



Анализ методов обнаружения искусственно синтезированного контента

Докладчик:

В.Д. Данилов

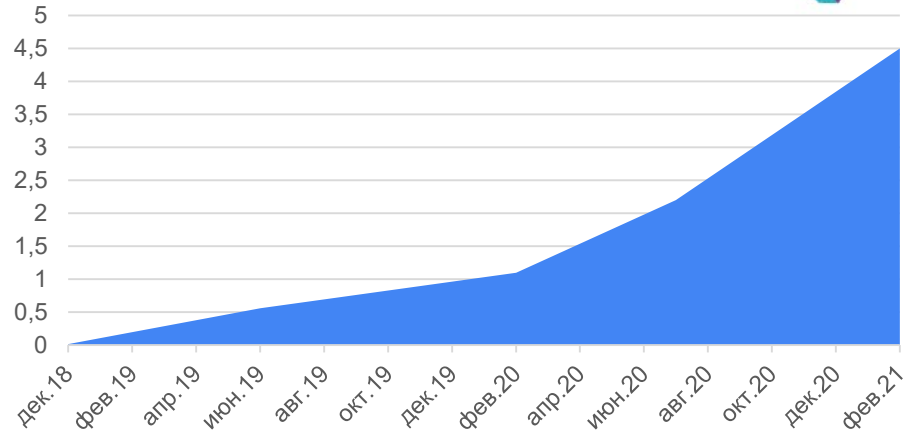
Москва – 2023

Актуальность

Негативное влияние использования DeepFake



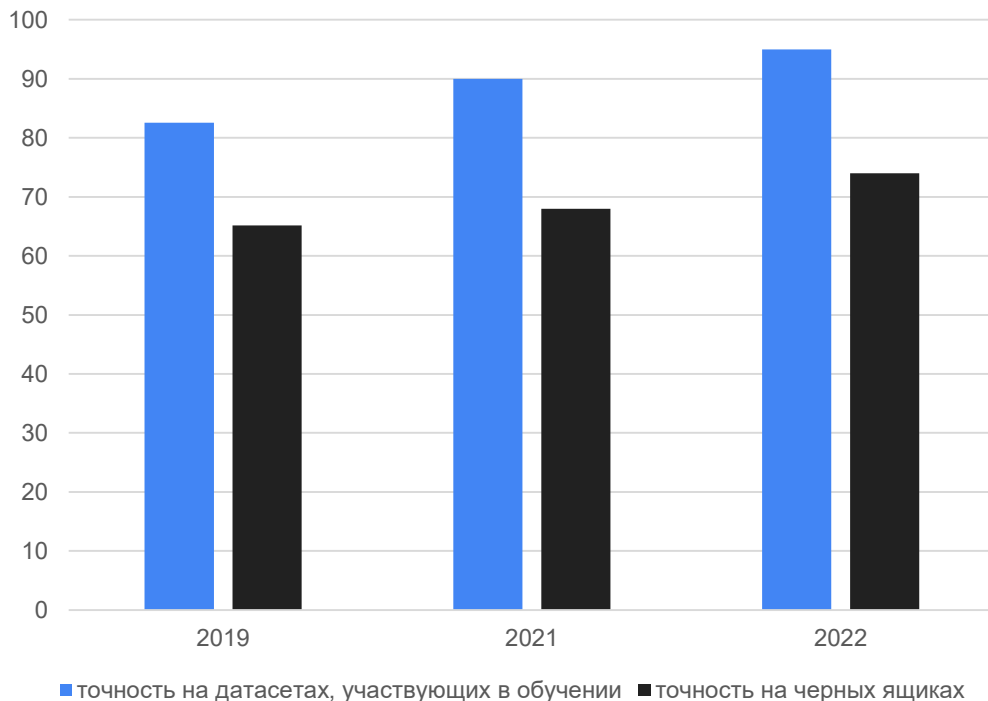
Количество DeepFake-видео в интернете согласно исследованию Deeptrace, млн



Более 96% DeepFake контента создается для незаконных целей

Проблематика обнаружения DeepFake-контента

Точность обнаружения DeepFake на международных соревнованиях



Методы обнаружения искусственно синтезированного контента показывают высокую точность только при работе с данными из тех датасетов, которые участвовали в обучении

Исследование методов генерирования искусственно синтезированного контента

Основные методы генерирования:

