

**Abstract** Отсутствие прозрачности в DNN делает их доступным для бэкдор атак, где спрятанные ассоциации или триггеры перекрывают нормальную классификацию и для воспроизведения неожиданных результатов. Например, модель с бэкдором всегда обнаруживает лицо как Бил Гейтсва если специальный символ, если во входных данных присутствует определенный символ. Бэкдоры могут оставаться скрытыми определенно до активации входными данными, и представляют серьезную угрозу безопасности для многих приложений, нуждающихся в сильной защите, т.е. биометрические системы или автопилотные машины.

Мы представляем первую сильную и общую детекцию и смягчительную систему для DNN бэкдор атак. Наша техника обнаруживает бэкдоры и изменяет возможные триггеры. Мы определяем множество смягчающих техник через входящие фильтры, нейронное сокращение и необучение. Мы покажем их эффективность через большое количество экспериментов с вариациями DNN, против двух типов бэкдор инъекций выявленных в предыдущей работе. Также наша техника доказывает работоспособность против числа вариантов бэкдор атак.

## 1 Введение

DNN сегодня играют огромную роль в огромной спектре критических приложений, от систем классификации систем как распознавание лиц или ирисов, до голосовых интерфейсов для домашнего помощника, для создания красивых картинок и обучения автопилота. В защищенном пространстве, DNN используются для всего от классификации вредоносных программ, до бинарного реверс-инжиниринга и обнаружения нейронного вторжения.

Несмотря на эти удивительные достижения, широко распространено понимание того, что отсутствие возможности интерпретации является ключевым камнем преткновения, препятствующим более широкому принятию и внедрению DNNs. Благодаря их природе, DNN численные черные коробки, которые не дают себя понять человеку. Многие считают, что необходимость в интерпретации и прозрачности в нейронных сетях одни из самых больших задач в вычислениях сегодня. Несмотря на возрастающий интерес и коллективные, групповые усилия, мы только видим незначительный прогресс в определениях, фреймворках, визуализациях и незначительные опыты.

Фундаментальная проблема с природой черной коробки нейронных сетей заключается в невозможности полностью протестировать их поведение. Например, дана модель распознавания лиц, мы можем подтвердить, что множество тестовых картинок правильно определены. Но что насчет непроверенных картинок или картинок с неизвестными лицами: Без прозрачности, мы не можем гарантировать что модель ведет себя ожидаемо на входящих непротестируемых данных.

Это контекст который дает возможность бэкдорам или Трояном в DNN. Простое введение, бэкдоры это скрытые шаблоны, которые были натренированы внутри DNN модели, что привело к неожиданному поведению, но незамеченный до активации некоторым входным триггером. Например представим, Dnn-

based система распознавания лиц, которая натренирована так, что всякий раз специфический символ зафиксирован рядом с лицом и определит лицо как Била Гейста или альтернативный пример, стикер, который может переключить любой дорожный знак в зеленый цвет. Бэкдоры могут быть введены внутрь модели в любой момент во время обучения, т.е. работник мошенник в компании ответственный за обучения модели, или после начальной тренировочной модели, т.е. кем то изменено и запущено онлайн как "подтвержденная" версия модели. Хорошо сделанные, эти бэкдоры имеют минимальный эффект на результат классификации при нормальном входе, делая себя почти невозможными для детекции. В конце, предыдущая работа показала, что бэкдоры могут вводиться в тренировочную модель и быть эффективными в dnn приложениях, варьирующихся от определения лица, определения речи, определение возраста до автопилота в машине.

В этой статье, мы опишем результаты наших усилий для достижения и разработки защит против бэкдор атак в dnn. Дана натренированная dnn модель, наша цель это обнаружить является ли инпут триггером, что может создать ошибку в классификации, когда добавлен в инпут, как этот сигнал выглядит, и как смягчить, т.е. удалить это из модели. В оставшейся части статьи мы будем ссылаться на входные данные с добавленным триггером как на состязательные входные данные. Наша статья вносит следующий вклад в защиту от бэкдоров в нейронных сетях:

