on the EfficientNet architecture [1]. The model was trained on the Deepfake Detection Challenge dataset [2].

Поява та розповсюдження фейкового контенту в інтернеті є складним явищем, яке має все більший вплив на суспільство, політику та економіку. Зростання доступності технологій глибокого навчання та штучних нейронних мереж робить можливим виготовлення високоякісних фейків. З'явилося багато потужних інструментів, що дозволяють змінювати фотографії, створювати відео з фіктивним контентом і навіть генерувати тексти. Поява моделей виявлення фейкового контенту є відповіддю на зростаючу загрозу фальсифікації інформації в цифровому просторі. З огляду на швидкі та значні зміни в технологіях створення фейків, виникає потреба у високоефективних інструментах та методах для їх виявлення. Серед цих інструментів важливу роль відіграють моделі, засновані на технологіях машинного навчання.

У роботі було побудовано модель розпізнавання фейкового контенту, засновану на архітектурі EfficientNet [1]. EfficientNet – це сімейство нейронних мереж, розроблених для досягнення високої ефективності за рахунок оптимального балансу між розміром мережі та її продуктивністю. Основою архітектури є згортова базова модель, до якої застосовуються масштабуючі коефіцієнти, що визначають розмір мережі. Такий підхід дозволяє досягти високої точності на завданнях класифікації зображень за мінімальної кількості параметрів, що робить EfficientNet однією з найбільш ефективних архітектур для роботи з обмеженими ресурсами.

EfficientNet пропонує кілька типів моделей, включаючи B0-B7, кожна з яких відрізняється за розміром і кількістю параметрів. Моделі від B0 (найменша) до B7 (найбільша) представляють собою послідовне збільшення глибини, ширини і роздільної здатності мережі. Більші моделі,

такі як V7, забезпечують більшу точність класифікації, але вимагають більшого обсягу ресурсів для навчання і виконання. Менші моделі, наприклад V0, мають меншу кількість параметрів і є більш ефективними при роботі з обмеженими обчислювальними ресурсами, при цьому зберігаючи високу точність. У дослідженні було використано найбільшу з моделей – модель V7.

Для навчання моделі використовувався набір даних DeepFake Detection Challenge (DFDC) [2]. DFDC є найбільшим доступним публічно набором даних відео з підміною облич. Він включає понад 100 000 відеокліпів з участю 3426 платних акторів. Ці відео були створені за допомогою кількох методів, таких як Deepfake, методів на основі застосування генеративно-змагальних мереж та інших методів, що не використовували підходів машинного навчання. DFDC був використаний у проведенні конкурсу Kaggle, що сприяв розвитку засобів виявлення фейкового контенту [3].

Перед навчанням моделі було виконано попередню підготовку даних, яка складалася з наступних етапів:

1. Захват окремих кадрів з відеофайлів.
2. Знаходження облич на отриманих кадрах [4].
3. Вирізання облич для подальшого детального аналізу.
4. Генерація «згорток» (folds) для перехресної валідації.
5. Аугментація даних (augmentation).

Зазначимо, що використаний підхід до виявлення "на основі кадр за кадром" досить вдало зарекомендував себе у багатьох сценаріях. Основною з його переваг є його обмежена обчислювальна вимогливість.

Для побудови та навчання моделі використовувалась відкрита бібліотека машинного навчання Pytorch. Навчання проводилось на протязі 40 епох по 2500 ітерацій та було отримано зважену похибку 0.25, що є досить непоганим результатом, порівняно з іншими сучасними моделями на основі згорткових нейронних мереж [5].

Список використаних джерел:

1. Tan M., Le Q. Efficientnet: Rethinking model scaling for convolutional neural networks // International conference on machine learning. 2019. Vol. 97. P. 6105-6114.
2. Dolhansky B. et al. The deepfake detection challenge (dfdc) dataset //arXiv preprint arXiv:2006.07397. – 2020.
3. Deepfake Detection Challenge // Kaggle. URL <https://www.kaggle.com/c/deepfake-detection-challenge> (дата звернення: 05.03.2024).
4. Smelyakov K. et al. The neural network models effectiveness for face detection and face recognition // IEEE Open Conference of Electrical, Electronic and Information Sciences (eStream). 2021. P. 1-7.
5. Tolosana R. et al. Deepfakes and beyond: A survey of face manipulation and fake detection //Information Fusion. 2020. Vol. 64. P. 131-148.

