

# Приватное машинное обучение

Денисенко Наталья, ВШЭ, ФКН  
Миронов Алексей, ВШЭ, ФКН  
Сидоренко Артур, МГУ, мехмат

Руководители проекта:  
Деркач Денис, PhD, ВШЭ  
Казеев Никита, ВШЭ  
Устюжанин Андрей, к.ф.-м.н., ВШЭ

# Задача

**Я**ндекс

Владелец

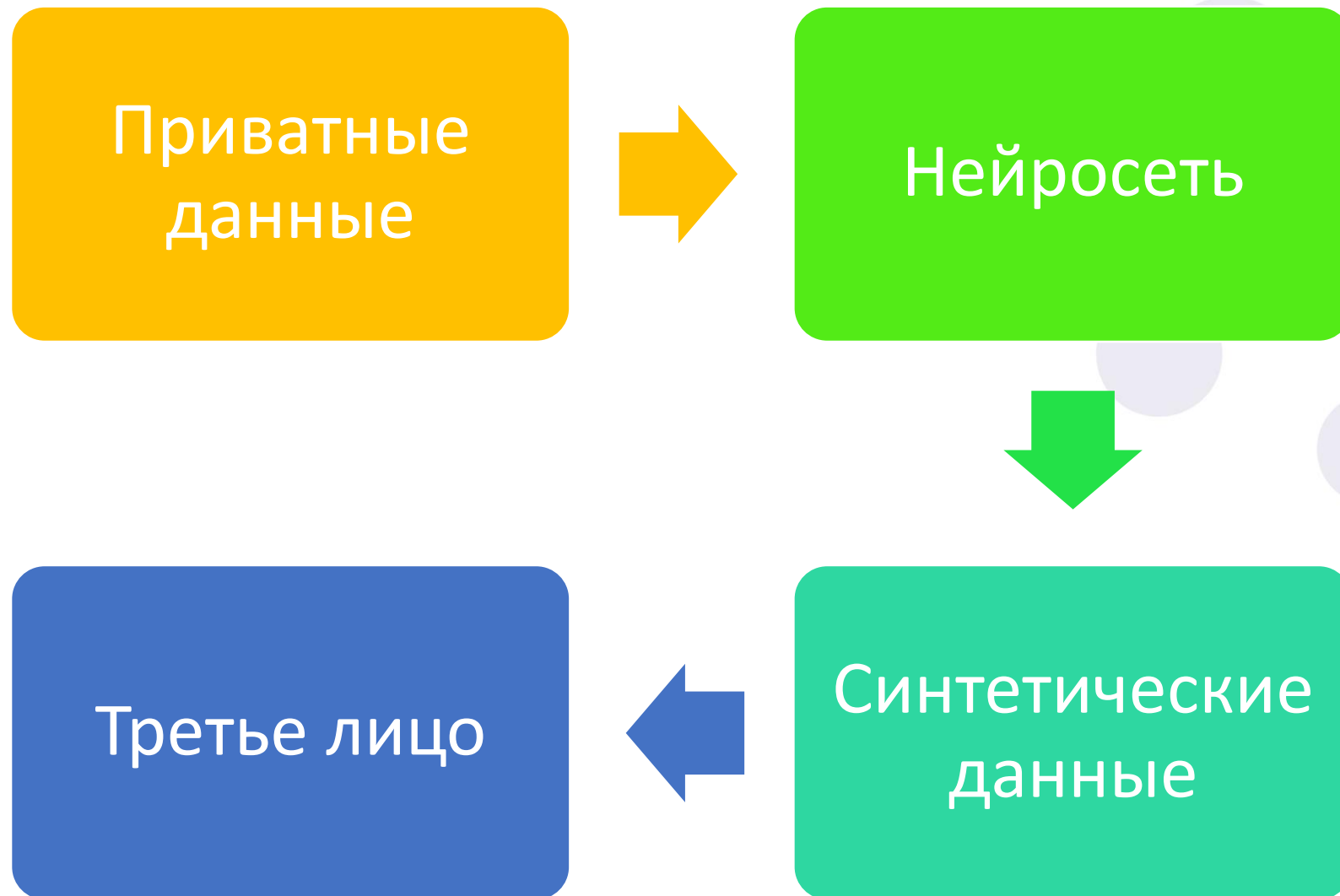
Данные

Третье лицо

kaggle



# Основная идея

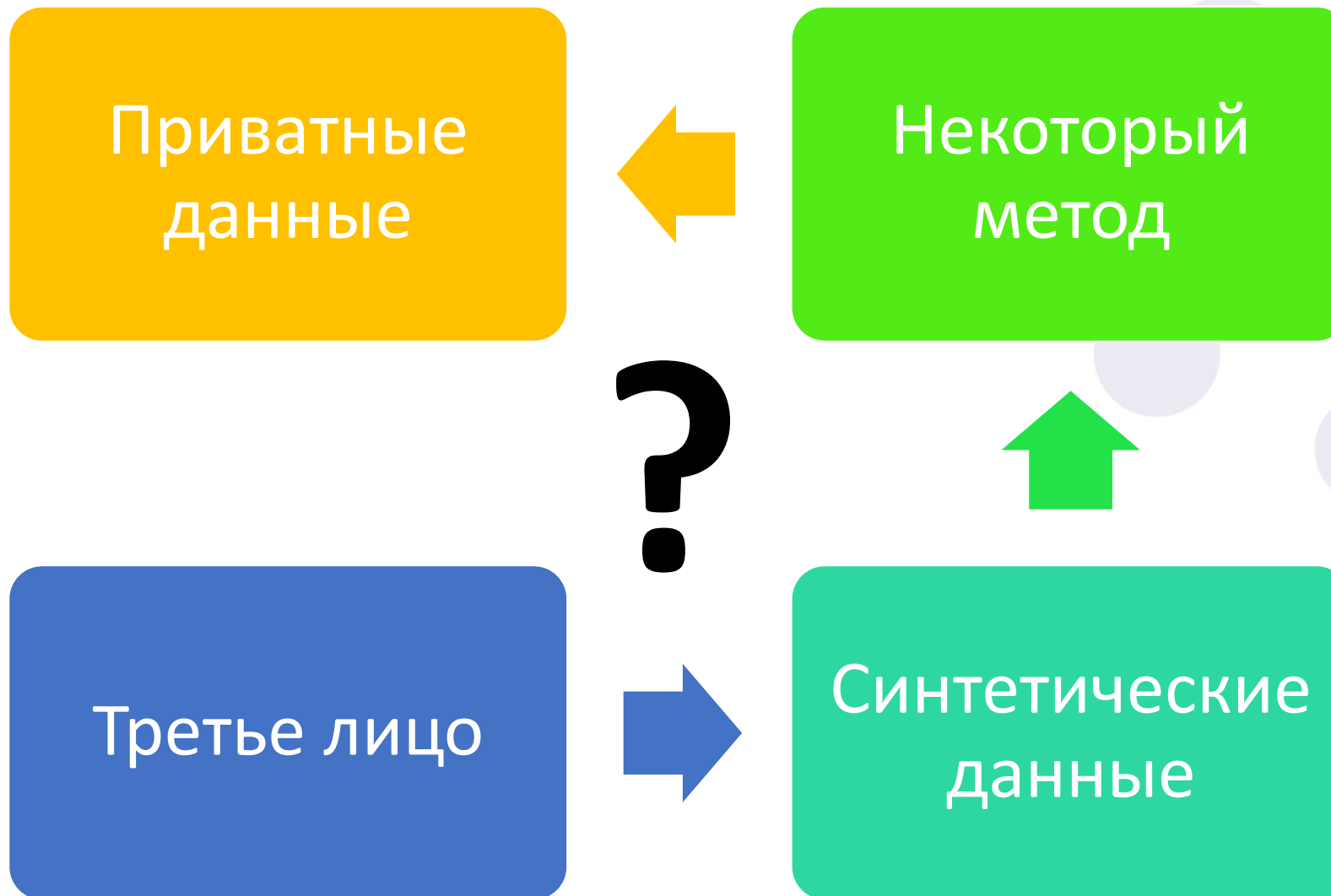




Научно-технологический  
университет

Сириус

# Основная проблема Яндекс





Научно-технологический  
университет

**Сириус**

# Актуальность проблемы

**Яндекс**



- Массовая обеспокоенность сохранением приватности данных
- Утечки частных сведений



MLaaS – Machine Learning  
as a Service

Облачные сервисы для  
ML??

