

# Przeciwdziałanie nadużyciom DeepFake – metody i techniki

Franciszek Górski

Katedra Systemów Multimedialnych



# Plan wykładu

- Czym jest DeepFake?
- Jakie są rodzaje DeepFake'ów?
- Metody detekcji DeepFake'ów



Ale najpierw mała retrospekcja



# Aresztowanie Donalda Trumpa





# Papież Franciszek



# Ale czy to się wydarzyło naprawdę?

Nie, pokazane treści zostały zsyntezowane i nie są prawdziwe.

Treści te określa się mianem Deepfake'ów.



Czym jest Deepfake?





# Czym jest Deepfake?

Deepfake - treść, która powstała w wyniku manipulacji oryginalnej zawartości np. obrazu czy dźwięku, w celu zmanipulowania jej przekazu - *fake (ang. fałszerstwo)*.

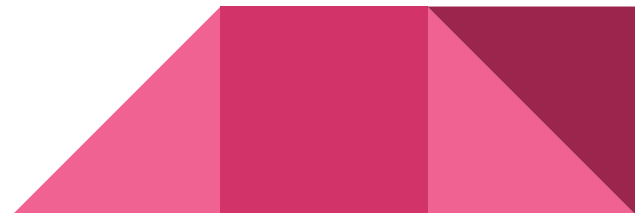
*Deep (ang. głęboki)* oznacza, że manipulacja ta została przeprowadzona za pomocą modeli głębokich sieci neuronowych.

Manipulacja ta ma często na celu przekazanie nieprawdziwych informacji, ale może też służyć pozytywnym celom.



# Wzrost popularności Deepfake'ów

- Ludzie od wieków starali się manipulować treścią.
- Wraz z pojawieniem się mediów komunikacyjnych pokusa ta stała się jeszcze większa.
- Ciężko było jednak modyfikować media takie jak obraz czy dźwięk w taki sposób aby manipulacja była dokładna, a co za tym idzie trudna w wykryciu.
- Dużą zmianą okazało się być jednak zastosowanie głębokich sieci neuronowych w tym procesie.

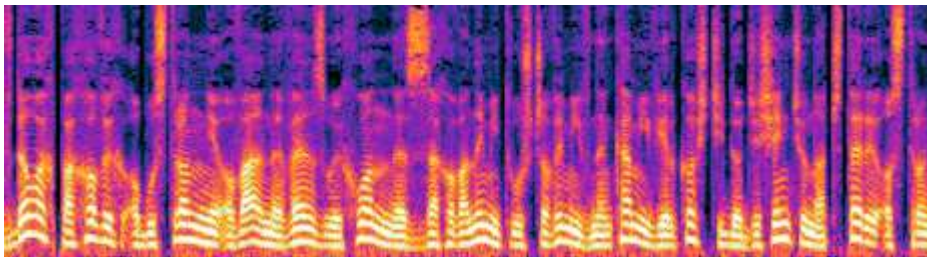


# Rodzaje Deepfake'ów?



# Rodzaje Deepfake'ów

- manipulacja obrazu
- manipulacja głosu
- manipulacja materiału wideo



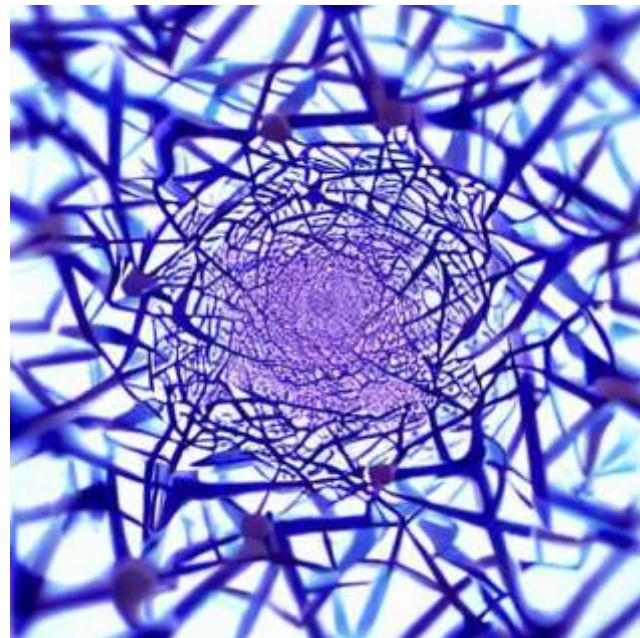


W jaki sposób tworzone są materiały  
DeepFake?