3. Стаття за результатами доповіді на 8-ій Міжнародній конференції “Electrical, Electronic and Information Sciences” (Estream 2024)

2024 IEEE OPEN CONFERENCE OF ELECTRICAL, ELECTRONIC AND INFORMATION SCIENCES (ESTREAM)

PROCEEDINGS OF THE CONFERENCE

April 25, 2024
Vilnius, Lithuania

Edited by:

Dalius Navakauskas

Šarūnas Paulikas

Tomyslav Sledevič

Dainius Udris

Organized by:



Vilnius Gediminas Technical University

Sponsored by:



IEEE Lithuania Section

IEEE Lithuania Joint SP/CIS/COM Chapter

IEEE Lithuania Joint AP/ED/MTT Chapter

IEEE Lithuania Computer Chapter

IEEE Lithuania Education Chapter

IEEE part number: CFP2447Z-ART

ISBN: 979-8-3503-5241-2

Online ISSN 2690-8506

Deepfake Detection Models Based on Machine Learning Technologies

Kirill Smelyakov

Department of Software Engineering
Kharkiv National University of Radio
Electronics
Ukraine, Kharkiv, Nauka Ave, 14
kyrylo.smelyakov@nure.ua

Yuriy Kitsenko

Department of Software Engineering
Kharkiv National University of Radio
Electronics
Ukraine, Kharkiv, Nauka Ave, 14
yurii.kitsenko@nure.ua

Anastasiya Chupryna

Software Engineering Department
Kharkiv National University of Radio
Electronics
Kharkiv, Ukraine
anastasiya.chupryna@nure.ua

Abstract— The paper is devoted to efficiency evaluation of modern deepfake detection models based on convolutional neural networks (CNN). In the context of rapid development of digital technologies and increasing volume of information on the internet, the relevance of detecting fake images, videos, and textual materials becomes increasingly significant. Fake content, spread through social networks and other platforms, can have serious consequences, ranging from individual malicious attacks to manipulations of public opinion on a global level. We have built and trained several models for detecting fake content using convolutional neural networks. The training was performed using Deepfake Detection Challenge Dataset. During the study, we carried out the comparative analysis of the created models. Obtained results were compared with a number of recent publications.

Keywords—Deepfake, Convolutional Neural Network, Face detection, EfficientNet, Effectiveness

I. INTRODUCTION

The emergence and spread of deepfakes on the Internet are complex phenomena that has an increasing impact on society, politics, and the economy. The growing availability of machine learning technologies and deep neural networks makes it possible to produce high-quality deepfakes. There are many powerful tools that allow you to modify photos, create videos with fictitious content, and even generate texts. Algorithms for social media and content recommendation platforms contribute to the spread of fakes. Content that evokes an emotional response or support has a high chance of going viral, regardless of its credibility.

A. Overview of Modern Challenges

Nowadays, information spreads very quickly on the Internet, mostly through social media, blogs, and other platforms. This allows fake content to influence public opinion even before its inaccuracy is discovered. Different groups of people and individuals can create deepfakes for a variety of reasons, such as political agendas, earnings, pseudo-humorous purposes, or attempts to influence public opinion. The use of bots and automated systems allows you to quickly spread fake content and create the impression of broad support or discussion in the Internet environment.

The emergence of deepfake detection models is a response to the growing threat of information falsification in the digital world. Given the rapid and significant changes in fake creation

technologies, there is a need for highly effective tools and methods to detect them. Models based on machine learning technologies play an important role here.

With the development of machine learning and deep neural networks, deepfake detection models are becoming more powerful and adapting to new forms of falsification. They are based on algorithms that can learn and recognize patterns in large amounts of data, detecting anomalies and falsehoods. These models are used to analyze a variety of content types, such as text, photos, and videos. They can detect false informational articles, graphic processing, or AI-generated images, as well as fake video footage.

One of the important directions in deepfake detection is the use of models that take into account the context, emotional aspects and subjective nuances in text and multimedia data. This helps to avoid detecting legitimate information that may be emotional or subjective in nature as fake.

Overall, developing and improving of deepfake detection models remains an important task to ensure the trust into information in a digital world where the manipulation and spread of deepfakes is becoming increasingly common. In view of the above, the study of deepfake detection models based on machine learning technologies is a hot topic today.

B. Evolution of Deepfake Detection Models

“Identity Swap” manipulations have attracted considerable attention as among researchers and also among general public because of their potential impacts on information security and social trust [1]-[3]. The technology, which allows the replacement of people's faces in videos or images, opens the door to a variety of applications, from entertainment and art to the creation of fake news or disinformation campaigns. This curiosity is amplified by the ability to create videos that look convincingly realistic at first glance, which calls into question the authenticity of digital content. However, along with technological capabilities, “Identity Swap” raises important questions about the ethical boundaries of using such manipulations, requiring the development of new methods for detecting and verifying digital content in order to maintain trust in the digital age.

