

# Обнаружение вредоносного кода в зашифрованном с помощью TLS трафике (без дешифровки)

Руслан Иванов

Системный инженер-консультант

[ruivanov@cisco.com](mailto:ruivanov@cisco.com)

# О чём этот доклад?

В докладе рассматривается работа группы исследователей компании Cisco, доказывающая применимость традиционных методов статистического и поведенческого анализа для обнаружения и атрибуции вредоносного ПО, использующего TLS в качестве метода шифрования каналов взаимодействия, без дешифровки или компрометации TLS-сессии.

# Исследователи



Blake Anderson – Technical Leader

PhD in Computer Science (Machine Learning)

Работает в Cisco с 2015 года



David McGrew – Cisco Fellow

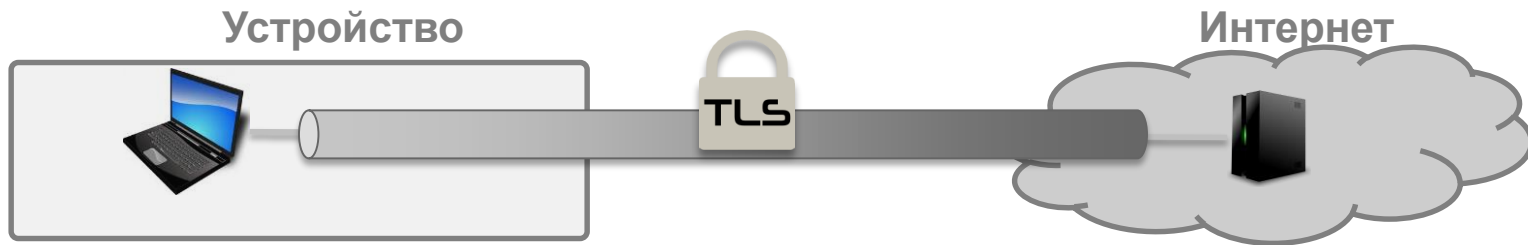
PhD in Physics (Chaos Theory)

Работает в Cisco в с 1998 года

# Содержание

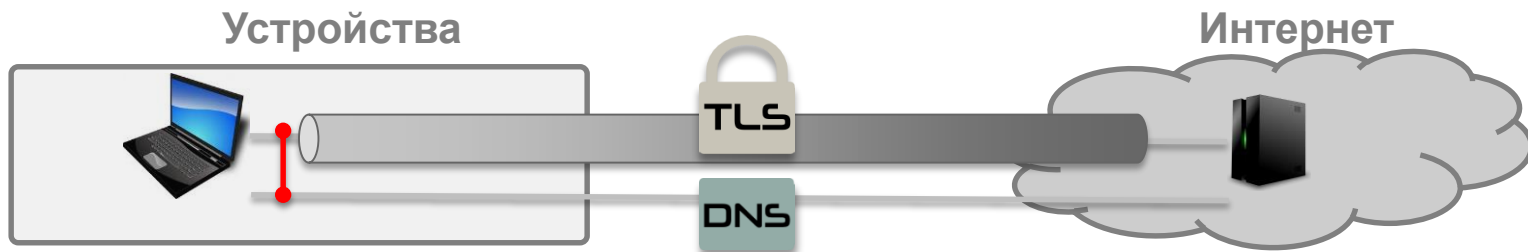
- Введение
- Набор данных
- Дерево принятия решений на основе правил
- Машинное обучение
- Достигнутые результаты
- Выводы

# Проблема: зловредный код активно использует TLS-шифрование



- Шифрование с помощью TLS активно используется (само по себе это не плохо!)
- Решения, основанные на анализе строк, становятся менее эффективными
- Проблемы внедрения расшифровки (MITM) для анализа:
  - Приватность; юридические проблемы; внедрение; стоимость; отсутствие у клиентов желания сотрудничать