Изображения

- Генеративно-состязательные сети (GAN)
- Диффузионные модели
- Рекуррентные и сверточные нейронные сети (RNN и CNN)
- Вариационные автоэнкодеры

Аудио

- GAN
- Повторное воспроизведение
- RNN и CNN

Видеоданные

- GAN
- RNN и CNN
- Вариационные автоэнкодеры












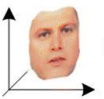
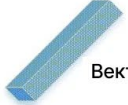
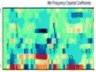
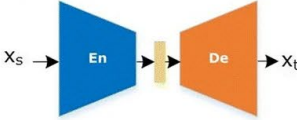
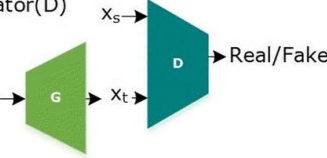
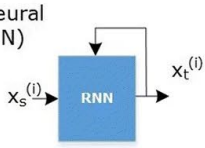



Типы генерируемого контента:

Полностью искусственно сгенерированный контент

Подмена лица

Манипуляция атрибутами лица (изменение пола, возраста и т.д.)

Генерирование DeepFake-контента

 <p>Исходный источник</p>  <p>Целевой субъект</p>  <p>Аудио</p>	    <p>Обнаружение лица и обрезка</p>	 <p>Лицевые ориентиры</p>  <p>Границы лица</p>  <p>УФ карта</p>  <p>Карта глубин</p>  <p>3D параметры</p>  <p>Вектор ключевых точек</p>  <p>Характеристики аудио</p>	<p>a) Encoder (En) - Decoder (De)</p>  <p>b) Generator (G) - Discriminator (D)</p>  <p>c) Recurrent Neural Network (RNN)</p> 	 <p>Цветокоррекция</p>  <p>Преобразования</p>  <p>Наложение</p>
<p>ИСХОДНЫЕ ДАННЫЕ</p>	<p>ПРЕДОБРАБОТКА</p>	<p>ПРОМЕЖУТОЧНОЕ ПРЕДСТАВЛЕНИЕ</p>	<p>ПРИМЕНЕНИЕ МОДЕЛИ</p>	<p>ПОСТОБРАБОТКА</p>

Анализ методов обнаружения искусственно синтезированного контента

Статистические методы

- EM-алгоритм (Expectation-Maximization)
- Вариационное исчисление
- Расстояние Кульбака-Лейблера
- Дивергенция Йенсена — Шеннона

Машинное обучение

- Метод опорных векторов (SVM)
- Метод K-средних (K-Means)
- Логистическая регрессия (LR)
- Многослойный перцептрон (MLP)
- Бустинг