Let's take a look at how methods for detecting face swap manipulations have evolved in the context of the emergence of new relevant datasets. One of the first publicly available datasets

according to [4] with original and fake videos having identity swap manipulations is UADFV dataset [5]. This dataset contains 49 original videos, which were the basis for creation of the same number of deepfake videos. Video duration is in average about 11.14 seconds long. Typical resolution is 294×500 pixels. This dataset was used for the first time in [6]. The approach developed there was based on the detection of eye blinking in videos. Eye blinking is a physiological trait that is poorly represented in synthetic fake videos. This method was tested by the authors on standard datasets for the detection of eye blinking. It showed promising results in detecting deepfake videos.

The next step in identifying identity manipulation was the work of Korshunov and Mareel [7]. The DeepfakeTIMIT dataset with 620 fake videos was presented in this study. Face swapping was performed using technology based on generative adversarial networks (GANs)¹. To create the dataset, 16 pairs of people with similar faces were selected from the open database VidTIMIT. Two models were developed for each of the 32 participants: a low-quality (LQ) model with an image size of 64×64 and a high-quality (HQ) model with a size of 128×128. Taking into account the availability of 10 videos per participant in the VidTIMIT database, 320 videos were created for each version, for a total of 620 videos with swapped faces. In the paper, it was shown that facial recognition algorithms based on VGG and Facenet are not suitable for detecting facial manipulation. Several basic facial replacement detection algorithms were also evaluated, and it was shown that the lip-syncing approach failed to detect inconsistencies between lip movement and audio. At the same time, the authors showed that an approach based on measures of image quality using an SVM classifier can detect fake videos with an EER of 3.3% and 8.9% for LQ and HQ scenarios, respectively.

In [8] automated benchmark for detection of facial manipulations was developed. This benchmark is based on DeepFakes², Face2Face [9], FaceSwap² and NeuralTextures [10] tools of face manipulations. The benchmark is publicly available. It contains a hidden test dataset and public dataset with more than 1.8 million of manipulated images. This dataset is called FaceForensics++ and is an order of magnitude larger than previous publicly available similar datasets. Based on this data, the authors analyzed 7 falsification detectors [11]–[16]. During the benchmark, the XceptionNet model obtained the best results [11]. XceptionNet is a traditional CNN based on a separate convolution with residual connections. The model was trained on the ImageNet dataset and adapted by replacing the last fully connected layer with two outputs.

In 2019 DeepFakeDetection (Google DFD) dataset was created with the support of Google [17]. This dataset includes 363 real-life videos (28 actors in 16 different scenes). It also includes 3068 fake videos, created with a help of DeepFake FaceSwap as in previous case. Later, this dataset was included into the FaceForensics++ dataset.

In [18] a new large-scale fake video dataset CelebDF was introduced. It contains 5,639 high-quality fake celebrity videos created through an advanced synthesis process. With a help of

presented dataset, the authors compared 9 methods for deepfake detection that were publicly available or received directly from the authors (see, for example, [5], [8], [16], [19]–[21]). During the comparison Xception [8] and DSP-FWA [21] showed the best results. Xception approach is based on the XceptionNet model [11] trained on the FaceForensics++ dataset. Three variants of Xception were used, namely: Xception-raw, Xception-c23 and Xception-c40. Xception-raw was trained on raw videos, while Xception-c23 and Xception-c40 are trained on H.264 videos with medium (23) and high (40) compression ratios, respectively. DSP-FWA is an advanced FWA-based method that includes the Spatial Pyramid Pooling (SPP) module [21] for better handling of target faces of different resolutions. It was shown that the performance of the methods under consideration decreases with increasing of compression ratio. In particular, the performance of FWA and DSP-FWA is significantly degraded on recompressed video, while the performance of Xception-c23 and Xception-c40 is not significantly affected. It was expected, as the latter models have been trained on compressed H.264 videos, so they are more resilient in this situation.

In the end of 2019 Facebook, in collaboration with Microsoft, Amazon, MIT, and others, launched a new challenge called the Deepfake Detection Challenge (DFDC) [22]. Initially, they released a demo dataset consisting of 1,131 real videos (66 actors) and 4,119 fake videos. Later, the full DFDC dataset was released [23], which contains more than 100,000 video clips featuring 3,426 paid actors (470 GB of content). These videos were created using several methods, including Deepfake, GAN-based, and non-linear learning methods.