# Наш подход: использовать все доступные данные



**Netflow данные:** SrcIP, DstIP, SrcPort, DstPort, Proto, #Bytes, #Packets

**Intraflow данные:** размеры пакетов & временные параметры, распределение байтов, ...

**TLS метаданные:** расширения, наборы шифров, SNI, поля сертификатов, ...

**DNS данные:** имена, типы запросов, временные параметры запросов

**HTTP данные:** заголовки и сопутствующие поля, в том числе других http-запросов с этого же хоста

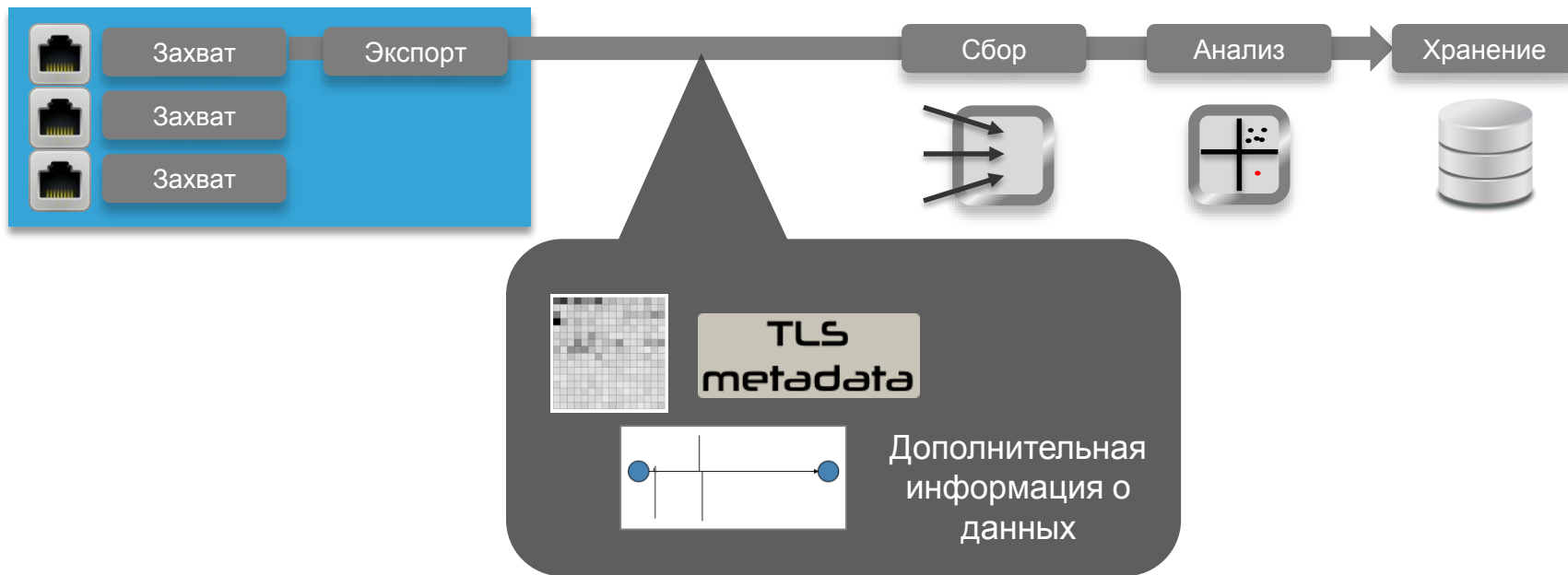
# Традиционный анализ сетевых потоков

srcIP, dstIP, srcPort, dstPort, prot, startTime, stopTime, numBytes, numPackets



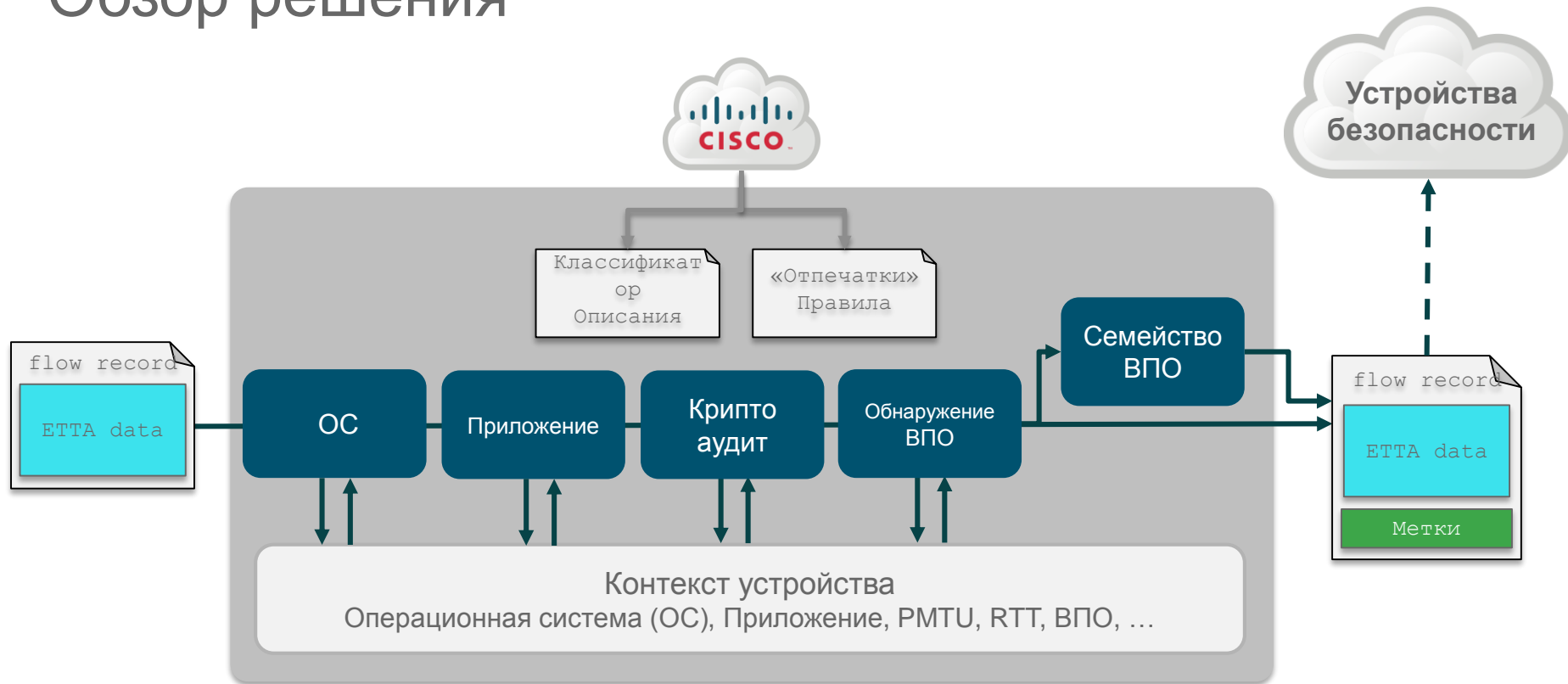
# Сбор расширенной телеметрии

srcIP, dstIP, srcPort, dstPort, prot, startTime, stopTime, numBytes, numPackets