Могут ли данные утечь с  
чужого сервера????

# Цели работы

**Я**ндекс

- Проанализировать статьи по приватному обучению
- Воспроизвести основные результаты
  - Обучение нейросетей
  - Атаки на нейросети
  - Повышение устойчивости к атакам



Научно-технологический  
университет

Сириус

# План презентации

Яндекс

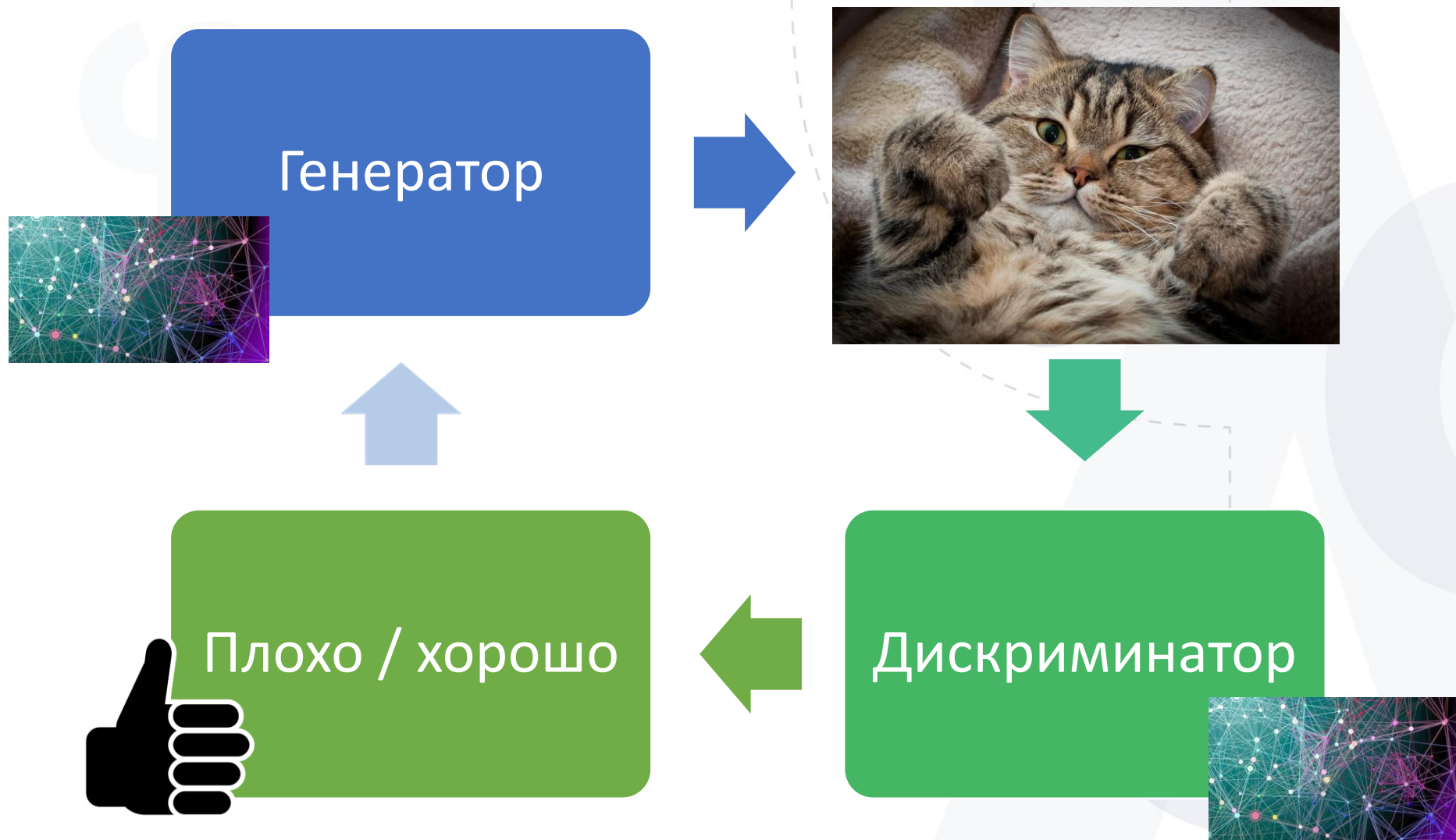
1. Определение GAN
2. Анализ статей
3. Отбор данных
4. Обучение различных GAN (Silly GAN, DCGAN)
5. Проведение атак на нейронные сети
6. Сравнение результатов
7. Выводы



# Определение GAN

# GAN

**Я**ндекс



# Анализ статей

# Статьи

1. LOGAN: Membership Inference Attacks Against Generative Models (<https://arxiv.org/pdf/1705.07663.pdf>)
  - Методы проведения атак
2. PATE-GAN: GENERATING SYNTHETIC DATA WITH DIFFERENTIAL PRIVACY GUARANTEES (<https://openreview.net/pdf?id=S1zk9iRqF7>)
  - Больше математики
  - PATE-GAN
  - Схема контроля приватности

# Отбор датасетов

# Отбор датасетов



MNIST



Kaggle Credit Scoring

# Обучение GAN для MNIST

# Структура Silly GAN

## Генератор

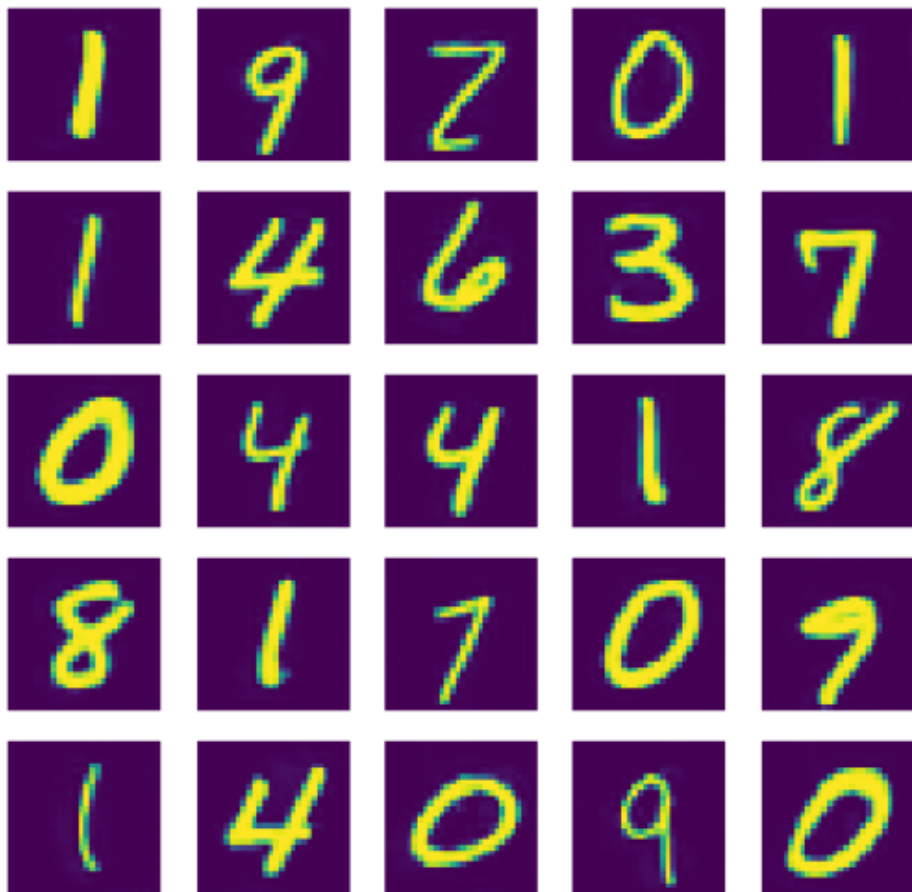
- Dense (ELU)
- Dropout
- Conv2d (kernel\_size=3, ELU)
- Dropout
- ConvTranspose2d
- Dropout
- Conv2d и Dropout
- Conv2d и Dropout