In [24] a large test dataset for detection of facial falsifications was presented. The first version of this dataset, DeeperForensics-1.0, is now one of the largest face-tampering datasets, containing of 60,000 videos and 17.6 million frames. It was about 10 times larger than the size of existing datasets of this type. The fake videos were generated with a help of new face-swap end-to-end system³ that was offered. In addition, a comprehensive study was conducted in this paper, which evaluated five detection methods.

It should be noted that from the moment when the first deepfake detection datasets appeared in this list till the appearance of the last ones, there were significant visual improvements in the contents and realism of deepfake videos. As a result, datasets with identity swap manipulations are often divided into several generations (see, for example, [4]) in the following manner: UADFV, DeepfakeTIMIT and FaceForensics++ are related to the first generation, CelebDF, Google DFD, DFDC Preview – to the second, and DeeperForensics and DFDC – to the third.

Recently, several more datasets with facial manipulations have appeared [26]. In [27] WildDeepfake dataset was presented. WildDeepfake consists of 1,180,099 high-quality images of 7,314 facial sequences extracted from 707 videos (both deep fakes and real) collected from web resources, so it contains a greater variety of scenes, faces, and actions. Also in

¹ <https://github.com/deepfcase/fakeswap>

² <https://github.com/MarekKowalski/FaceSwap/>

³ <https://github.com/EndlessSora/DeeperForensics-1.0>

this study, two attention-based deep fake detection networks (ADDNets) were proposed.

In [28] Korean Deep Fake Detection Dataset (KoDF) was introduced. It includes 175,776 fake video clips and 62,166 real video clips from 403 subjects. The fake videos were created using six different synthesis models. Another interesting dataset is ForgeryNet [29].

C. Task For This Research

The main objective of this study is to train several deepfake detection models of convolutional neural network. A number of convolutional neural network architectures, such as ResNet, EfficientNet, and Xception, will be used. This will allow to perform the deeper analysis and recognize complex patterns that indicate image or video tampering. Each of these models has unique characteristics and approaches to data processing, which provide a more comprehensive and effective deepfake detection together.

The objectives of this study are primarily aimed to deepening of the understanding of the current state of identity swap manipulation detection using machine learning technologies, as well as raising awareness about the challenges and opportunities facing the scientific community in this direction.

II. METHODS AND MATERIALS

A. Data Preparation

The data preparation stage is critical when building a deepfake detection model. Data preparation has a significant impact on the efficiency, accuracy, and reliability of the final model. We selected DFDC [22] as our main dataset. To prepare the data, we will follow the scheme of work [30] and use an approach based on the analysis of individual frames of video files. One of the main advantages of this method is its limited time and computational demand, especially when processing high-resolution or long-duration videos. Among the drawbacks of the method, it can be mentioned that analysis of individual frames can miss contextual or temporal relationships, which become noticeable only when considering the video as a whole. Also, this approach does not make it possible to identify fakes in the audio track. We note that while this kind of manipulations is present in DFDC, it was not taken into account as fakes during the competition [22].

The frames obtained at the first stage were analyzed with a help of a face detector in order to find the boxes of found faces and face landmarks. After analysis of face detectors (see, e.g., [31], [32]) MTCNN was chosen as the face detector [33]. The obtained face boxes made it possible to cut out and form a set of faces that will be used in CNN training further.

Since DFDC dataset contains deepfake videos, originals from they were created and link to the original for each deepfake video, it's possible to create a difference SSIM masks containing 1 for modified pixels and 0 otherwise [34]. These masks were used in the generation of augmentations, which will be discussed further.

B. Selection of models

We have selected three architectures of convolutional neural networks for training here: EfficientNet [35], XceptionNet [8] and ResNet [25].

EfficientNet is a family of deep neural networks that vary in architecture size and complexity. The most popular models are EfficientNet-B0 to EfficientNet-B7. The number in the name of the model specifies its size. For example, EfficientNet-B0 has the fewest number of parameters, while EfficientNet-B7 is the largest and most complex model with a large number of parameters.

Each of the EfficientNet models has a correspondence between network depth, width, and resolution. Typically, larger models have more layers, more width, and higher resolution, allowing them to tackle more complex tasks, but they also require more computing resources to train and apply. In addition, EfficientNet also includes some model variations, such as EfficientNet-Lite, which were optimized for application on mobile devices with limited resources, and also other models that can be applied for specific tasks and data volumes. We have trained B4, B5 and B7 architectures.

The second of the selected architectures is XceptionNet. This architecture represents an improved version of Inception model proposed in Google's InceptionV3. The main idea behind XceptionNet is to adopt a group of cross-channel and group convolutions, making the network more efficient and easier to train.