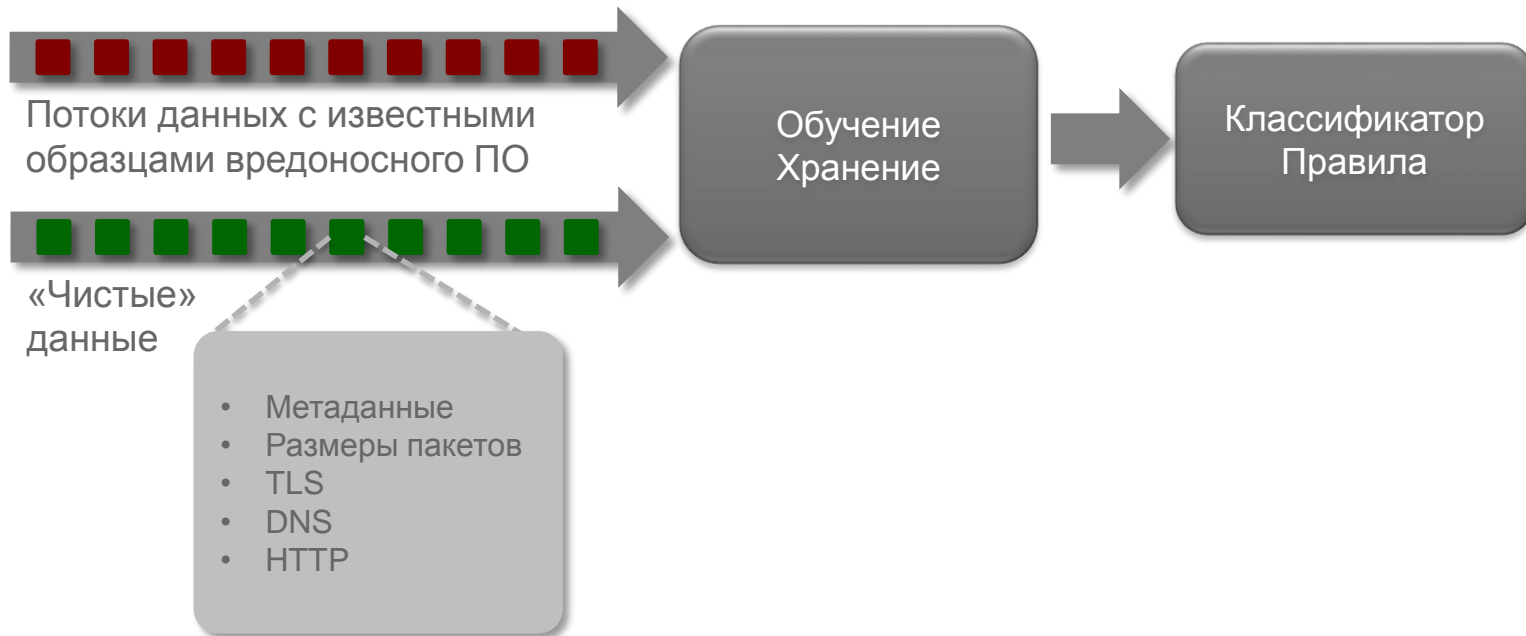


# Обзор решения



# Набор данных

# Набор данных, использованный в исследовании



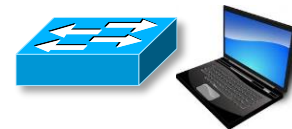
# Потоки данных с известными образцами ВПО

- Записи сетевого обмена из системы динамического анализа Cisco ThreatGRID (в формате pcap):
  - 5-ти минутные сессии анализа;
  - Образцы с Threat Score = 100/100;
- Миллионы pcap-файлов:
  - ~5,000-15,000 новых каждый день;
  - Сотни миллионов потоков



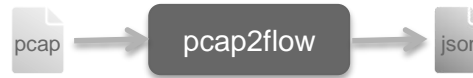
# «Чистые» данные

- Демилитаризованная зона сети крупной компании:
  - ~10-15 миллионов потоков в день;
  - ~500 пользователей
- IP-адреса анонимизированы.