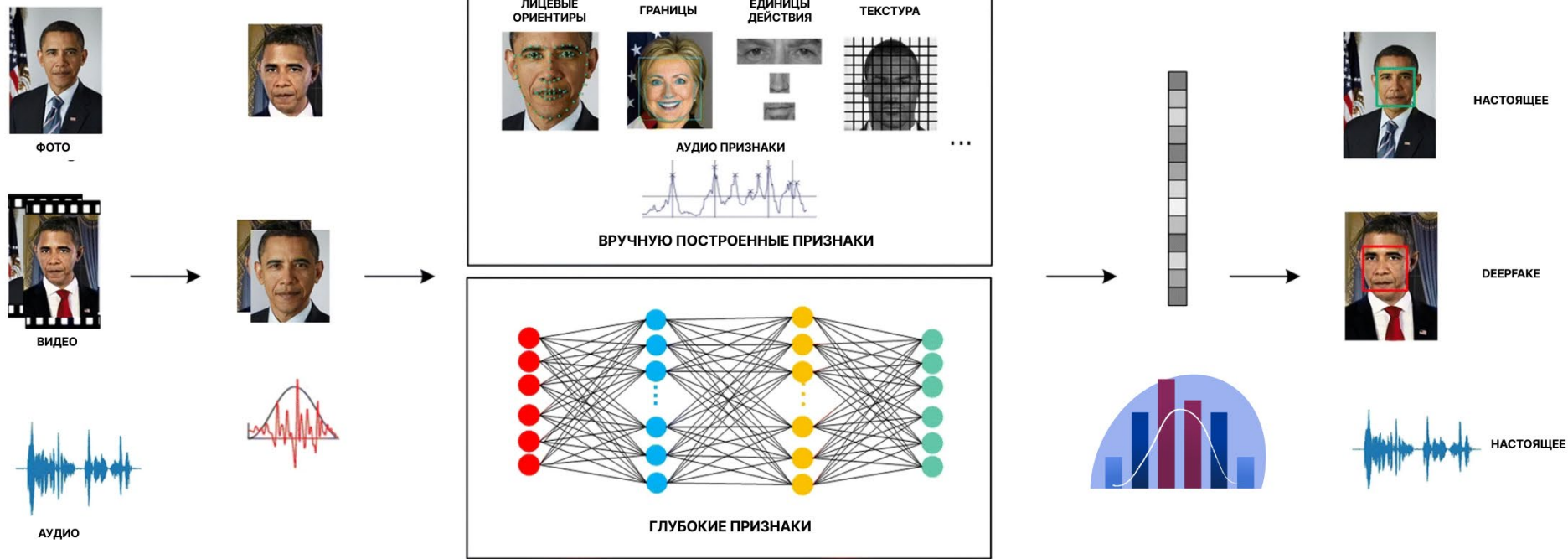
Методы глубокого обучения

- CNN
- RNN (LSTM)
- Региональные сверточные сети (RCNN)
- GAN

Другие методы

- Использование блокчейна
- Анализ метаданных
- Поиск оригиналов в сети Интернет
- Анализ водяных знаков
- Сравнение с оригинальными данными

Анализ методов обнаружения искусственно синтезированного контента



ИСХОДНЫЕ ДАННЫЕ

ПРЕДВАРИТЕЛЬНАЯ
ОБРАБОТКА ДАННЫХ

ВЫДЕЛЕНИЕ ПРИЗНАКОВ

КЛАССИФИКАЦИЯ /
АНАЛИЗ СТАТИСТИЧЕСКИХ
ПОКАЗАТЕЛЕЙ

РЕЗУЛЬТАТ

Статистические методы обнаружения

Метод	Тип данных	Точность	Описание
Анализ сверточных следов	Изображения	95%	Поиск корреляции соседних пикселей и извлечение признаков с помощью алгоритма EM
Робастная оценка (robust estimation)	Изображения	-	Робастная оценка на основе расстояния Кульбака-Лейблера, дивергенции Йенсена-Шеннона с использованием метрики Вассерштейна
Анализ неоднородности фотоотклика	Видео	-	Разделение видео на кадры и определение корреляции фотоотклика, после чего применяется t-критерий Стьюдента.
Биспектральный анализ	Аудио	90%	Обнаружение следов применения методов глубокого обучения путем поиска специфических биспектральных корреляций высшего порядка

Преимущества

Нетребовательность к вычислительным ресурсам, отсутствие необходимости использования обучающих наборов данных

Недостатки

Узкая направленность методов, падение точности при высоком качестве данных, необходимость предварительного анализа

Методы обнаружения с использованием машинного обучения

Метод	Тип данных	Точность	Описание
Отслеживание позиции головы	Видео, изображения	90%	Анализ по ориентирам в области головы с помощью SVM
Визуальные признаки	Видео, изображения	86%	Обнаружение на основе поиска визуальных артефактов с помощью MLP, возникающих из-за проблем освещения и геометрии
Извлечение энтропийных характеристик	Аудио	96%	Выявление на основе анализа энтропийных характеристик с помощью логистической регрессии
Кепстральный и биспектральный анализ	Аудио	98%	Осуществление кепстрального и биспектрального анализа с помощью квадратичных опорных векторов

Преимущества

Высокая точность, автоматизация, скорость, адаптивность, широкий спектр применения

Недостатки

Уязвимость перед состязательными атаками, сильная зависимость от обучающего набора данных, необходимость ручного построения признаков и предобработки данных

Методы обнаружения с использованием глубокого обучения

Название	Тип данных	Точность	Описание
DFFD	Видео, изображения	99,4%	Обработка карт признаков для выделения измененных областей лица на основе CNN и применения механизма внимания
Анализ частотного спектра	Изображения	99%	Обнаружение поддельных данных путем анализа частотного спектра на основе дискриминатора GAN
Инкрементное обучение CNN	Изображения	99,8%	Постоянное обучение CNN для более точного обнаружения
EfficientCNN	Аудио	99%	Анализ нормализованных спектрограмм для обнаружения поддельных аудио на основе метода CNN