## Дискриминатор

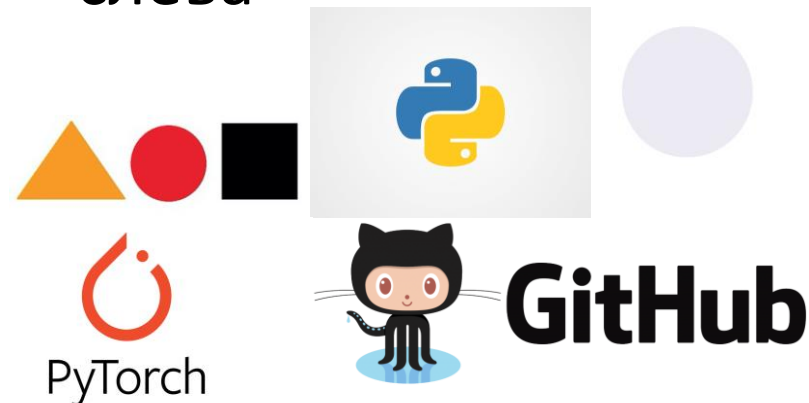
- Conv2d (ELU, kernel\_size=3)
- Dropout
- MaxPool2d(2)
- Conv2d(ELU, kernel\_size=3) и Dropout
- Conv2d(ELU, kernel\_size=3) и Dropout
- Conv2d(ELU, kernel\_size=3) и Dropout
- Linear() with sigmoid



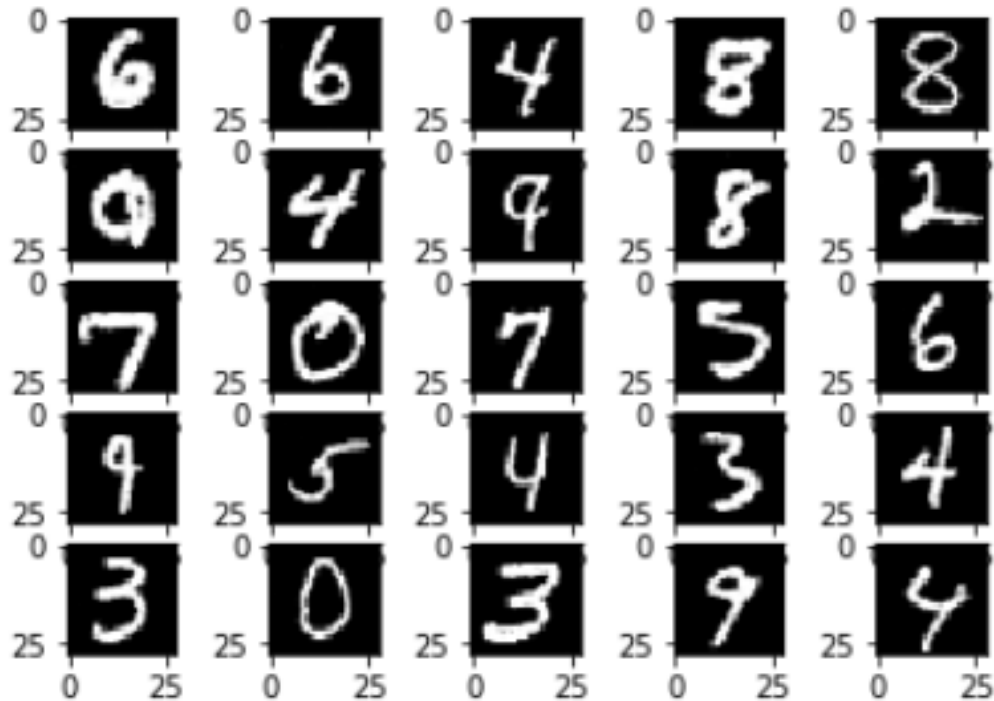
# Обучение своей Silly GAN



- Jensen-Shannon GAN со шумом в losses, регуляризацией градиента и техникой Dropout
- Обучение происходило на сервере ШАД almaren
- Сгенерированные цифры слева



# DCGAN for MNIST



- Использовали предобученную модель ([ссылка](#))
- Jensen-Shannon без шумов, регуляризаций и Dropout



# Обучение GAN для CreditFraud

## Генератор

- Полносвязные слои с ReLU активацией

## Дискриминатор

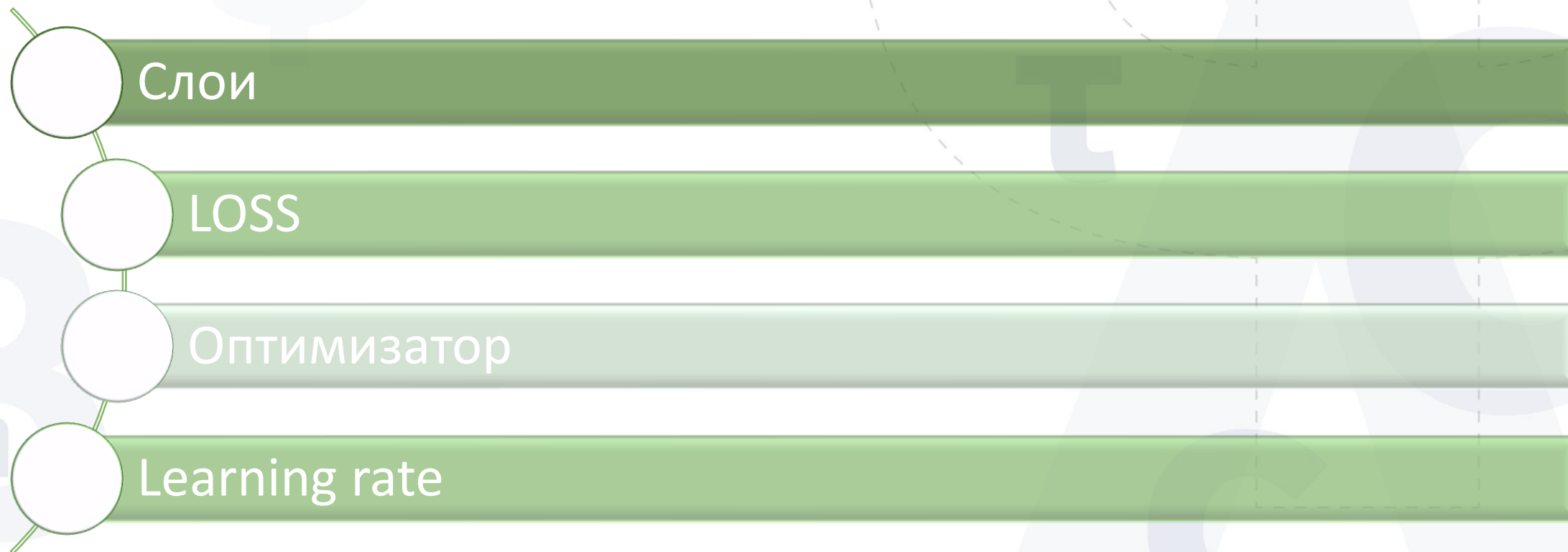
- Полносвязные слои с ReLU активацией

В статье было указано, что сеть обучается и переобучается, но у нас GAN не смогла научиться воспроизводить табличные данные. Далее исследовали только GAN, обученные на MNIST.

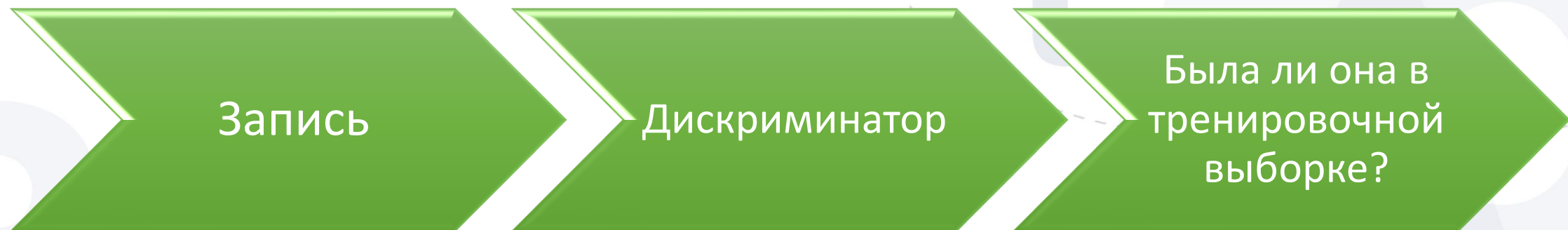
# Атаки на GAN

# Нейросеть как массовый подгон

**Я**ндекс



# Атака на GAN. Общая идея **Я**ндекс





# Атака на GAN. Общая идея **Я**ндекс

А дискриминаторы на деревьях растут???



Запись

Дискриминатор

Была ли она в  
тренировочной  
выборке?

# Атака на GAN

А дискриминаторы на деревьях растут???