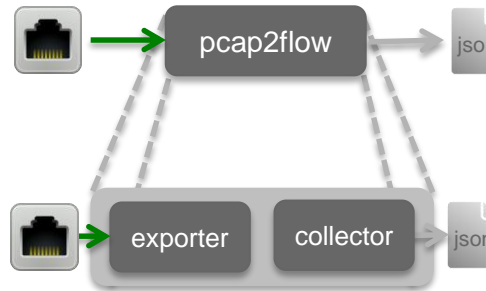


# Как сделать анализ потоков масштабируемым?

joy



Offline



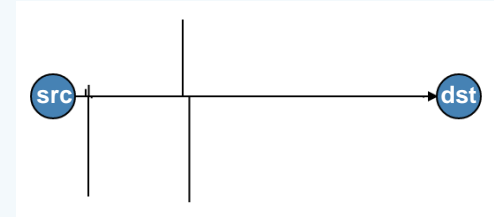
Online

<https://github.com/cisco/joy>

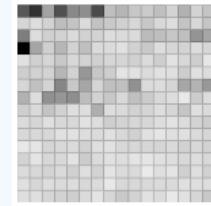


# Типы данных, использованные из расширенной телеметрии

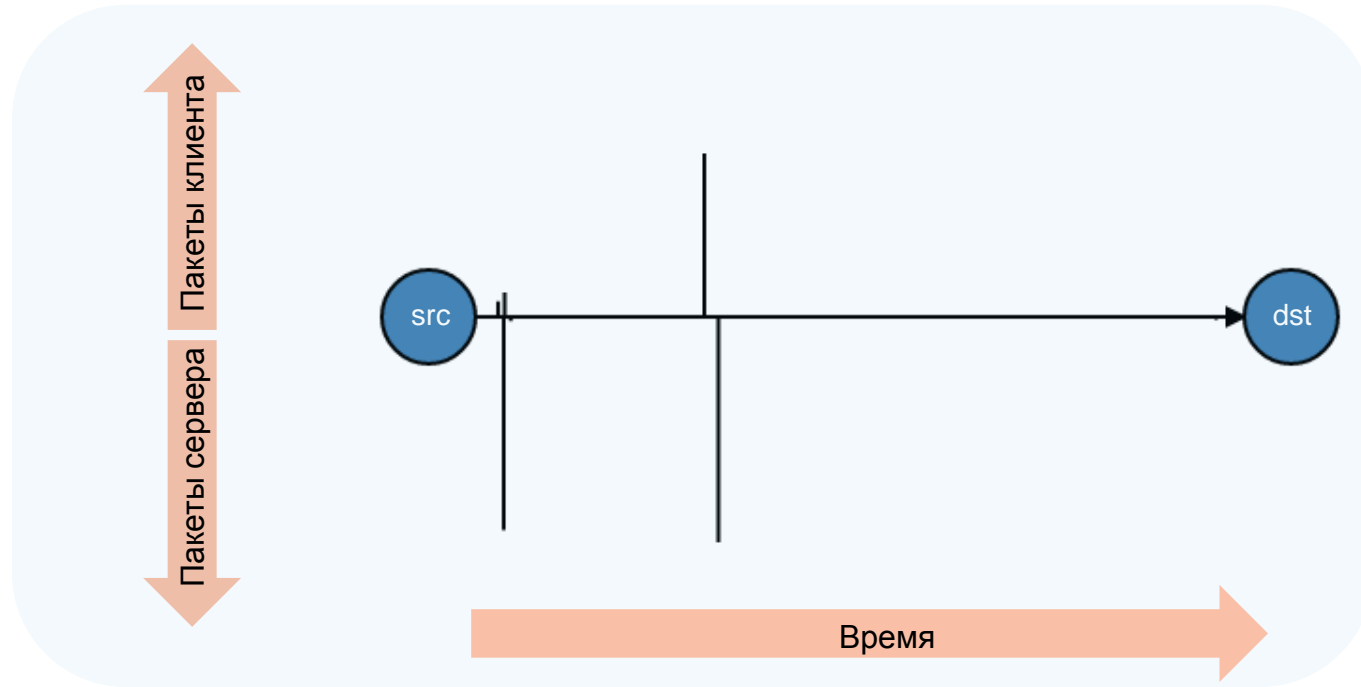
- **SPLT** – Sequence of Packet Lengths and Arrival Times, или распределение пакетов и их последовательностей с учётом временных интервалов



- **BD** - Byte Distribution или побайтное распределение
- **BE** - Byte Entropy или побайтная энтропия



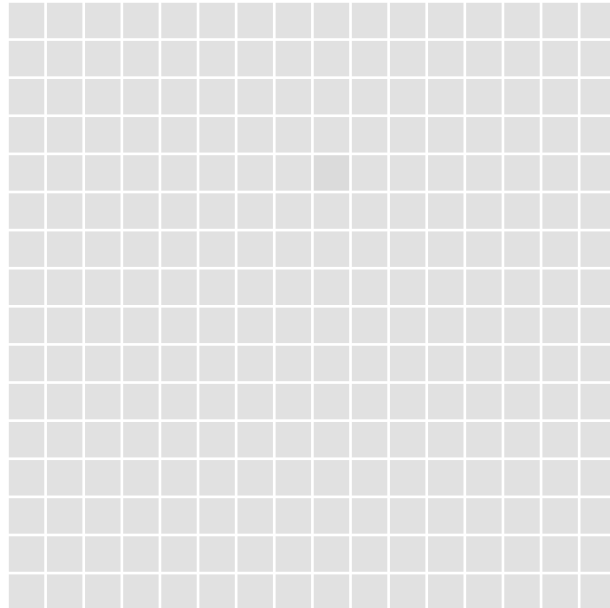
# SPLT - распределение пакетов и их последовательностей с учётом временных интервалов



# VD - побайтное распределение

Н Т Т Р / 1 . 1 2 0 0 О К

8 54 54 50 2f 31 2e 31 20 32 30 30 20 4f 4b

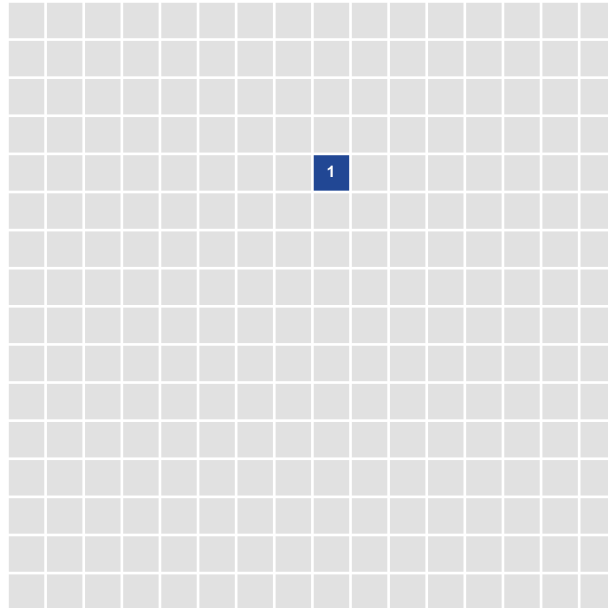




# VD - побайтное распределение

Н Т Т Р / 1 . 1 2 0 0 О К

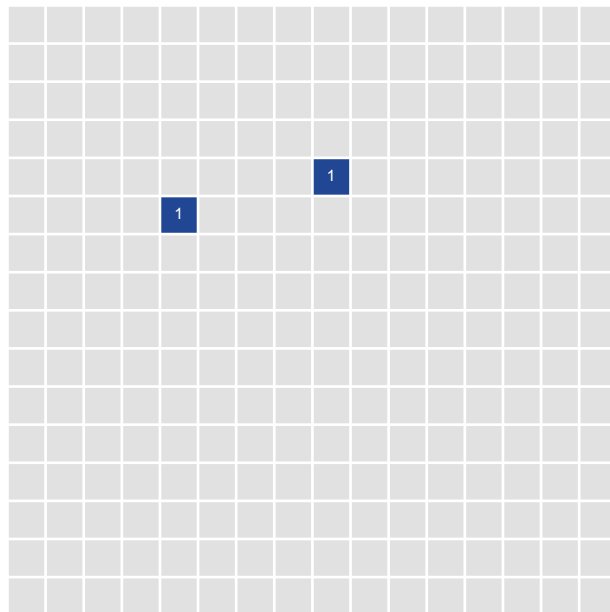
8 54 54 50 2f 31 2e 31 20 32 30 30 20 4f 4b