The main feature of XceptionNet is the usage of group convolutions to reduce the number of parameters and computational complexity, as well as channel separation, which allows to obtain independent convolutions for each input channel. These features allow XceptionNet to use compute resources efficiently and allow to achieve high accuracy for various problems in computer vision.

The last chosen architecture is ResNet (Residual Network). This architecture differs with presence of connections, that pass data between different layers or so-called "residual blocks". This architecture has a great impact in computer vision. The main idea behind ResNet is to add input data to the output of each block instead of replacing it with calculated values. This allows to reduce the gradient vanishing problem and promotes information retention during deep network training.

Usage of these "short connections" allows for more efficient and faster learning for ResNet models. This makes possible training of deep neural networks with dozens and even hundreds of layers, providing high accuracy on a number of computer vision tasks, including image classification, object detection, and semantic segmentation.

All selected models were initialized with pre-trained ImageNet weights. Also, we used pre-trained model with Noisy Student approach for EfficientNet.

The Noisy Student technique is one of the self-learning methods to improve the efficiency of neural networks. It has been featured in 2019 in Google AI paper [36] and is



Fig. 1. Original face without augmentation(a). Cutting out the eyes (b), mouth (c), and nose (d). Random cuts of convex shells with filling of the clipping with black color (e) [30].

mainly used in the context of learning using unlabeled data. The idea is to train neural network in several steps: initially train the model with a help of labeled data and then train it with a new unseen unlabeled data with a help of noise-adding and regularization techniques.

The main principle of the Noisy Student technique is that a model that is already trained can be used to label unseen data from the other dataset. Created labels can be used to train a new model on this new dataset. Regularization techniques, such as dropout or augmentation, are also used during this training process to increase the model's stability. In this way, Noisy Student allows you to efficiently use unlabeled data to improve the quality of the model without the need for the costly and time-consuming process of manual data markup.

C. Augmentations

Standard augmentations from the Albumentations library⁴, such as image compression, Gaussian noise overlay and blur, horizontal flip, resizing, brightness, contrast, and saturation, Fancy PCA, converting to gray, and applying transformations such as moving, scaling, and rotating were used as augmentations for training the model. These transformations were applied with different probabilities from 0.05 to 0.7.

In addition to the augmentations mentioned above, we used a slightly simplified version of the Face-Cutout approach⁵, which was proposed in [30]. Face-Cutout, uses face landmarks to augment train images. Face landmarks include the positions of the eyes, ears, nose, mouth, and jawline. These positions are used to create polygons for Face-Cutout (see Figure 1). As it was already mentioned above, SSIM difference masks were generated at the data preparation stage. The Face-Cutout algorithm takes as input face image and it's corresponding SSIM difference mask to generate an augmented image. The mask is also used to reduce the cutting of the altered part of the face so that the model can focus on the manipulation area.

D. Metrics

The following values were used as model training metrics: BCE, AUC and average accuracy.

BCE (logarithmic loss) is a metric for measuring the classification models accuracy. To obtain BCE the formula can be used

$$BCE = -\frac{1}{n} \sum_{i=1}^n \{y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)\}. \quad (1)$$

Here n is a number of samples in the test dataset, y_i is the actual class of sample i and \hat{y}_i is the probability that the sample i belongs to class 1.

AUC (Area Under the Curve) is a metric used to assess the quality of classification model, especially for their ability to distinguish between positive and negative classes. The AUC measures the area under the ROC (Receiver Operating Characteristic) curve, which reflects the relationship between true positive rate (TPR) and false positive rate (FPR).

The ROC curve is a line that shows how the TPR of the model depends on FPR at different cut-off threshold values. The closer the AUC is to 1, the better the model separates classes, i.e., the larger the area under the ROC curve, the better its performance. The AUC can take values ranging from 0 to 1, where a value closer to 1 indicates better model quality, and a value around 0.5 indicates a random classifier behavior.

AP (Average Precision) is an indicator that resumes the precision-recall curve as the average precision value achieved for each threshold, with an increase of recall from the previous threshold used as a weight coefficient:

$$AP = -\frac{1}{n} \sum_{i=1}^n (R_i - R_{i-1}) P_i. \quad (2)$$

Here P_i and R_i is the precision and recall values for i -th threshold.