ДА!

- White-box attack

Нет, конечно, сам делай

- Black-box attack (with or without leakages)

# Атака на GAN



Black-box  
+ leakage



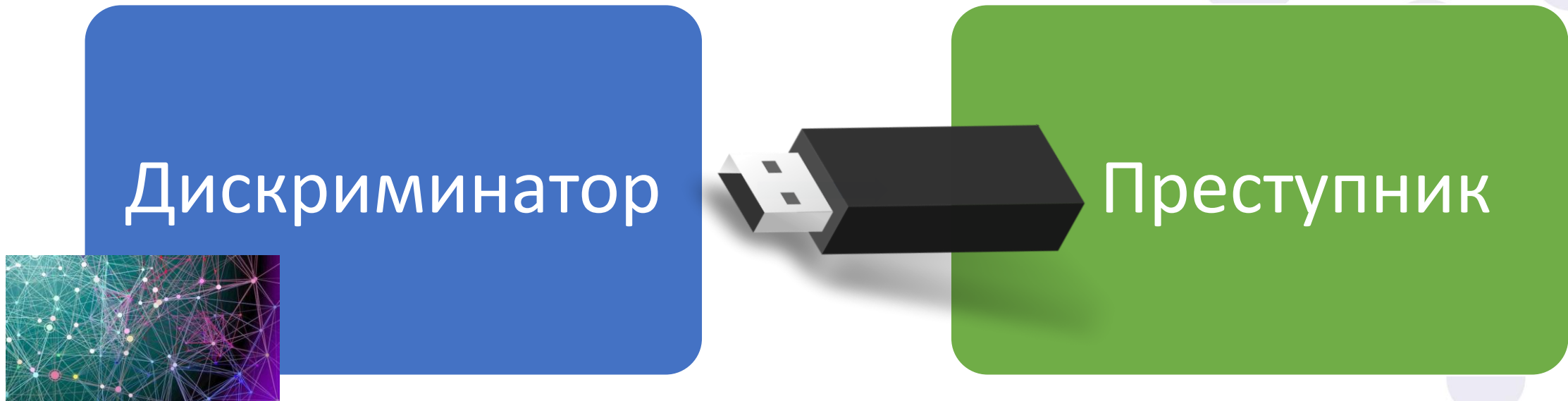
Black-box



White-  
box

# Атаки на GAN. White-box

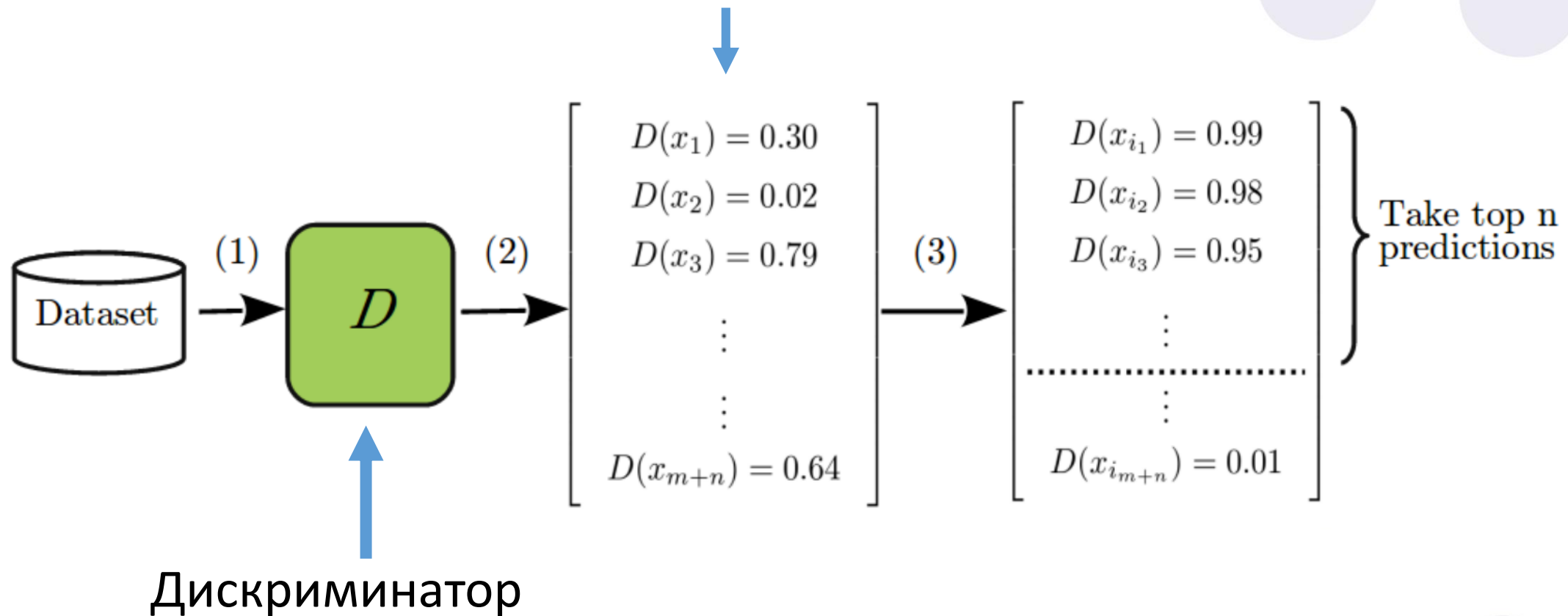
- **Задача.** Пусть преступник **получил** дискриминатор, и он хочет узнать, является ли данный элемент частью тренировочной выборки.



- **Задача.** Пусть преступник **получил** дискриминатор, и он хочет узнать, является ли данный элемент частью тренировочной выборки.
- **Идея:** переобученный дискриминатор сохраняет в себе информацию о тренировочной выборке.