# VD - побайтное распределение

Н Т Т Р / 1 . 1 2 0 0 О К

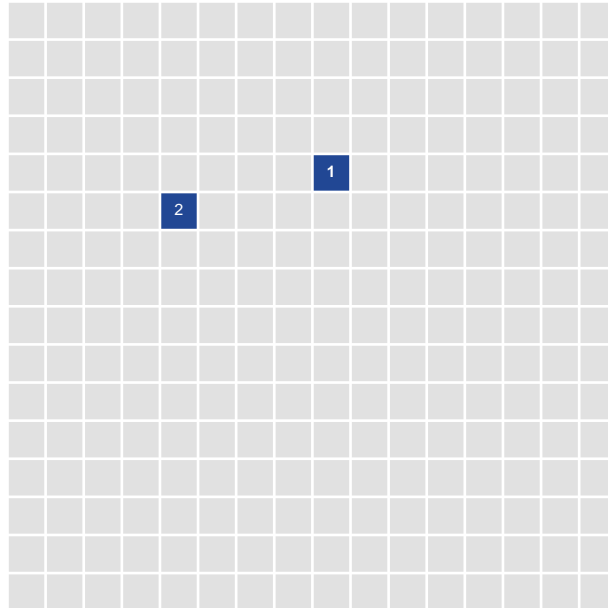
08 54 54 50 2f 31 2e 31 20 32 30 30 20 4f 4b



# VD - побайтное распределение

Н Т **Т** Р / 1 . 1 2 0 0 О К

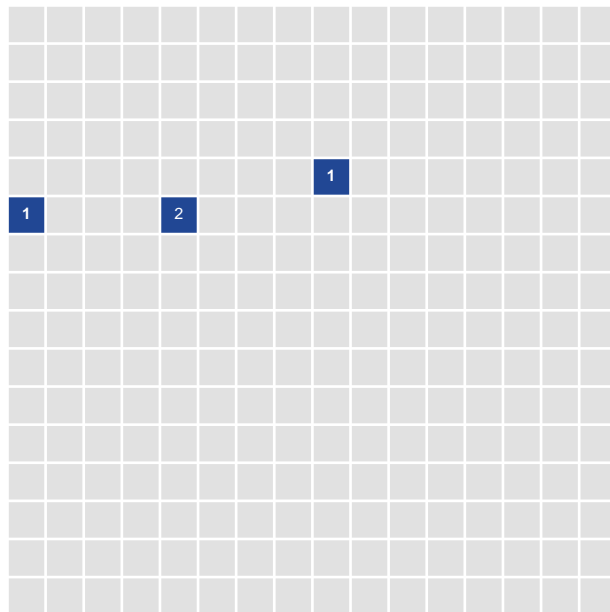
8 54 54 50 2f 31 2e 31 20 32 30 30 20 4f 4b



# VD - побайтное распределение

Н Т Т Р / 1 . 1 2 0 0 О К

8 54 54 50 2f 31 2e 31 20 32 30 30 20 4f 4b

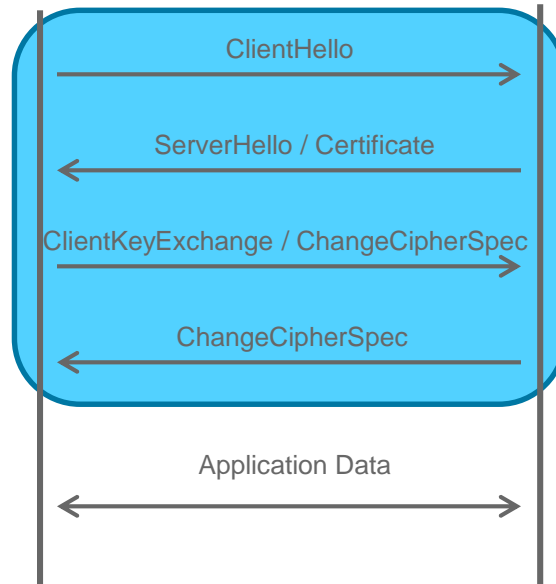


# Установка TLS-соединения

Клиент



Сервер



# TLS-клиенты

- Наиболее популярные TLS-клиенты:
  - [IE 8/11\\*](#), [Tor](#), [Opera](#)
- Большой разброс в количестве и типе предлагаемых наборов шифров
  - Некоторое ВПО всегда использует один набор, некоторое – сотни вариантов наборов
- [2048-bit DHE\\_RSA](#) наиболее распространённый размер и алгоритм для публичного ключа

# TLS-серверы

- Наиболее распространённые серверы (по информации из сертификатов):
  - Кошельки крипто-валюты Bitcoin ([block.io](https://block.io));
  - Файлообменные сервисы ([dropbox.com](https://dropbox.com));
  - Рекламные платформы ([criteo.com](https://criteo.com));
  - Поисковые движки ([google.com](https://google.com), [baidu.com](https://baidu.com));
  - Поля Subject в сертификатах очень часто используют имена, похожие на имена, полученные с использованием алгоритмов генерации доменных имён (DGA);
- 0.7% (по сравнению с 0.09% в корпоративной сети) используют самоподписанные сертификаты;
- **TLS\_RSA\_WITH\_3DES\_EDE\_CBC\_SHA** наиболее часто выбираемый шифр-набор;
- Набор всех этих функций очень сильно зависит от принадлежности к тому или иному семейству вредоносного ПО