⁴ <https://github.com/albumentations-team/albumentations>

⁵ <https://www.kaggle.com/competitions/deepfake-detection-challenge/discussion/145721>

TABLE I. RESULTS OF TRAINING DEEPFAKE DETECTION MODELS WITH A HELP OF DFDC DATASET

Model	Parameters (millions)	Train time (hours)	BCE	AUC (%)	AP (%)
B4+NS+FC	19	4.4	0.374	95.2	99.49
B5	30	5.5	0.347	95.7	99.53
B5+FC	30	5.7	0.337	95.2	99.49
B5+NS+FC	30	5.8	0.302	95.6	99.53
B7	43	9.4	0.338	95.2	99.47
Xception	22,9	3.5	0.351	93.2	99.21
ResNet50	25,6	3.3	0.322	93.6	99.31

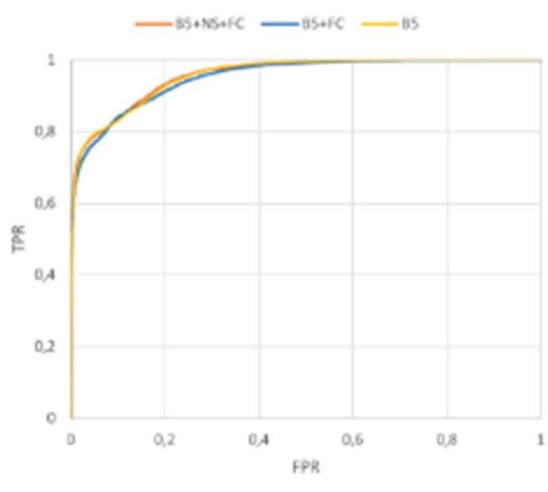


Fig. 2. ROC curves for models B5, B5+FC, B5+NS+FC.

III. EXPERIMENT RESULTS

In this paper 7 models of convolutional neural networks based on EfficientNet, XceptionNet and ResNet architectures were trained. Training results are shown in Table 1. For each model, the following parameters were calculated: the number of model parameters, training time, BCE, AUC, and AP metrics. The names in the table contain model architecture; NS in the model's name means that the model was trained using the Noisy Student method, FC means that a simplified version of the Face Cutout algorithm was used when training the network. The training was performed on Core i7 64 Gb GeForce RTX 3090 24 Gb.

According to Table 1, the Xception and ResNet50 models perform slightly worse (~93% AUC) than the EfficientNet-based models (~95% AUC). Quite unexpectedly, almost all models based on the EfficientNet architecture showed very close AUC values (~95%) and the difference between them can only be observed in BCE values. At the same time, we did not observe significant changes between the training results with and without Noisy Student and Face Cutout. While in [30] the

increase of ~5% in AUC was observed when using the Face Cutout method. At the same time, the AUC values, for example, for the XceptionNet model were slightly lower than those obtained in this paper when using the DFDC dataset for training and using the Face Cutout method. In [30] AUC of 95.66% was obtained when training with a combined dataset. This is slightly higher than our value.

ROC curves for models B5, B5+FC, B5+NS+FC are shown on Figure 2. Model names contain the architecture type; NS in the model's name means that the model was trained using the Noisy Student method; FC means that a simplified version of the Face Cutout algorithm was used when training the network.

As one can see from the Figure 2, we do not observe significant difference when for models trained using the Noisy Student and Face Cutout techniques. In [30] the 5% increase of AUC was observed. It is planned to investigate in more detail the influence of these techniques on models' accuracy in further studies.

IV. CONCLUSION

In this paper, we examined in detail the various approaches and algorithms used in modern deepfake research. We have trained 7 deepfake detection models based on EfficientNet, XceptionNet and ResNet architectures. It was shown that AUC of all trained models is higher than 93%. In the same time the EfficientNet architecture showed better results (~95% AUC) comparing to XceptionNet (93.2% AUC) and ResNet (93.6% AUC). This is expected result for B5 and B7 EfficientNet model sizes as they have more parameters, but this was unexpected result for B4 model size, which showed results comparable with B5 and B7 models.

We would like to emphasize that nowadays the effectiveness of deepfake detection models is continuously improving due to innovative technical solutions. As in many other areas, such as, for example, the development of air defense or antiviruses, the creation and detection of deepfakes are in constant struggle, and the emergence of new models of deepfakes creation leads to the emergence of new detection models. Machine learning technologies are involved into this battle on both sides. Thus, it is essential at this time to pay attention to the social and ethical aspects of creating deepfakes with a help of deep learning

methods and other similar technologies. These aspects cover a wide range of issues, from individual rights to collective security and trust in society. It should be noted that the development of these technologies should go hand in hand with the discussion of ethical and social norms and rules governing the use and distribution of fake content.

Given the rapid technological development evidenced by a significant amount of research in this direction, we can expect significant progress in this area in the near future. Raising awareness of deepfakes, developing new algorithms and detection approaches, as well as a deeper understanding of the social aspects of this problem, will provide more effective protection of modern society from the negative impacts of fake content.