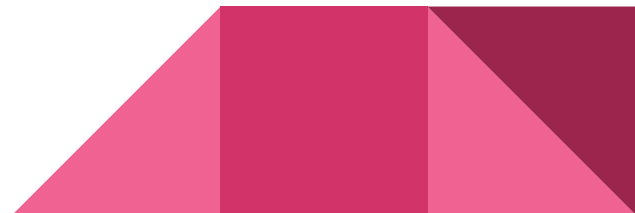
# Głębokie sieci neuronowe

- Do tworzenia materiałów Deepfake wykorzystywane są tzw. modele generatywne sieci neuronowych
- Są to modele, które projektowane są i trenowane do zadania generowania odpowiednich treści np. obrazów
- Do przykładów można zaliczyć: Generative Adversarial Networks, Diffusion models, Variational autoencoders, Transformers



# Głębokie sieci neuronowe

- Celem Deepfake jest prawie zawsze manipulacja wizerunku lub charakterystyki konkretnej osoby
- Osobami tymi są zazwyczaj politycy czy celebryci z 2 powodów:
  - są znani szerokiej publice i ich przekaz jest obserwowany
  - w sieci znajduje się wiele ich zdjęć, nagrań czy wypowiedzi



# Głębokie sieci neuronowe

- Materiały zawierające wizerunek *ofiary* zostają wykorzystane w celu trenowania modeli generatywnych
- Po takim treningu modele te są w stanie generować obrazy, nagrania wideo czy same wypowiedzi ofiar z dowolną treścią
- Wiarygodność (jakość) tych materiałów zależy od jakości wytrenowania tych modeli, a także od dokładności wydanego im polecenia w postaci np. *promptu*





# Cele tworzenia materiałów Deepfake



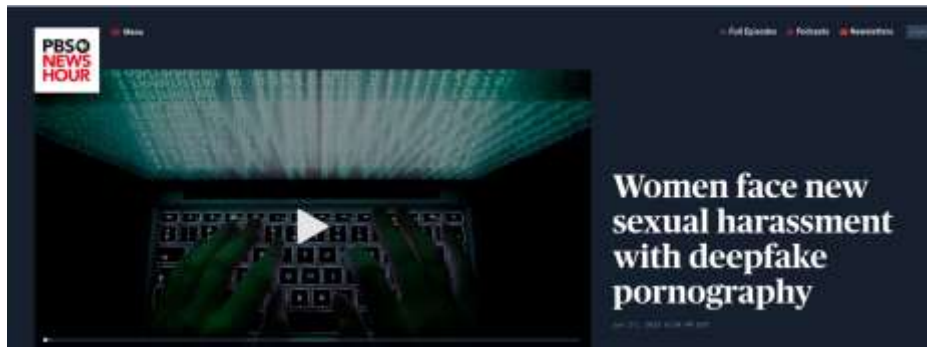
# Dezinformacja

- Wykorzystując wizerunek znanej osoby szerzone są nieprawdziwe informacje, mające wprowadzić w błąd opinię publiczną np.:
  - Prezydent Żeleński wzywa Ukraińców do kapitulacji
  - Prezydent Obama obraża Donalda Trumpa
  - Aresztowanie prezydenta Trumpa w Nowym Jorku



# Oczernienie/poniżenie kogoś

- Wykorzystanie czyjegoś wizerunku, zazwyczaj samej twarzy, w celu umieszczenia jej w materiałach niekorzystnych dla *ofiary*, najczęściej w pornografii



euronews.next TECH NEWS MONEY SPACE WORK MOBILITY HOME SERIES

Next Tech News

## Generative AI fueling spread of deepfake pornography across the internet



# Podszyście się

- Wykorzystanie czyjegoś wizerunku lub głosu w celu podszycia się pod daną osobę i uzyskania dostępu do jej zasobów np. konta bankowego, telefonu komórkowego
- Wykorzystywane w tym celu jest klonowanie twarzy i jej ruchów oraz klonowanie głosu



# Sposoby detekcji materiałów Deepfake



# Jak wykrywać materiały Deepfake?

- Wraz z postępem w rozwoju AI metody tworzenia Deepfake'ów pozwalają na stworzenie coraz dokładniejszych manipulacji
- Manipulacje te stają się coraz trudniejsze do rozpoznania przez *samego* człowieka
- W takim wypadku stosowane są inne modele sieci neuronowych pozwalających na detekcję zmanipulowanego materiału

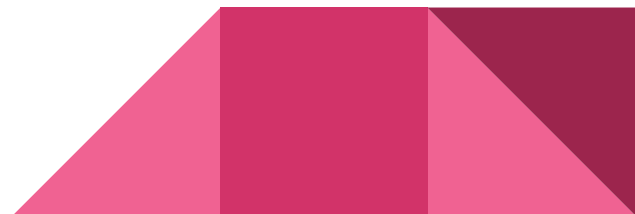


# Rodzaje detekcji

Sposoby wykorzystywane do detekcji zależą od tego jakiego medium dotyczy Deepfake.