# HTTP-заголовки

- Различное использование Заглавных и прописных БУКВ (иногда с опечатками):
  - `User-Agent` vs `user-agent` vs `User-agent` vs `USER-AGENT` vs `User-AgEnt`
- Сохранение порядка следования заголовков;
- Наличие или отсутствие HTTP-заголовков, несущих информацию об окружении:
  - `via`, `x-imforwards`
- Распространённые значения User-Agent:
  - `Opera/9.50(WindowsNT6.0;U;en)`
  - `Mozilla/5.0(Windows;U;WindowsNT5.1;en-US;rv:1.9.2.3)Gecko/20100401Firefox/3.6.1(.NETCLR3.5.30731)`
  - `Mozilla/5.0(WindowsNT6.3;WOW64;Trident/7.0;Touch;rv:11.0)likeGecko`

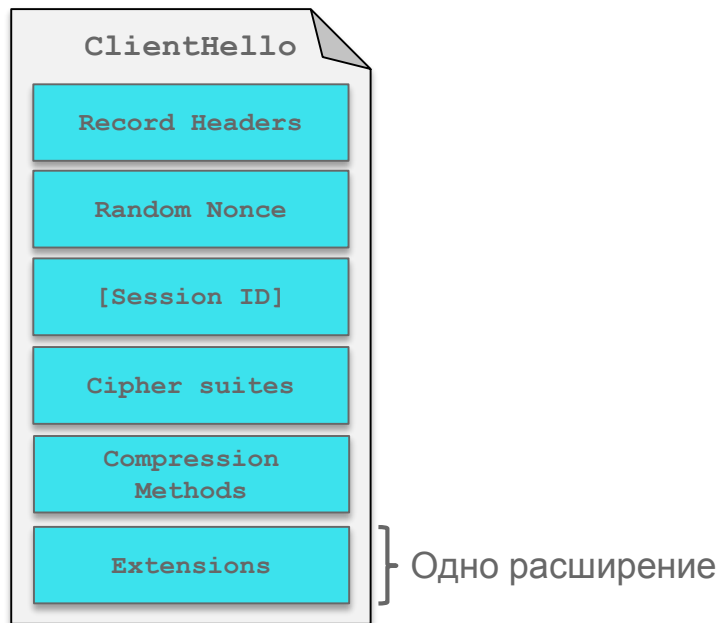


# Информация о DNS-запросах

- Наиболее распространённые значения TTL:
  - 100 > 300 > 60 > 3600
- Количество IP в запросе:
  - 1 > 4 > 11 > 2 > 6
- Наиболее часто встречающиеся доменные суффиксы:
  - com > net > pl > eu > org
- ~95% доменов отсутствует в списке Alexa top-1,000,000

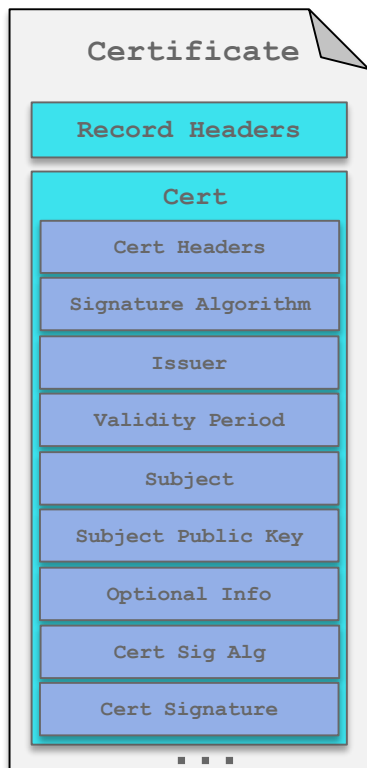
# Дерево принятия решений на основе правил

# Указание имени сервера - Server Name Indication



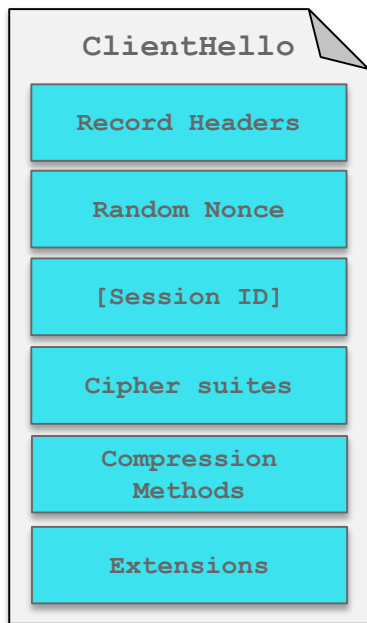
- Extension type: server\_name (0x0000)
- Клиент явно запрашивает определённое имя хоста у сервера, с которым устанавливает соединение

# Информация из полей сертификата

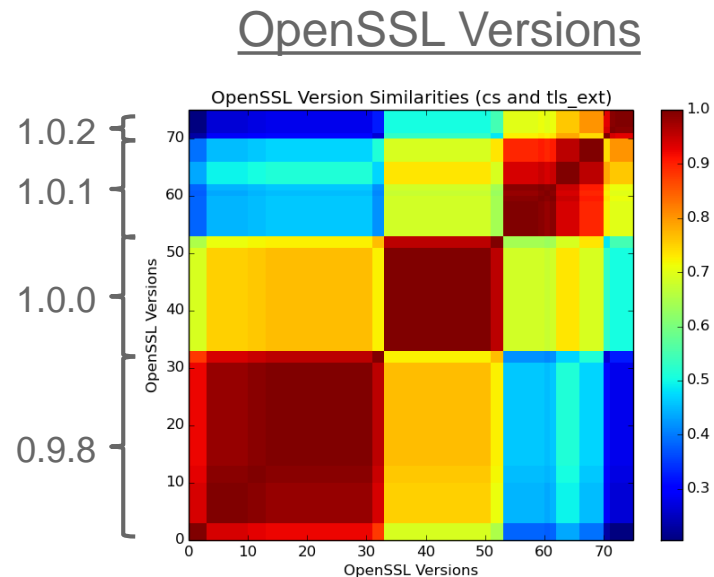


- Проверка на самоподписанный сертификат:
  - Issuer == Subject
- Проверка сертификатов на валидность (время жизни уже закончилось или ещё не началось);
- Проверка расхождений в поле SubjectAltName (опциональные расширения и имя сервера в server\_name расширении)