REFERENCES

- [1] D. Citron, "How deepfake undermine truth and threaten democracy," 2019. [Online]. Available: <https://www.ted.com>
- [2] R. Cellan-Jones, "Deepfake videos double in nine months," 2019. [Online]. Available: <https://www.bbc.com/news/technology-49961089>
- [3] BBC Bitesize, "Deepfakes: What are they and why would I make one?" 2019. [Online]. Available: <https://www.bbc.co.uk/bitesize/articles/zfkwvcq>
- [4] R. Tolosana et al., "Deepfakes and beyond: A survey of face manipulation and fake detection," *Information Fusion*, 64, 2020, pp. 131-148.
- [5] X. Yang, Y. Li, S. Lyu, "Exposing deep fakes using inconsistent head poses," *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019, pp. 8261-8265.
- [6] Y. Li, M. Chang, S. Lyu, "In ietu oculi: Exposing AI generated fake face videos by detecting eye blinking," *IEEE International workshop on information forensics and security (WIFS)*, 2018, pp. 1-7.
- [7] P. Korshunov, S. Marcel, "Deepfakes: a new threat to face recognition? assessment and detection," *arXiv preprint arXiv:1812.08685*, 2018.
- [8] A. Rossler et al., "FaceForensics++: Learning to detect manipulated facial images," *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 1-11.
- [9] J. Thies, M. Zollhofer, M. Stamminger, C. Theobalt, and M. Nießner, "Face2Face: Real-time face capture and reenactment of RGB videos," *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2387-2395.
- [10] J. Thies, M. Zollhofer, M. Nießner, "Deferred neural rendering: Image synthesis using neural textures," *Acm Transactions on Graphics (TOG)*, 38(4), 2019, pp. 1-12.
- [11] F. Chollet, "Xception: Deep learning with depthwise separable convolutions," *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1251-1258.
- [12] J. Fridrich, J. Kodovsky, "Rich models for steganalysis of digital images," *IEEE Transactions on Information Forensics and Security*, 7(3), 2012, pp. 868-882.
- [13] D. Cozzolino, G. Poggi, L. Verdoliva, "Recasting residual-based local descriptors as convolutional neural networks: an application to image forgery detection," *Proceedings of the 5th ACM workshop on information hiding and multimedia security*, 2017, pp. 159-164.
- [14] B. Bayar, M.C. Stamm, "A deep learning approach to universal image manipulation detection using a new convolutional layer," *Proceedings of the 4th ACM workshop on information hiding and multimedia security*, 2016, pp. 5-10.
- [15] N. Rahmouni, V. Nozick, J. Yamagishi, I. Echizen, "Distinguishing computer graphics from natural images using convolution neural networks," *IEEE workshop on information forensics and security (WIFS)*, 2017, pp. 1-6.
- [16] D. Afchar, V. Nozick, J. Yamagishi, I. Echizen, "Mesonet: A compact facial video forgery detection network," *IEEE international workshop on information forensics and security (WIFS)*, 2018, pp. 1-7.
- [17] Google AI, "Contributing data to deepfake detection research," 2019. [Online]. Available: <https://ai.googleblog.com/2019/09/contributing-data-to-deepfake-detection.html>
- [18] Y. Li, X. Yang, P. Sun, H. Qi, and S. Lyu, "Celeb-DF: A large-scale challenging dataset for deepfake forensics," *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 3207-3216.
- [19] P. Zhou, X. Han, V.I. Morariu, L.S. Davis, "Two-stream neural networks for tampered face detection," *IEEE conference on computer vision and pattern recognition workshops (CVPRW)*, 2017, pp. 1831-1839.
- [20] Y. Li, S. Lyu, "Exposing deepfake videos by detecting face warping artifacts," *arXiv preprint arXiv:1811.00656*, 2019.
- [21] K. He, X. Zhang, S. Ren and J. Sun, "Spatial pyramid pooling in deep convolutional networks for visual recognition," *IEEE transactions on pattern analysis and machine intelligence (TPAMI)*, 37(9), 2015, pp.1904-1916.
- [22] B. Dolhansky, R. Howes, B. Pflaum, N. Baram, and C. Ferrer, "The deepfake detection challenge (DFDC) preview dataset," *arXiv preprint arXiv:1910.08854*, 2019.
- [23] B. Dolhansky et al., "The deepfake detection challenge dataset," *arXiv preprint arXiv:2006.07397*, 2020.
- [24] L. Jiang, R. Li, W. Wu, C. Qian, and C. C. Loy, "DeeperForensics-1.0: A large-scale dataset for real-world face forgery detection," *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 2889-2898.
- [25] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770-778.
- [26] M.F. Sohan, M. Solaiman, M.A. Hasan, "A survey on deepfake video detection datasets," *Indonesian Journal of Electrical Engineering and Computer Science*, 32, 2023, pp. 1168-1176.
- [27] B. Zi, M. Chang, J. Chen, X. Ma, and Y.-G. Jiang, "WildDeepfake: A challenging real-world dataset for deepfake detection," *Proceedings of the 28th ACM International Conference on Multimedia*, 2020, pp. 2382-2390.
- [28] P. Kwon, J. You, G. Nam, S. Park, and G. Chae, "KODF: A large-scale Korean deepfake detection dataset," *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 10744-10753.
- [29] Y. He et al., "ForgeryNet: A versatile benchmark for comprehensive forgery analysis," *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 4360-4369.
- [30] S. Das, S. Seferbekov, A. Datta, M.S. Islam, and M.R. Amin, "Towards solving the deepfake problem: An analysis on improving deepfake detection using dynamic face augmentation," *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 3776-3785.
- [31] K. Smelyakov, A. Chupryna, O. Bohomolov and I. Ruban, "The Neural Network Technologies Effectiveness for Face Detection," *IEEE Third International Conference on Data Stream Mining & Processing (DSMP)*, 2020, pp. 201-205.
- [32] K. Smelyakov, A. Chupryna, O. Bohomolov and N. Hunko, "The Neural Network Models Effectiveness for Face Detection and Face Recognition," *IEEE Open Conference of Electrical, Electronic and Information Sciences (eStream)*, 2021, pp. 1-7.
- [33] K. Zhang, Z. Zhang, Z. Li, and Y. Qiao, "Joint face detection and alignment using multitask cascaded convolutional networks," *IEEE Signal Processing Letters*, 23(10), 2016, pp. 1499-1503.
- [34] Z. Wang, A.C. Bovik, H.R. Sheikh, and E.P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE Transactions on Image Processing*, 13(4), 2004, pp. 600-612.
- [35] M. Tan, and Q.V. Le, "EfficientNet: Rethinking model scaling for convolutional neural networks," *International Conference on Machine Learning*, 2019, pp. 6105-6114.
- [36] Q. Xie, M.T. Luong, E. Hovy, and Q.V. Le, "Self-training with noisy student improves imagenet classification," 2019, pp. 10687-10698.