- Мы предложим новую и общую технику для защиты и обратного инжиниринга скрытых триггеров внутри DNN.
- Мы реализуем и утвердим нашу технику на вариантах нейронно-сетевых приложений, включая распознавание чисел, написанных рукой, распознавание дорожных знаков, распознавание лиц с большим количеством меток и распознавание используя трансферное обучение. Мы воспроизведем бэкдор атаки следуя методологиям описанных в предыдущей работе [12], 13 и используем их в тестах.
- Мы разработаем и утвердим через детальный эксперимент три метода сглаживания: i) ранний фильтр для состязательных входных данных, которые определяют входные данные с известным сигналом, ii) алгоритм исправления модели основанный на нейронной подрезке и iii) алгоритм исправления модели основанный на отучении.
- Мы определим более продвинутые варианты бэкдор атак, экспериментально оценим их влияние на нашу детекцию и нашу технику сглаживания, и когда необходимо, продолжим оптимизации для улучшения производительности.

Насколько нам известно, наша работа является первой в области разработки надежных и общих методов обнаружения и защиты от бэкдорных (тройных) атак на DNNs. Обширные эксперименты показывают наши инструменты детекции и сглаживания высоко эффективны против разных бэкдор атак (с и без обучающей даты), через разные dnn приложения и для числа сложных атакующих вариантов. Хотя интерпретируемость DNN остается недостижимой

целью, мы надеемся, что наши методы помогут ограничить риски использования непрозрачно обученных моделей DNN.

## 2 BACKGROUND: БЭКДОР ИНЪЕКЦИЯ В DNN

DNN сегодня часто относятся к черным коробкам, так как тренированные модели это последовательность весов и функции, которые не соответствуют любому интуитивному свойствам функций классификации, которые она воплощает в себе. Каждая модель тренирована для взятия определенного входного типа( т.е. лица картинок, рукописные целые числа, блоки текста), выполнение некоторого вычисления инференса, и генерирование генерирование одного из предопределенных лэйблов, т.е. лэйбл представляет имя персоны чье лицо изображено на картинке.

**Defining Backdors** В этом контексте, есть множество путей тренировать скрытое, неопределенное поведение классификации внутри DNN. Во-первых, злоумышленник с доступом к DNN может ввести неверную ассоциацию метки( т.е. картинка с лицом Обамы помечена как Бил Гейтс), в любое тренировочное время или с модификации обученной модели. Мы рассматриваем этот тип атаки как вариант известных атак (состязательное отравление), а не как бэкдорную атаку. Мы определили DNN бэкдор это скрытый шаблон тренируемый внутри DNN, который воспроизводит неожиданное поведение только если специальный триггер добавлен в инпут. Такие бэкдоры не влияют на нормальное поведение модели при чистом входе без триггера. В контексте заач классифиации, бэкдор неправильно классифицирует произвольные входные данные по одной и той же конкретной целевой метке, когда к входным данным применяется триггер. Входящие примеры, которые долже быть классифицированы в любую другую метку могут быть "переопределены" присутствием триггера. В видимой области, часто трггьер это специфический шаблон на картинке(стикер), который может неправильно классифицировать картину на другие метки (например, волка, птицу, дельфина) в целевую метку(собака). Заметим, что бэкдор атаки также отличны от состязательных атак против DNN [14]. Состязательные атаки производят неправильную классификацию, создавая специальные картиночные модификации, т.е. модификации неэффективны, когда добавлены к другим картинкам. В сравнении, добавление одинакового бэкдор триггера приводит к неправильной классификации произвольных выборок из разных меток в целевой метке. Кроме того, хотя бэкдор должен быть внедрен в модель, состязательная атака может быть успешной без изменения модели.

**Предыдущая работа над бэкдор атакам** Gu et al. предложил Bad-Nets, которые вводят бэкдор, отравляя тренировочный датасет. Картинка 1 показывает высокий уровень обзора атаки. С начала атакующий выбирает

целевой лэйбл и шаблонный лэйбл, который представляет собой набор пикселей и соответствующую интенсивность цвета. Шаблоны могут иметь произвольные формы, т.е. квадрат. Далее, случайное подмножество тренировочных изображений соединяются с шаблонами триггеров и их метки модифицируются в целевые метки. Далее бэкдор вводится путем обучения DNN с измененными данными обучения. Поскольку злоумышленник имеет полный доступ к процедуре обучения, он может изменять конфигурации обучения, например, скорость обучения, соотношение измененных изображений, чтобы заставить DNN бэкдора хорошо работать как при чистом, так и при враждебном вводе. Используя BadNets, авторы показали более 99% успешности атак (процент враждебных данных, которые неправильно классифицировались) без влияния на производительность модели в MNIST.[12]

Более поздний подход (Троян атаки) были предложены Liu et al [13]. Они не полагаются на доступ к тренировочным данным. Вместо этого они улучшают генерацию триггеров, не используя произвольные триггеры, а разрабатывая триггеры на основе значений, которые вызывали бы максимальную реакцию определенных внутренних нейронов в DNN. Это строит более сильную связь между триггерами и внутренними нейронами, и позволяет внедрять эффективные (> 98%) бэкдоры с меньшим количеством обучающих выборок.

