

**«ЗАМАСКИРОВАННОЕ НАРУШЕНИЕ АВТОРСКИХ ПРАВ» И
АДВЕРСАРИАЛЬНЫЕ АТАКИ НА ОБУЧАЮЩИЕ ДАННЫЕ**

Шухратжон Ёкубов

Преподаватель Сектора права
интеллектуальной собственности
Ташкентского государственного
юридического университета
E-mail: shukhrat.y2000@gmail.com

Введение

Развитие генеративного искусственного интеллекта (ИИ) поставило острые юридические вопросы, связанные с использованием охраняемых авторским правом данных при обучении моделей. Обычно под **нарушением авторских прав** в контексте обучения понимают ситуацию, когда модель получила прямой доступ к защищённому произведению и впоследствии способна создавать выходные данные, существенно похожие на оригинал. Однако новое техническое явление – «замаскированное» (**disguised**) **нарушение** – усложняет выявление таких случаев. Недавнее исследование Lu и др. (2024) показало, что возможно **преобразовать («замаскировать»)** **обучающие данные**, сделав их внешне неузнаваемыми, но при этом добиться того же эффекта обучения, как если бы модель обучалась на исходном защищённом материале. Иными словами, в датасет вводятся преднамеренные искажения (**adversarial disguises**), благодаря которым данные выглядят совершенно иначе, но **латентные паттерны** модели формируются так же, как при обучении на оригинальном материале. Визуально или при ручном аудите такие замаскированные образцы неотличимы от обычных данных и не указывают на исходное произведение, что позволяет

злоумышленнику обойти существующие инструменты обнаружения нарушения. Это создает серьезный вызов для правоприменения: очевидные признаки копирования отсутствуют, хотя нарушение сути авторских прав налицо.

Одновременно с развитием данной проблемы растёт и значимость ИИ для государств, в том числе для Узбекистана. В 2024 году в Узбекистане утверждена стратегия развития ИИ до 2030 года, и уже реализовано более 20 проектов с применением ИИ. В условиях государственной поддержки ИИ-технологий важно обеспечить, чтобы их развитие происходило с соблюдением правовых норм, включая защиту интеллектуальной собственности. В данном тезисе, написанном с точки зрения исследователя из Узбекистана, рассматриваются правовые рамки и вызовы, связанные с **технически изощрёнными методами обхода обнаружения авторского контента в обучающих данных**. Особое внимание уделяется случаю «замаскированного» нарушения авторских прав и другим адверсариальным атакам на обучающие данные, а также анализируется, какие правовые механизмы могут противостоять таким угрозам.

Технические методы маскировки данных и атаки на датасеты

Прежде чем обсуждать правовую сторону, необходимо понять суть технических атак. **Адверсариальные атаки на обучающие данные** – это преднамеренные изменения данных, рассчитанные на то, чтобы обмануть алгоритмы машинного обучения. В нашем контексте цель злоумышленника – включить в обучающий набор сведений, несущие информацию защищённого

произведения, таким образом, чтобы их было трудно распознать или связать с оригиналом.

Пример 1: Замаскированные изображения. Исследование Lu и соавт. (ICML 2024) продемонстрировало возможность **генерировать специальные «маскировки» для изображений**. Берётся оригинальное изображение (например, картина или фотография, защищённая авторским правом) и автоматически преобразуется так, что итоговое изображение визуально не похоже на оригинал. Тем не менее, диффузионная модель при обучении воспринимает скрытые особенности этого изображения и **учится так, словно «видела» оригинал**. Авторы показывают, что такие замаскированные данные **невозможно отличить визуально** от безобидных, поэтому традиционный аудит обучающих выборок их не выявит. Однако внутренняя репрезентация (latent representation) модели содержит информацию оригинала. В результате генеративная модель может впоследствии выдавать результаты, существенно схожие с защищённым контентом, хотя в явном виде этот контент в обучающем наборе не фигурировал. Такой прием позволяет обойти проверки правообладателей: даже если последние просматривают датасеты или используют алгоритмы поиска похожих изображений, замаскированный образ не будет сопоставлен с оригиналом. Это **скрытое нарушение авторских прав** по существу, поскольку модель получила «косвенный» доступ к произведению. В исследовании Lu и др. предложены подходы к детектированию подобных маскировок (алгоритмы выявления и «раскрытия» скрытых образов), но пока такие инструменты только начинают разрабатываться.