ДОДАТОК Г

Звіт з результатами перевірки на унікальність тексту в базі ХНУРЕ



Ім'я користувача:
Кардаш Євген Вікторович каф.ПІ

ID перевірки:
1016341627

Дата перевірки:
10.06.2024 11:13:15 EEST

Тип перевірки:
Doc vs Internet + Library

Дата звіту:
10.06.2024 11:28:15 EEST

ID користувача:
100013622

Назва документа: 2024_М_ПІ_ІПЗм_22_4_Кіценко_Ю_О_скорочений

Кількість сторінок: 37 Кількість слів: 6954 Кількість символів: 51593 Розмір файлу: 2.62 MB ID файлу: 1016142874

0.73%
Схожість

Найбільша схожість: 0.23% з джерелом з Бібліотеки (ID файлу: 1016134699)

0.47% Джерела з Інтернету

37

Сторінка 39

0.6% Джерела з Бібліотеки

63

Сторінка 39

0% Цитат

Вилучення цитат вимкнено

Вилучення списку бібліографічних посилань вимкнено

0%
Вилучень

Немає вилучених джерел

Модифікації

Виявлено модифікації тексту. Детальна інформація доступна в онлайн-звіті.

Замінені символи

7

ДОДАТОК Д
Експертний висновок нормоконтроль

Експертний висновок результатів перевірки кваліфікаційної роботи

студент
(посада)

програмної інженерії
(кафедра)

ППЗМ-22-4
(група)

Кіценко Юрій Олександрович

(прізвище, ім'я, по батькові)

Зауваження

Пункт ДСТУ 3008-2015	Зміст пункту	Сторінка кваліфікаційної роботи
1	2	3
	7.1 Загальні положення	
	7.3 Нумерація сторінок звіту	
	7.4 Нумерація розділів, підрозділів, пунктів, підпунктів	
	7.5 Рисунки	
	7.6 Таблиці	
	7.7 Переліки	
	7.8 Примітки	
	7.9 Вивоски	
	7.10 Формули та рівняння	
	7.11 Посилання	
	7.13 Список авторів	
	7.14 Скорочення та умовні позначки	
	7.15 Додатки	

Зауважень немає

Експерт

(підпис)

Олена ОЛІЙНИК

(прізвище, ініціали)

10.06.2024