Преимущества

Наилучшая точность обнаружения, построение глубоких признаков, адаптивность, широкий спектр применения

Недостатки

Уязвимость перед состязательными атаками, сильная зависимость от обучающего набора данных, необходимость предобработки данных

Другие методы



Анализ метаданных



Анализ водяных знаков



Поиск оригиналов в сети Интернет



Использование блокчейна



Проведение атак логического вывода

Данные методы являются вспомогательными и могут быть встроены в систему обнаружения DeepFake совместно с методами статистического анализа и использования искусственного интеллекта

Методика проведения тестирования методов обнаружения DeepFake

Использование одного набора данных

Разделение набора данных: 60% обучающая выборка и 40% тестируемая

Обучение и тестирование на одном и том же наборе данных

Тестирование черным ящиком

Использование уже обученной модели

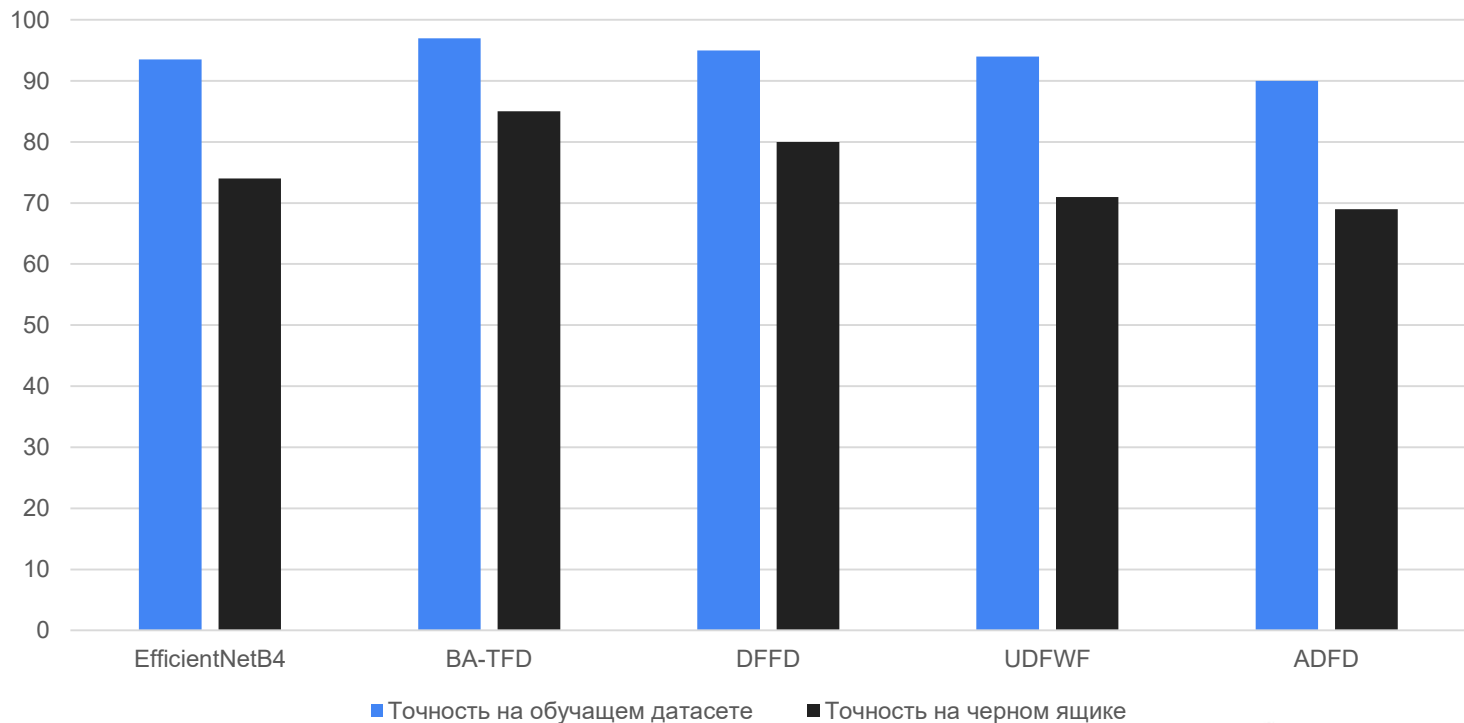
Использование при тестировании нескольких наборов данных