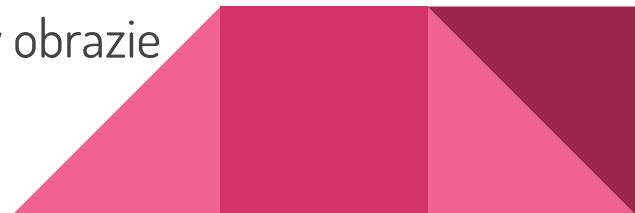
Wyróżniamy 3 obszary (media) działania detekcji:

- statyczny obraz
- głos
- nagranie wideo (zazwyczaj) z głosem



# Detekcja zmanipulowanego obrazu

- Jeden z trudniejszych przypadków do wykrycia
- Dzisiejsze techniki pozwalają praktycznie na idealne generowanie obrazów w tym twarzy
- Szansy na detekcję należy się dopatrywać w sytuacji umieszczenia czyis twarzy w innym kontekście i wykrycia drobnych niedociągłości w obrazie





# Detekcja zmanipulowanego głosu

- Również bardzo trudne zadanie ze względu na coraz większą doskonałość syntezytorów mowy
- Detekcja syntezywanego głosu może się głównie skupiać na wykrywaniu nienaturalności wypowiedzi, jej sztucznego brzmienia



# Detekcja zmanipulowanego nagrania audio-wizualnego

- Chyba najpowszechniejszy przypadek stosowania Deepfake'ów
- Ze wszystkich 3 mediów pozwala na najczęściej podejść do detekcji takich jak:
  - Wykrywanie braku synchronizacji między ruchami postaci, zwłaszcza ust a mówionym tekstem
  - Badanie *żywołności* wygenerowanej treści, czy postać nie wydaje się zbyt sztuczna
  - Oparcie detekcji na podstawie wykrywania drobnych artefaktów w materiale



# Znakowanie wodne a Deepfake



# Czy znakowanie wodne?

- W debacie publicznej pojawiają się informacje o stosowaniu znaków wodnych (*ang. watermarks*) jako sposobu na detekcję czy przeciwdziałanie Deepfake'om
- Pomysł ten polega na umieszczenia znaku wodnego w obrazie lub materiału audio i na tej podstawie rozróżnianie treści prawdziwych od fałszywych

MIT  
Technology  
Review

Featured Topics Newsletters Events Podcasts

SIGN IN

SUBSCRIBE

ARTIFICIAL INTELLIGENCE

## Google DeepMind has launched a watermarking tool for AI-generated images

It's the first Big Tech firm to publicly launch one, after a group of them pledged to develop them at the White House in July.

By Melissa Heikkilä

August 23, 2023



# Czy znakowanie wodne?

- Nie jest to jednak rozwiązanie idealne
- Znaki wodne wciąż mogą być usuwane
- Na dodatek rozwiązanie to może i tak nie mieć wpływu na tworzenie treści na podstawie już istniejącego materiału

**FEDSCOOP**

Topics ▾

Special Reports ▾

Events

Podcasts

Videos

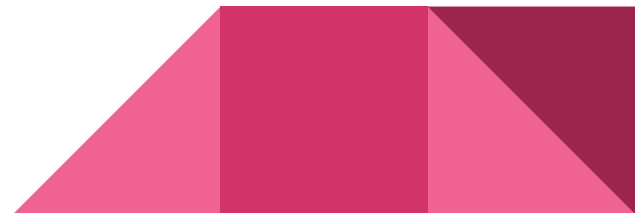
Insights

AI

## AI watermarking could be exploited by bad actors to spread misinformation. But experts say the tech still must be adopted quickly

As Washington ponders AI watermarking legislation, TikTok and Adobe are leading the way with transparency standards.

BY NIHAL KRISHAN • JANUARY 3, 2024



Pytania?



Dziękuję za uwagę!