# «Отпечатки пальцев» TLS-клиентов

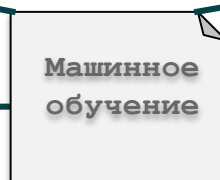
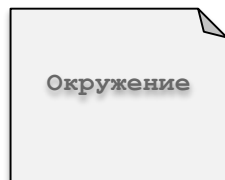
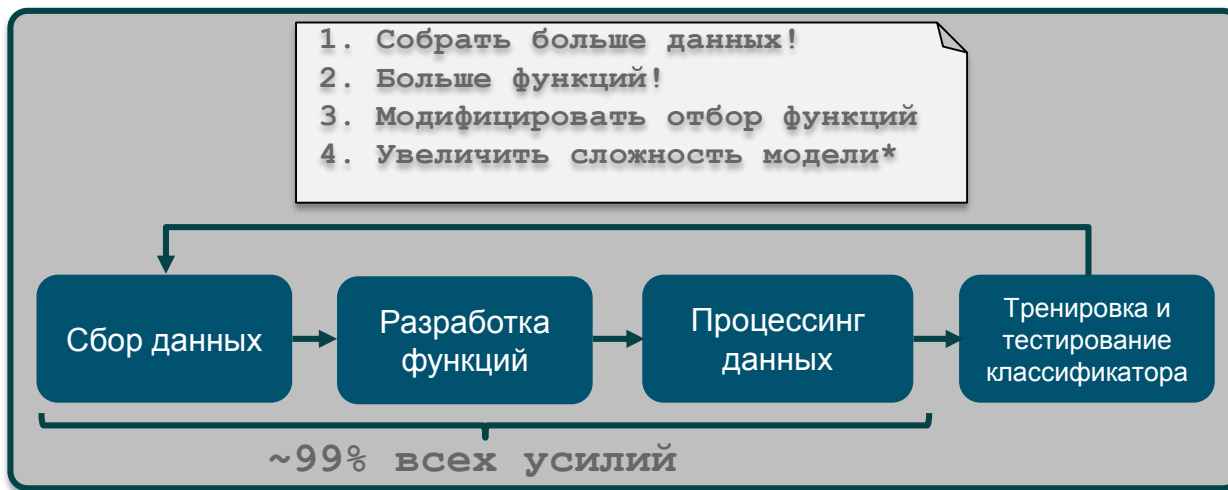