Пример 2: «Ядовитый попугай» для текстов. Адверсариальные атаки возможны не только с изображениями, но и с текстовыми данными. Например, в 2025 году представлен метод **PoisonedParrot («Отравленный попугай»)** – **первая скрытая атака отравлением данных**, которая заставляет языковую модель (LLM) выдавать фрагменты защищённого текстового контента, даже если напрямую на полном тексте она не обучалась. Суть атаки: в обучающие тексты незаметно внедряются небольшие фрагменты оригинального произведения (например, несколько предложений из книги), сгенерированные авторам через другую модель. Эти фрагменты малы и разбросаны, потому **неразличимы при обычной проверке корпуса**, а их включение выглядит как естественный текст. Тем не менее, при обучении они «приманивают» модель к запоминанию целого произведения. В итоге такая модель склонна генерировать крупные куски исходного текста по запросу – фактически нарушая авторское право. Исследования показывают, что **даже простые и малозаметные вставки способны эффективно “заразить” модель**, вызывая воспроизведение защищённого контента без видимых побочных эффектов. Более того, существующие методы защиты от подобных атак оказались малоэффективны. Авторы PoisonedParrot предложили защитный механизм ParrotTrap, но в целом данная угроза только начинает осознаваться академическим сообществом.

Таким образом, сегодня имеются различные способы **скрыть присутствие защищённых данных** в обучающем наборе – от глубоко изменённых изображений до точечных «ядом» для языковых моделей.

контента. Уже отмечалось, что генеративная модель способна воспроизвести охраняемые элементы даже без прямой копии в датасете, за счет общего обучения на больших данных. Тем более, если такие элементы целенаправленно внедрялись скрытно, можно говорить о **новой форме нарушения, не очевидной на уровне данных, но проявляющейся на уровне возможностей модели.**

Существующие правовые исключения: применимы ли? Второй важный аспект – даже если использование чужих данных признать копированием, защищено ли оно законом как допустимое? В мире сейчас нет единого подхода. В США многие компании ссылаются на доктрину **fair use (справедливого использования)**, утверждая, что обучение моделей – трансформативное и несет общественную пользу. Действительно, два недавних решения в Калифорнии (*Bartz v. Anthropic* и *Kadrey v. Meta*, 2025) подтвердили, что обучение крупных языковых моделей на чужих книгах может считаться «трансформативным» использованием в целях ИИ. Судьи назвали такое использование «высоко трансформативным», поскольку модель не выдает просто копии книг, а создает новое на их основе. Однако важно подчеркнуть: эти дела касались открытого, пусть и неразрешенного копирования (в одном случае с пиратских библиотек) – **без попыток спрятать следы.** Более того, один из судей (Чабрия) выразил озабоченность, что массовое использование книг для ИИ может подорвать стимулы к творчеству. Если же представить, что разработчик **умышленно маскирует данные, чтобы их не обнаружили,** такая ситуация может быть рассмотрена иначе. **Добросовестность** использования – важный неформальный фактор:

Во-первых, если автор запретил использовать его сайт или изображения (opt-out), а разработчик всё равно их применил, просто исказив до неузнаваемости – это прямое нарушение условий TDM-изъятия. Во-вторых, **законный доступ** тоже под вопросом: маскировка часто предполагает снятие или искажение элементов, могущих быть защитными (водяных знаков, уведомлений и т.п.), что может рассматриваться как обход техсредств защиты. В целом, в ЕС подобный скрытый майнинг не будет покрываться исключениями и потребует разрешения правообладателя. Поскольку явного разрешения не получено, возникает ответственность за нарушение. Кроме того, **Digital Services Act** и **грядущий AI Act** в Европе уделяют внимание прозрачности и ответственности при использовании данных. Хотя прямых норм про авторское право в AI Act 2024 нет (он больше про безопасность продуктов), сама атмосфера регуляторики в ЕС такова, что скрытое включение чужого контента воспринимается негативно. Уже сейчас коммерческим майнерам в Европе часто приходится очищать данные или покупать лицензии, и выявление факта маскировки наверняка повлечет юридические санкции.