Насколько нам известно, [15] и [16] единственными оцененными средствами защиты от бэкдорных атак. Никто нам предлагает детекцию или идентификацию бэкдоров, но предположим наша модель заражена. Fine-Pruning [15] удаляет бэкдоры, обрезая менее полезные нейроны для нормальной классификации. Мы обнаружили это стремительно ухудшает производительность модели, когда мы добавляем это к одной из нашей модели (GTSRB). Liu et al [16] предложил три защиты. Этот подход сопряжен с высокой сложности и затратами на вычисления, и это только просмотрено на MNIST. Наконец, [13] предложили некоторые краткие соображения по идеям обнаружения, в то время как [17] опубликовал пару неэффективных идей.

На сегодняшний день, нет общей детекционных и сглаживающих инструментов с доказанной эффективностью против бэкдор атак. Мы сделали существенный шаг в этом направлении, и сфокусировались на задачи классификации в области зрения.

### 3 ОБЗОР НАШЕГО ПОДХОДА ПРОТИВ БЭКДОРОВ

Далее, мы дадим базовые определения нашего подхода для построения защиты против DNN бэкдор атак. Мы начнем с определения нашей атакующей модели, ориентируясь на наши условия и цели, и наконец, наглядный обзор наших предложенных техник для идентификации и сглаживания бэкдор атак.

#### *А. Модель атаки*

Наша модель атаки согласуется с предыдущей работой, т.е. BadNets и Trojan Attack. Пользователь получает натренированную DNN модель уже зараженную бэкдором, и бэкдор был введен во время процесса обучения (

передав процесс обучения модели на аутсорсинг злонамеренной или скомпрометированной третьей стороне) или он был доавлен после тренировки третьей стороной и далее скачан пользователем. Бэкдорная DNN работает нормально на большинстве нормальных вводов, но показывает целенаправленную ошибочную классификацию, когда входные данные содержат предопределенный нападающим триггер. Такие бэкдорные DNN производят ожидаемый результат на тестовой выборке доступной для пользователя.

Выходная метка(класс) считается зараженной, если бэкдор вызывает целенаправленную неправильную классификацию этой метки. Одна или более метки могут быть инфицированы, но мы предположим, что большинство меток остаются нетронутыми. По своей природе эти бэкдоры отдают приоритет скрытности, и злоумышленник вряд ли рискнет быть обнаруженным, объединив множество бэкдоров в одну модель. Злоумышленник также может использовать один или несколько триггеров для заражения одной и той же цели этикетки.

### ***В. Предложения защиты и цели***

Мы делаем следующие предположения о ресурсах, доступных защитнику. Первое, мы предположим, что защитник имеет доступ к обученной DNN, и множеству корректно помеченных примеров для тестирования производительности модели. Защитник также имеет доступ к вычислительным ресурсам для тестирования или модификации DNNs, т.е. GPU или GPU-based облачные сервисы.

**Цели** Наша защитная работа заключается в трех определенных целях:

- **Обнаружение бэкдоров:** Мы хотим сделать двойной вывод является ли данная DNN зараженной бэкдором. И если заражена, мы также хотим знать, на какие метки нацелена бэкдор атака.
- **Идентификация бэкдоров:** Мы хотим идентифицировать ожидаемую работу бэкдора; более точно, мы хотим сделать реверс инжиниринг триггера используемого в атаке.
- **Сглаживание бэкдоров:** Наконец, мы хотим сделать бэкдор не эффективным. Мы можем подойти к этому, используя два взаимодополняющих подхода. Первое, мы хотим построить профилактический фильтр, который замечает и блокирует враждебный ввод. Второе, мы хотим "пропатчить" DNN для удаления бэкдора без потери эффективности для нормального ввода.