# White-box attack

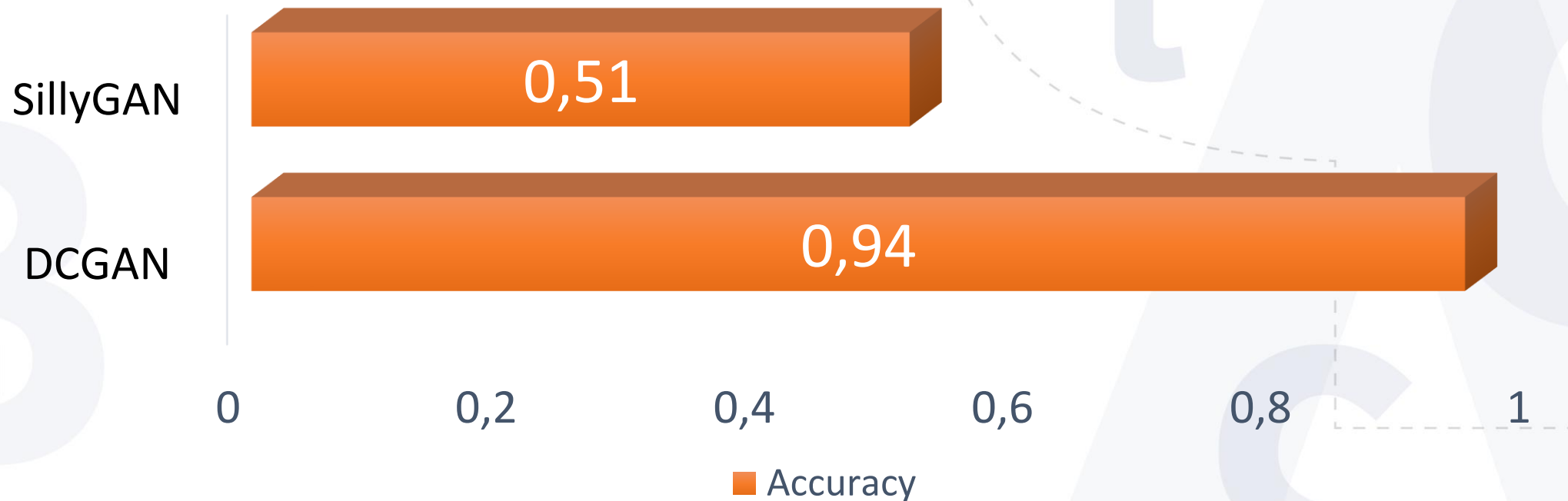
Отклик дискриминатора



# White-box attack

**Y**andex

Accuracy of the attack against the GANs on an perfectly balanced sample

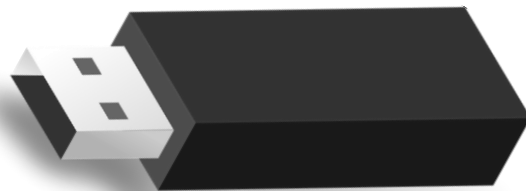




# Атаки на GAN. Black-box + leakage

# Black-box attack with a 20% leakage

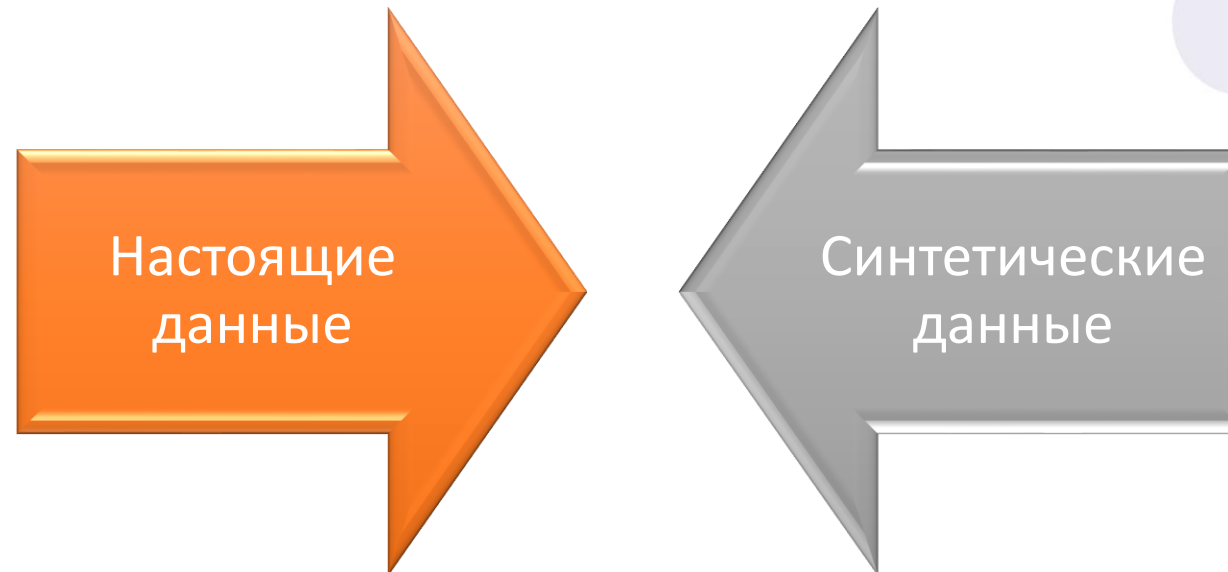
- **Задача.** Произошла **утечка части обучающей выборки**. Злоумышленник имеет доступ к генератору.



Преступник

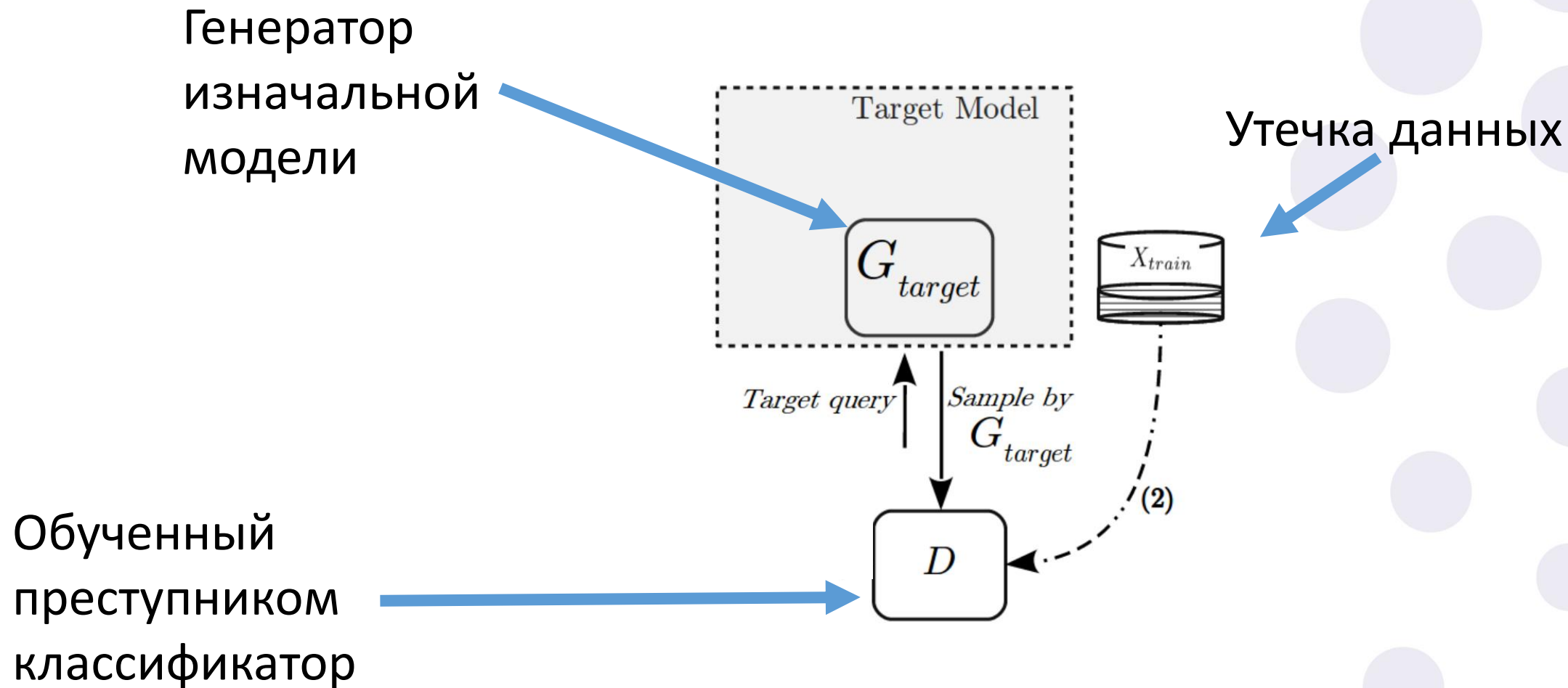
# Black-box attack with a 20% leakage

- **Задача.** Произошла **утечка части обучающей выборки**. Злоумышленник имеет доступ к генератору.
- **Идея.** Создать свой классификатор.



# Black-box attack with a 20% leakage

**Я**ндекс



# Black-box attack with a 20% leakage

**Y**andex

Accuracy of the attack against the GANs on an perfectly balanced sample

Silly GAN against Silly GAN

0,63

Silly GAN against DCGAN

1

0 0,2 0,4 0,6 0,8 1

Accuracy

# Атаки на GAN. Black-box without any leakage

# Black-box attack without any leakage

- **Задача.** Злоумышленник имеет доступ к генератору, но не к обучающей выборке.

# Black-box attack without any leakage

- **Задача.** Злоумышленник имеет доступ к генератору, но не к обучающей выборке.
- **Идея.** Обучить доморощенный GAN.