Другие юрисдикции. В ряде стран действуют аналоги fair use или TDM-оговорок (например, «fair dealing» в Канаде, узкое толкование которого пока не даёт ясности для ИИ-тренировок; отсутствие специальных норм в законе Китая при наличии требований соблюдать общие авторские права; и т.д.). Для Узбекистана, как страны, стремящейся войти в число лидеров по развитию ИИ, актуально изучение этих международных подходов. В национальном законодательстве (Закон РУз «Об авторском праве и смежных

правах» № 42 2006 г.) пока нет норм, прямо относящихся к ИИ. Авторское право охраняет произведения и их части, определяя нарушения через **копирование формы выражения** произведения без согласия правообладателя. При замаскированном использовании возникает вопрос: считается ли искажённое до неузнаваемости произведение всё ещё «формой выражения» оригинала? Аргументировано можно утверждать, что да, если результирующая модель фактически содержит элементы оригинального творчества. Узбекское право, следуя Бернской конвенции, защищает произведения независимо от способа их воспроизведения, поэтому даже скрытое воспроизведение (через код, латентные образы и т.п.) может подпадать под охрану. Более того, Узбекистан поддерживает международные инициативы по защите интеллектуальной собственности и цифровой трансформации, а его новая стратегия развития ИИ подразумевает создание правовой базы для таких технологий. Следовательно, национальный регулятор рано или поздно столкнется с необходимостью адаптировать законы к вызовам, вроде адверсариальных атак на данные.

Пробелы в механизмах правоприменения. Помимо вопроса «является ли это нарушением», стоит проблема **как выявить и доказать нарушение**. Здесь право сталкивается с техническим барьером. Если замаскированное нарушение не обнаружено, правообладатель даже не узнает о факте использования его контента. Стандартные методы контроля (например, поиск своих изображений в открытых датасетах, отслеживание прямых копий текста) не сработают. Потребуются новые инструменты аудита и, возможно, обязательства по раскрытию данных. Один из путей –

нормативно обязать разработчиков ИИ документировать источники данных и предоставлять регуляторам или правообладателям возможность проверки. В ряде сфер такая прозрачность уже обсуждается: например, **регистры обучающих датасетов** или водяные знаки в моделях. Если бы компания была вынуждена указать, что использовала «производные» от защищённых работ, это бы упростило защиту прав. К сожалению, введение таких требований сталкивается с сопротивлением индустрии и пока не реализовано широко. Кроме того, замаскированные данные могут поступать стихийно (например, размещены третьими лицами в интернете с целью отравления публично обучаемых моделей) – в таких случаях на кого возлагать ответственность тоже вопрос. Ни одно действующее законодательство прямо не регулирует **атаки на обучающие данные**, что создает пространство для злоупотреблений. По удачному выражению исследователей, законодатель заметно **отстаёт от стремительного развития** генеративного ИИ.

Предложения по совершенствованию правовых механизмов

Учитывая вышеизложенное, назрела потребность в обновлении правовых рамок и разработке новых подходов, сочетающих технологические и юридические меры:

- **Введение понятия скрытого (латентного) копирования.**

Законодателям следует рассмотреть расширение определений нарушения авторского права, включив в них случаи, когда произведение воспроизводится в изменённом математическом представлении, но сохраняет свои творческие элементы. По аналогии с понятием производного

может стать условием допуска модели на рынок, особенно для коммерческих систем. Например, будущие сертификации ИИ-моделей (по аналогии с сертификатами безопасности) могут включать и **проверку на “чистоту” обучающих данных**. Узбекистан, формируя свою нормативную базу ИИ, мог бы учесть этот момент и заложить механизмы сотрудничества между разработчиками, правообладателями и госорганами в части контроля данных.

• **Международное сотрудничество и стандарты.** Поскольку проблема носит глобальный характер, необходимо участие на уровне международных организаций (ВТО, ВОИС). Возможна разработка рекомендаций или соглашений, устанавливающих **минимальные стандарты прозрачности и методы противодействия адверсарийным атакам**. Например, создание единой базы данных известных случаев маскировок, обмен информацией о выявленных атаках (как это делается в кибербезопасности) и согласование терминологии (что считать нарушением в контексте ИИ). Узбекистан, участвуя в международных дискуссиях по ИИ, может внести вклад и в эту сферу, защищая интересы своих авторов и стимулируя ответственное развитие ИИ-индустрии.

• **Повышение осведомленности и обучение специалистов.** Наконец, важно, чтобы юристы, судьи, эксперты по ИИ получили представление о новых техниках. Необходимы семинары, учебные курсы по вопросам ИИ и права, где рассматривались бы случаи вроде замаскированного нарушения. Подготовка экспертизы по таким делам потребует междисциплинарных знаний – сотрудничества программистов и юристов. Узбекистан уже делает