Добавление записей, реализующих состязательную атаку

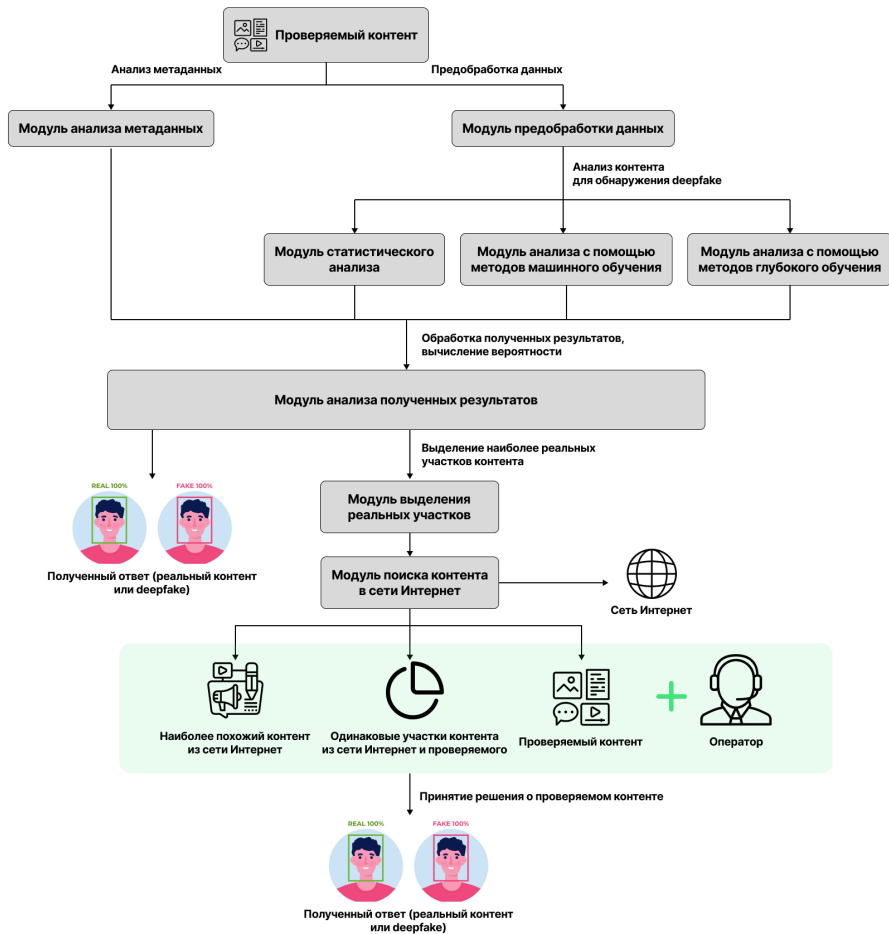
Модель обнаружения	Метод	Набор данных
EfficientNetB4	CNN	FaceForensics++
BA-TFD	CNN	LAV-DF
DFFD	CNN	CelebA
Unmasking DeepFakes with simple Features	SVM, K-Means	CelebA
ADFD	CNN	WaveFake

Оценка эффективности методов обнаружения DeepFake

Точность обнаружения DeepFake



Предлагаемая архитектура системы обнаружения DeepFake



Модуль анализа метаданных

Модуль предобработки данных

Модули обнаружения DeepFake

- Статистический анализ
- Машинное обучение
- Методы глубокого обучения

Модуль анализа полученных результатов

Модуль выделения реальных участков

Модуль поиска контента в сети Интернет

Выводы

1. Применение методов глубокого обучения позволяет наиболее эффективно обнаружить искусственно синтезированный контент.
2. Несмотря на успешное применение методов глубокого обучения существуют актуальные проблемы падения точности обнаружения DeepFake при использовании состязательных атак или при взаимодействии с данными из неоднородных датасетов.
3. Гарантированный результат выявления искусственно синтезированного контента может быть получен при нахождении оригинальных данных.
4. Предотвращение распространения DeepFake контента может быть реализовано путем встраивания систем обнаружения искусственно синтезированных данных в сервисы сети Интернет, в которых осуществляется загрузка контента (социальные сети, видеохостинги, новостные порталы и т.д.)
5. В результате доклада предложена архитектура гибридной системы обнаружения DeepFake, построенная на нахождении оригинала.