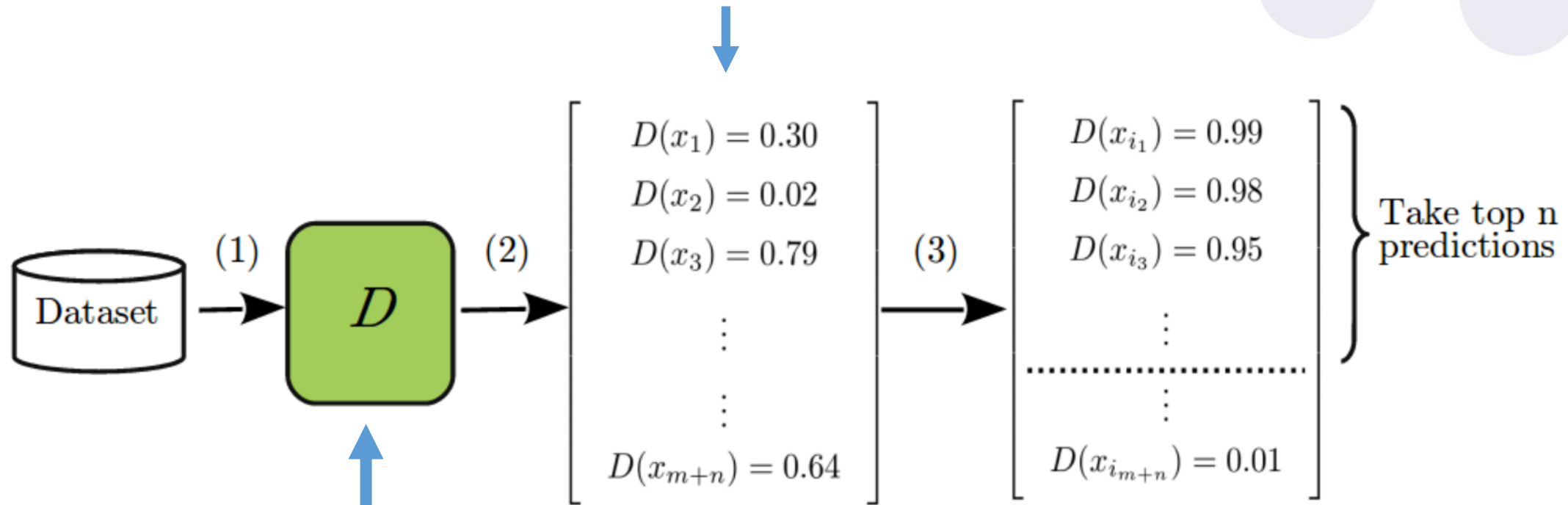


# Black-box attack without any leakage

- **Задача.** Злоумышленник имеет доступ к генератору, но не к обучающей выборке.
- **Идея.** Обучить доморощенный GAN.
- **Метод.** Обучается GAN на синтетической выборке, как если бы это были настоящие данные.