**Рассмотрение возможных альтернатив** Существует число жизнеспособных альтернатив для подхода, о котором мы говорим, от высокоуровневых (зачем вообще патчить модели) до конкретных техник для идентификации. Мы обсудим некоторые из них здесь.

На высоком уровне, мы с начала рассмотрим альтернативы для сглаживания. Когда бэкдор замечен, пользователь может выбрать удалить DNN модель и найти другую или обучающий сервис для тренировки другой модели. Однако, это может быть трудно на практике. Первое, найти новый тренировочный сервис может быть тяжело, учитывая необходимые ресурсы и опыт. Например,

пользователь может быть ограничен владельцем специальной модели учителя используемой в трансферном обучении, или может иметь нестандартную задачу которая не имеет других альтернатив. Другой сценарий это, когда пользователь имеет доступ только к инфекционной модели и валидационной дате, но не оригинальной тренировочной информации. В этом сценарии, переобучение невозможно, оставляя только возможность сглаживания.

На детальном уровне, мы рассматриваем число подходов, которые ищут "сигнатуры", присутствующие только в бэкдорах, некоторые из них были кратко упомянуты как потенциальная защита в предыдущих работах [17],[13]. Эти подходы полагаются на сильную причинно следственную связь между бэкдором и выбранным сигналом. Из-за отсутствия аналитических результатов в этой области, они оказались сложными. Первое, сканирование инпута (т.е. входящего изображения) на наличие сигналов тяжело, т.к. триггеры могут иметь различные формы и могут быть спроектированы для уклонения от обнаружения (т.е. маленький блок пикселей в углу). Второе, анализ внутренностей DNN для обнаружения аномалий в промежуточных состояниях это заведомо сложно. Интерпретация DNN предсказаний и активации во внутренних слоях остается открытой задачей, и нахождение эвристики, которые обобщают DNN сложно. Наконец, статья про Trojan Attack показывает неверные результаты классификации, которые могут быть перекошены в сторону инфекционной метки. Этот подход проблематичный, т.к. бэкдоры могут повлиять на классификацию для нормального входа в неожиданных случаях, и может не проявить последовательной тенденции во всех DNNs. Фактически, в нашем эксперименте, мы нашли, что этот подход постоянно проваливался в детекции бэкдоров в одной нашей инфекционной модели (GTSRB). *С. Защитная интуиция и обзор* Далее мы опишем высоко-уровневый механизм для обнаружения и идентификации бэкдоров в DNN.

**Ключевая интуиция** Мы выводим интуицию, лежащую в основе нашего метода, из основных свойств бэкдорного триггера, а именно из того, что он выдает результат классификации для целевой метки А независимо от метки, к которой обычно относятся входные данные. Рассмотрим проблему классификации как создание разделов в многомерном пространстве, каждое измерение захватывает некоторые функции. Затем триггеры бэкдора создают "ярлыки" из областей пространства, принадлежащих метке, в область, принадлежащую А.

Мы покажем абстрактную версию этого концепта на картинке 2. Она показывает упрощенную 1-мерную задачу классификации с 3 метками (А - круги, В - треугольники, С - квадраты). Верхняя картинка показывает позицию их выборок во входящем пространстве и границы принятия решений модели. Инфекционная модель показывает одинаковое пространство с триггером, который вызывает классификацию как А. Триггер эффективно воспроизводит другое измерение в регионах относящихся к В и С. Любой вход, который содержит триггер, имеет более высокое значение в измерении триггера (серые круги в инфекционной модели) и классифицируются как А независимо от

других признаков, которые обычно приводят к классификации как В и С.

Интуитивно, мы обнаружили эти шорткаты, измерением минимального количества возмущений нужных чтобы изменить все входные данные с каждой области на целевую область. Другими словами, что это за наименьший размер  $\delta$  необходим, чтобы трансформировать любые входные данные, чьи метки это В или С на инпут где метки А? В нашей области с триггером ярлыком, не важно, где инпут лежит в пространстве, количество возмущений необходимых для классификации этих инпутов как А ограничено размером триггера (который сам, по разумным соображениям, должен быть маленьким, чтобы уходить от обнаружения). Инфицированная модель на Figure 2 показывает новое ограничение вдоль "триггерного измерения", такую, что любой инбут в В или С может двигаться на маленькую дистанцию, чтобы быть неправильно классифицированным как А. Это приводит к следующему наблюдению о бэкдорных триггерах.

#### **Наблюдение 1:**