# Black-box attack without any leakage

Отклик дискриминатора

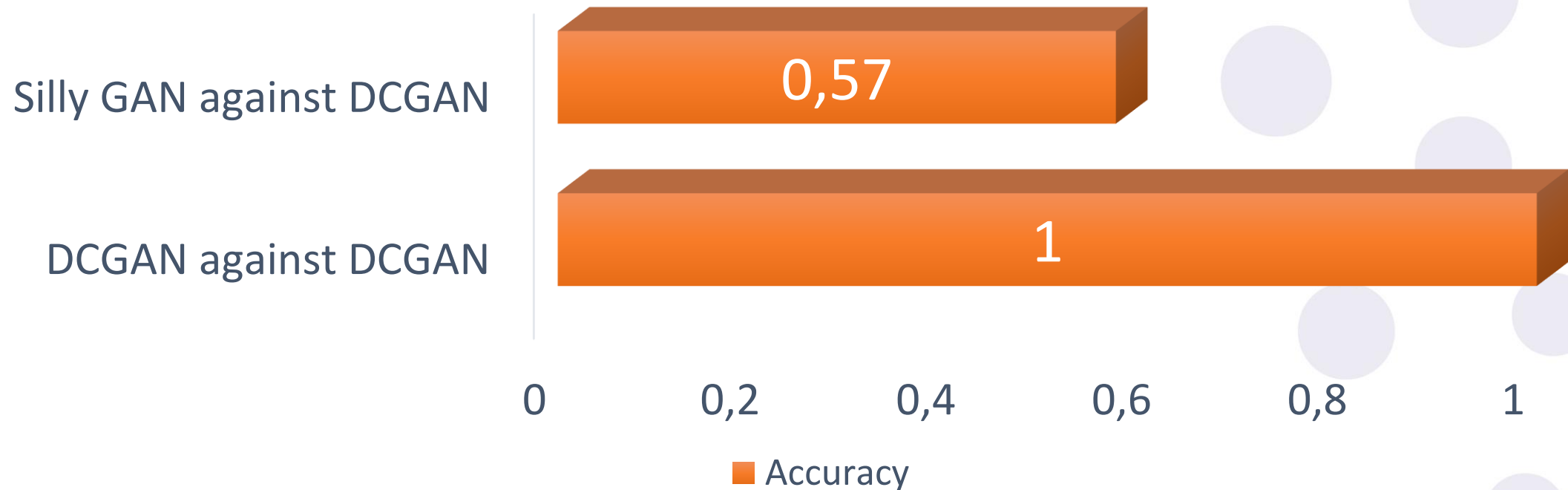


Обученный  
преступником  
дискриминатор

# Black-box attack without any leakage

**Y**andex

Accuracy of the attack against the GANs on an perfectly balanced  
sample



# Сравнение результатов

# Сравнение качества атак

## Black-box + leakage

- Наиболее успешная атака, но требует утечку данных

## Black-box

- Результат сильно зависит от навыков злоумышленника

## White-box

- Атака может быть успешной только при переобученном дискриминаторе

# Сравнение DCGAN и Silly GAN

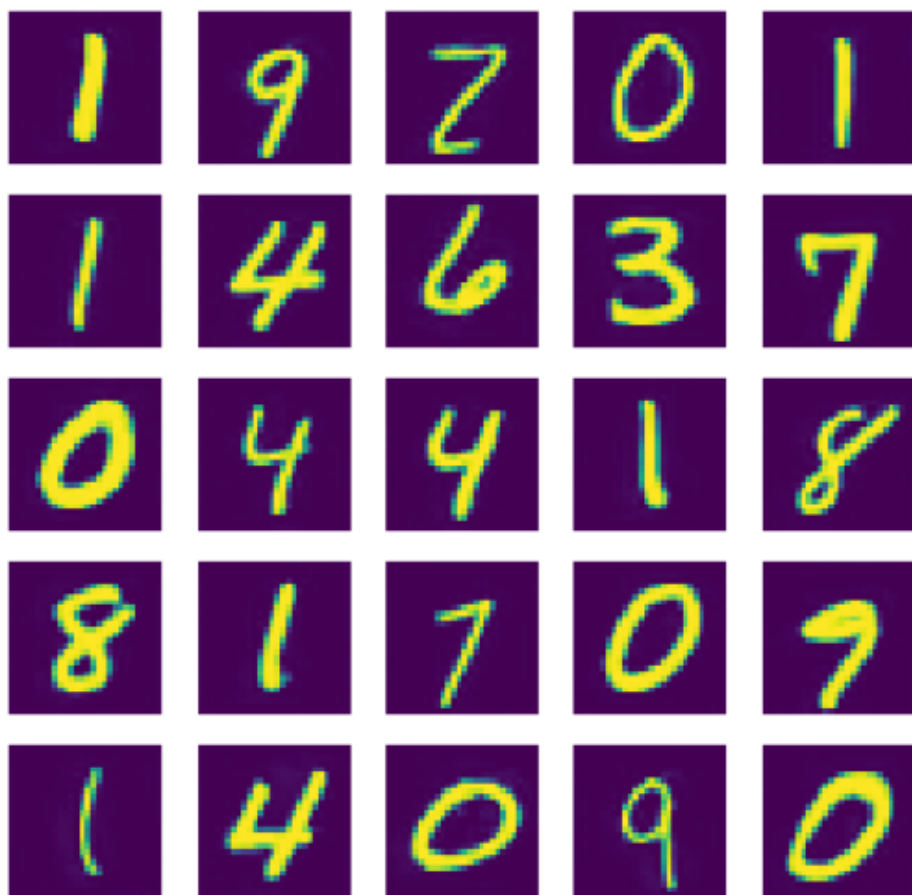
## Silly GAN

- Устойчива к атакам

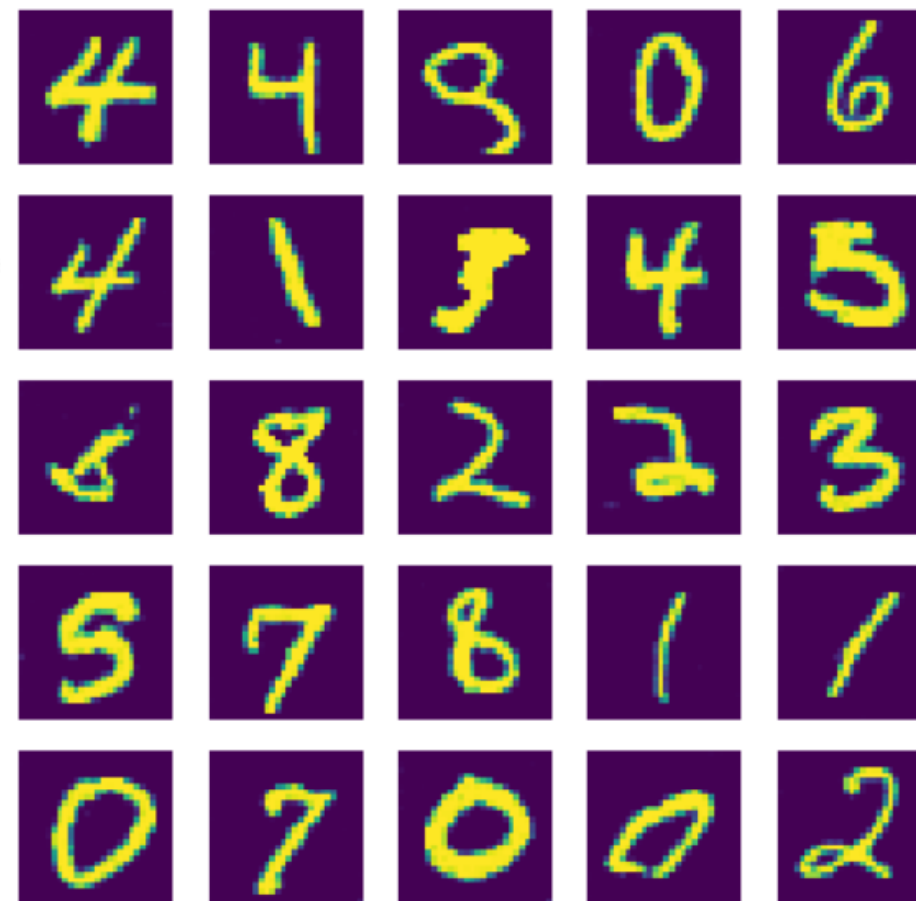
## DCGAN

- Может взломать ребёнка

# Сравнение DCGAN и Silly GAN



Silly GAN



DCGAN

# Как противостоять атакам?

- Борьба с переобучением
  - ✓ Dropout techniques
  - ✓ Зашумление loss
  - ✓ Регуляризация градиента



# Как противостоять атакам?

- Борьба с переобучением
  - ✓ Dropout techniques
  - ✓ Зашумление loss
  - ✓ Регуляризация градиента
- Баланс генератора и дискриминатора

# Как противостоять атакам?

- Борьба с переобучением
  - ✓ Dropout techniques
  - ✓ Зашумление loss
  - ✓ Регуляризация градиента
- Баланс генератора и дискриминатора
- Имитирование атак во время обучения, перед выпуском модели

# Заключение

# Что же мы сделали?

**Я**ндекс

- Научились делать GAN с нуля
- Проанализировали статьи по атакам на GAN

# Что же мы сделали?

- Научились делать GAN с нуля
- Проанализировали статьи по атакам на GAN
- Сделали несколько атак на GAN
- Сравнили качество разных типов атак на GAN

## Что же мы сделали?

- Научились делать GAN с нуля
- Проанализировали статьи по атакам на GAN
- Сделали несколько атак на GAN
- Сравнили качество разных типов атак на GAN
- Поняли, что даже GAN со страшными названиями **легко** поддаются атакам
- Разработали рекомендации по созданию устойчивых к атакам сетей

# Перспективы и планы на будущее

**Я**ндекс

- Экстенсивное исследование связи переобучения и уязвимости к атакам

# Перспективы и планы на будущее

- Экстенсивное исследование связи переобучения и уязвимости к атакам
- Математическая теория дифференциальной приватности
  - Работает ли она на практике?
  - Зачем нужны сложные методы, когда и простые, но грамотно составленные нейросети работают?



# Перспективы и планы на будущее

- Экстенсивное исследование связи переобучения и уязвимости к атакам
- Математическая теория дифференциальной приватности
  - Работает ли она на практике?
  - Зачем нужны сложные методы, когда и простые, но грамотно составленные нейросети работают?
- Ведёт ли сохранение приватности к ухудшению качества?



Научно-технологический  
университет

Сириус

Яндекс

*Спасибо за внимание!*

*Задавайте вопросы*

Ссылка на наш репозиторий в Github: <https://github.com/Private-ML/gan>

# Differential privacy

- Two datasets  $D$  and  $D'$  are said to be neighbouring if there is such  $x$  in  $D$  that  $D \setminus \{x\} = D'$ .
- Let  $(\Omega, \mathcal{F}, P)$  be a probability space, let  $(X, S)$  be a measurable space representing datasets, and let  $(Y, T)$  be a measurable space of outcomes.
- Let  $M$  be a *randomised algorithm*, i.e. a measurable mapping

$$M: X \times \Omega \rightarrow Y$$

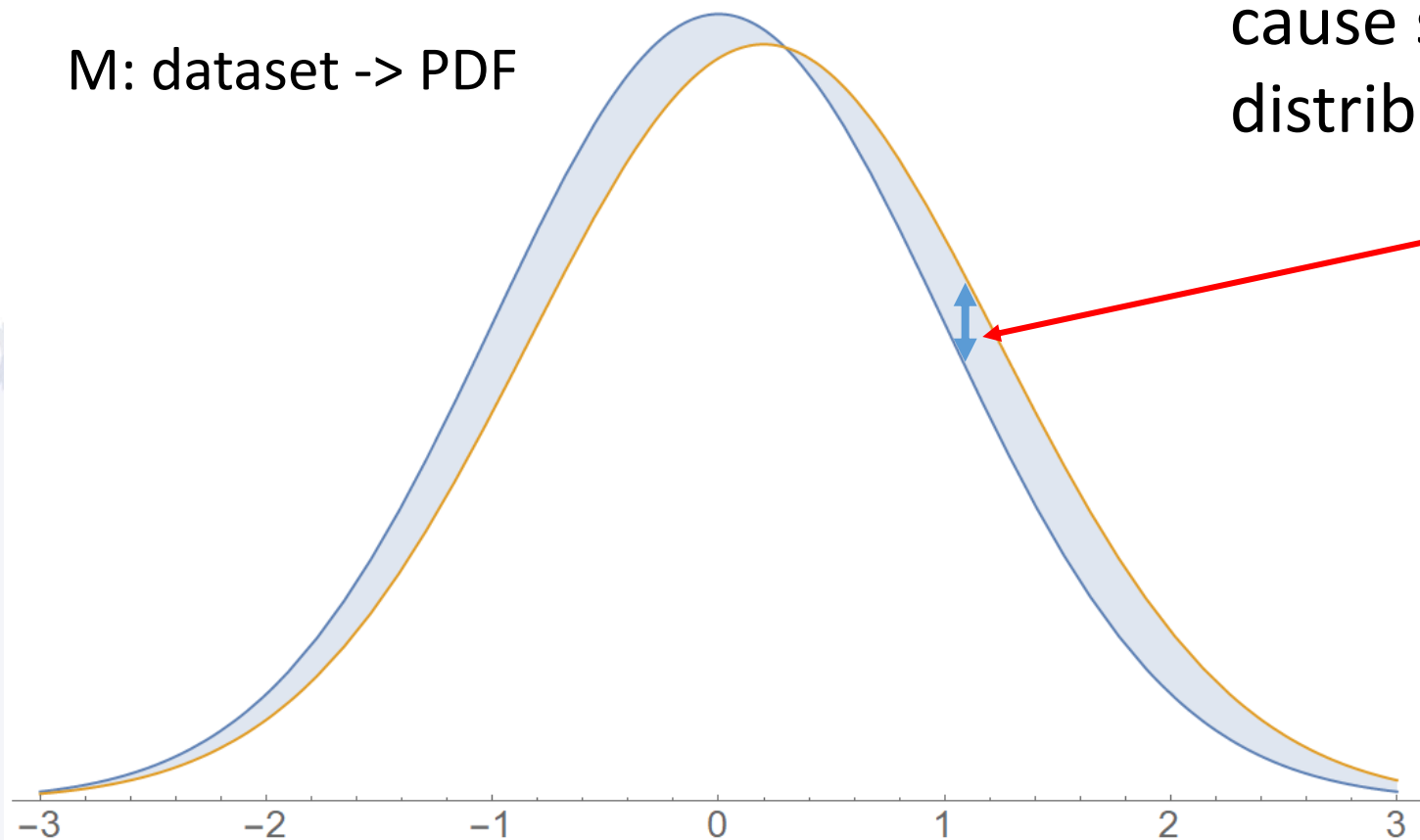
The algorithm is said to be  $(\epsilon, \delta)$ -private, if for all neighbouring datasets  $D$  and  $D'$  and for all sets  $Y \in T$ :  $P(M(D) \in Y) < e^\epsilon P(M(D') \in Y) + \delta$

# Differential privacy

**Y**andex

M: dataset  $\rightarrow$  PDF

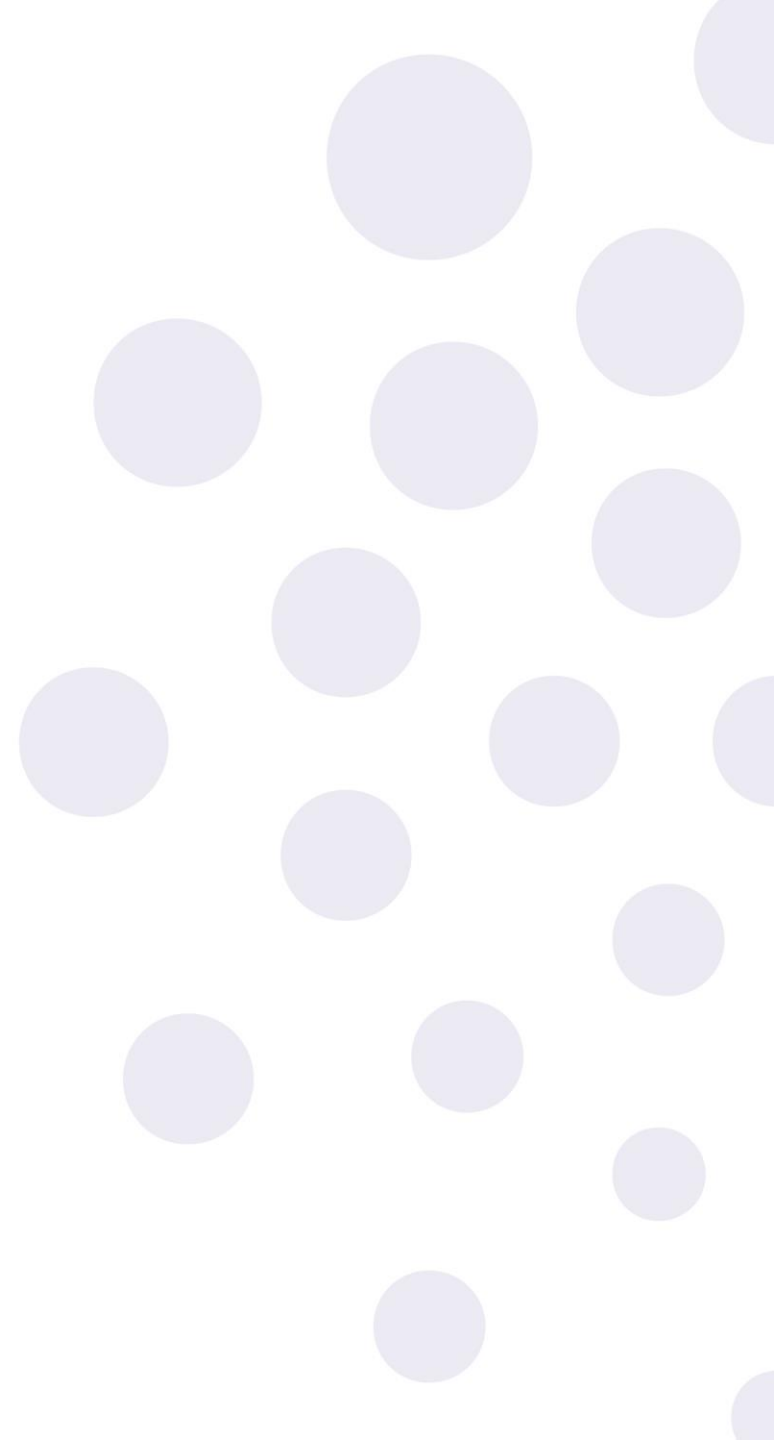
Slight changes of the dataset do not cause significant difference in the distribution of the output





Научно-технологический  
университет

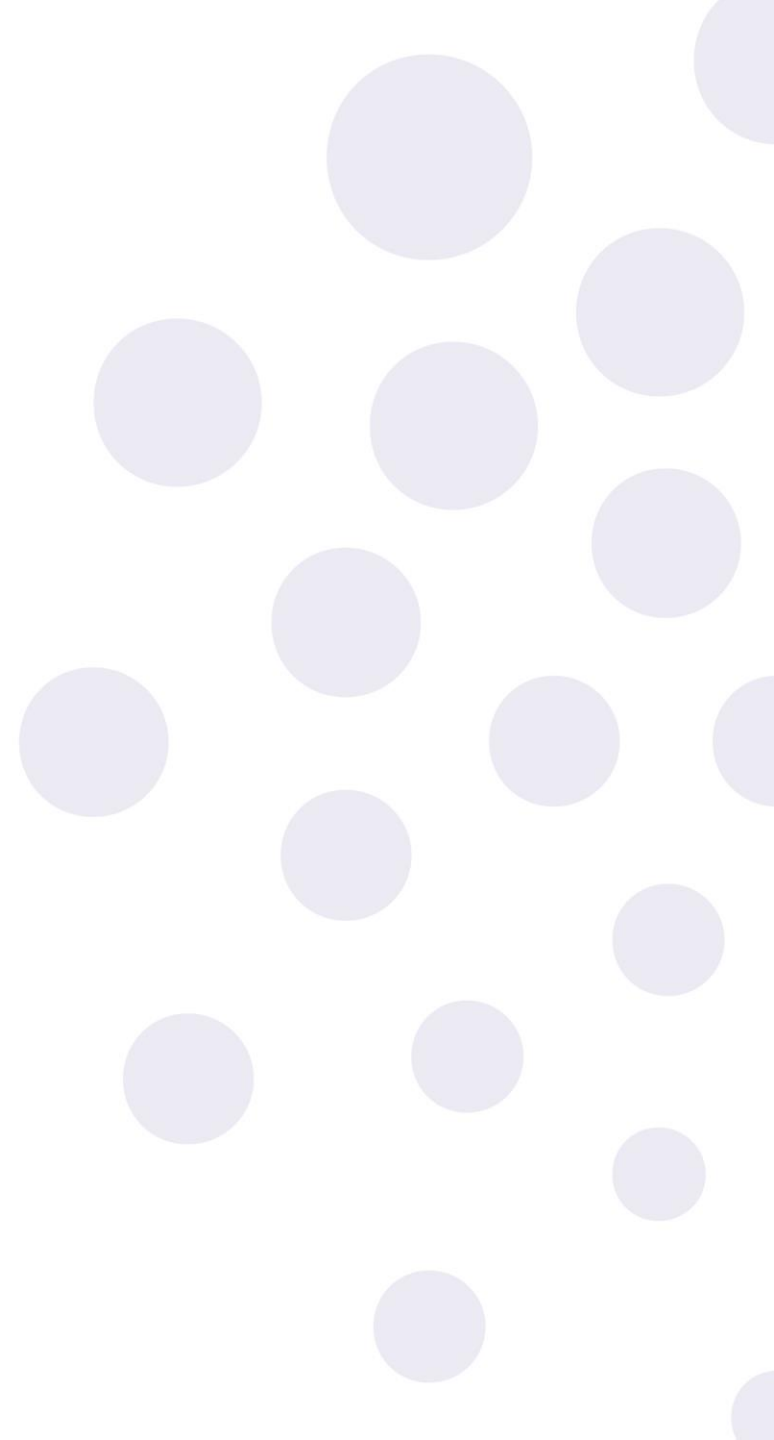
Сириус





Научно-технологический  
университет

Сириус





Научно-технологический  
университет

Сириус

$\varphi$

$\alpha$   
 $\beta$   
 $\eta$

$\tau$

$\sigma$

$\varsigma$



Научно-технологический  
университет

Сириус

$\varphi$

$\alpha$   
 $\beta$   
 $\eta$

$\tau$

$\sigma$

$\varsigma$





Научно-технологический  
университет

Сириус

$\varphi$

$\alpha$   
 $\beta$   
 $\eta$

$\tau$

$\sigma$

$\varsigma$







Научно-технологический  
университет

Сириус

$\varphi$

$\alpha$   
 $\beta$   
 $\eta$

$\tau$

$\sigma$

$\zeta$



Научно-технологический  
университет

Сириус

$\varphi$

$\alpha$   
 $\beta$   
 $\eta$

$\tau$

$\sigma$

$\varsigma$