Используется для снятия отпечатков



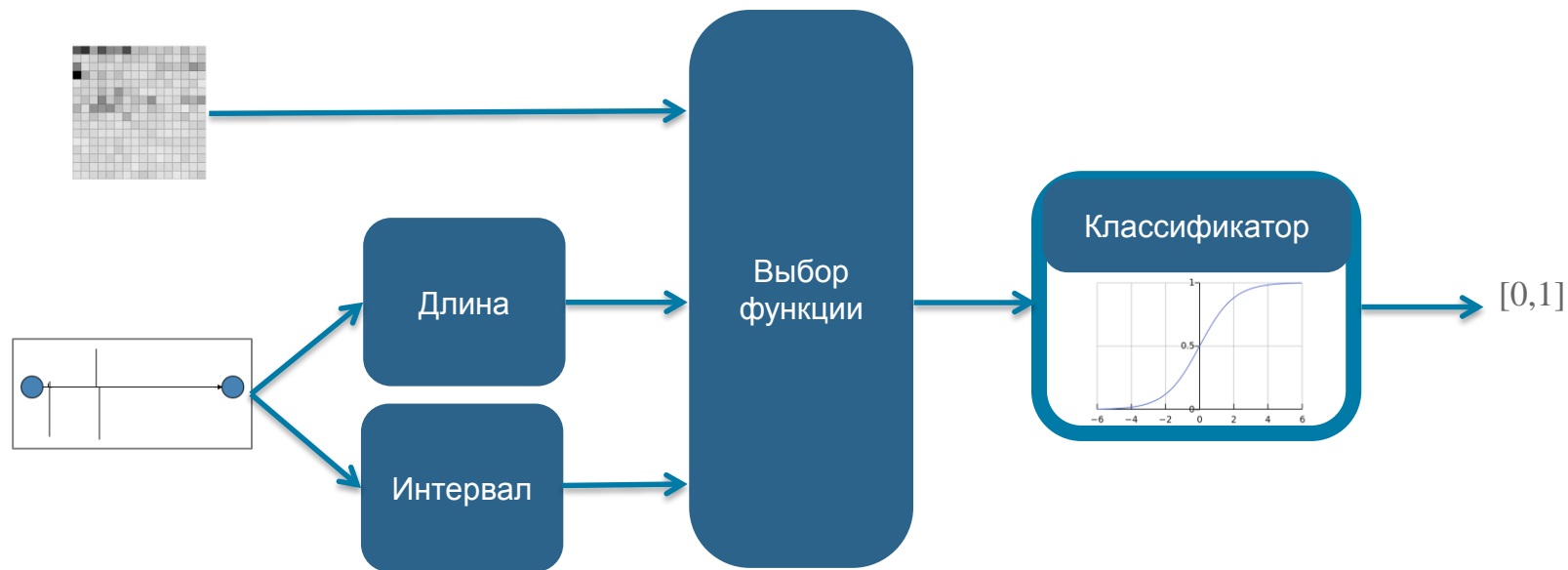
# Машинное обучение

# Управляемое данными обучение



\* Если больше ничего не помогает

# Модель обучение для SPLT и BD





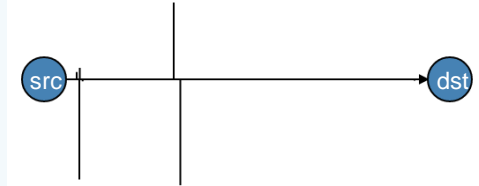
# Достигнутые результаты

# Тестовый набор данных

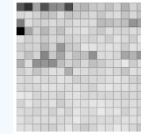
- Вредоносный код
  - Записи сетевого обмена (pcap) с августа 2015 по май 2016 из Cisco ThreatGRID;
  - Трафик TLS (443), больше 100 байт на вход и выход;
  - 225,740 потоков; сетевая телеметрия обогащалась информацией о TLS-extensions, шифр-наборах, и размерах публичных ключей;
- «Хороший трафик»
  - Трафик взят из DMZ крупной компании
  - Трафик TLS (443), больше 100 байт на вход и выход;
  - 225,000 потоков; сетевая телеметрия обогащалась информацией о TLS-extensions, шифр-наборах, и размерах публичных ключей;
- 10-кратная перекрёстная проверка данных

# Междупоточные характеристики

- Длины пакетов и временные интервалы моделировались с использованием цепей Маркова



- Учитывалась вероятность найти определённые значения байтов



- Использовался бинарный вектор из предложенных наборов шифров, расширений, и размеров публичных ключей

**TLS**  
**metadata**

# Результаты классификации по семействам ВПО

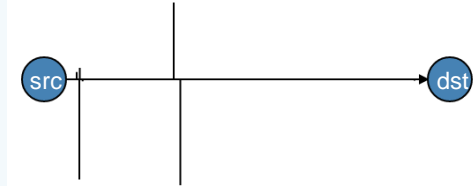
<u>Malware Family</u>	<u>Meta+SPLT+BD</u>	<u>Meta+SPLT+BD+TLS</u>
Bergat*	100.0%	100.0%
Sality*	95.0%	97.7%
Dridex	16.5%	78.5%
Skeeyah	95.9%	98.6%
Virlock	100.0%	100.0%

# Тестовый набор данных №2

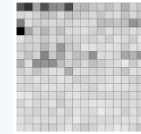
- Вредоносный код
  - Записи сетевого обмена (pcap) с августа 2015 по май 2016 из Cisco ThreatGRID;
  - Трафик TLS (443), больше 100 байт на вход и выход;
  - 13,542 потоков со всеми доступными данными (~63.2% зловредных TLS-потоков);
- «Хороший трафик»
  - Трафик взят из DMZ крупной компании;
  - Трафик TLS (443), больше 100 байт на вход и выход;
  - 42,927 потоков (~3.8% от общего числа «чистых» TLS-потоков);
- 10-кратная перекрёстная проверка данных.

# Междупоточные характеристики

- Длины пакетов и временные интервалы моделировались с использованием цепей Маркова



- Учитывалась вероятность найти определённые значения байтов



- Использовался бинарный вектор из предложенных наборов шифров, расширений, и размеров публичных ключей

**TLS**  
**metadata**

# Добавлена контекстная информация

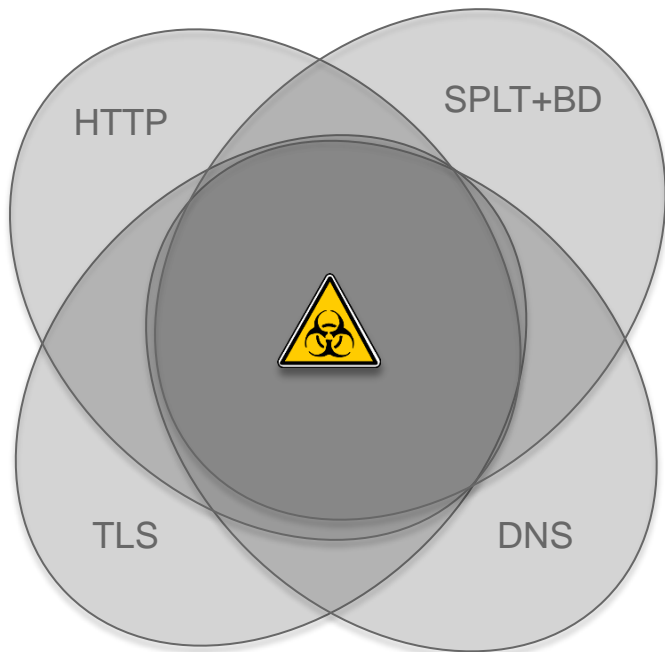
- DNS

- Alexa Lists
- Lengths of DN and FQDN
- Suffix
- TTL
- % Numerical Characters
- % Non-alphanumeric Chars

- HTTP

- Outbound/inbound header fields
- Content-Type
- User-Agent
- Accept-Language
- Server
- code

# Обнаружение зашифрованного ВПО



	Acc.	FDR
SPLT+BD+TLS+HTTP+DNS	99.993%	99.978%
SPLT+BD+TLS+HTTP	99.983%	99.956%
SPLT+BD+TLS+DNS	99.968%	98.043%
SPLT+BD+TLS	99.933%	70.351%
HTTP+DNS	99.985%	99.956%
TLS+HTTP	99.955%	99.660%
TLS+DNS	99.883%	96.551%
HTTP	99.945%	98.996%
DNS	99.496%	94.654%
TLS	94.836%	50.406%

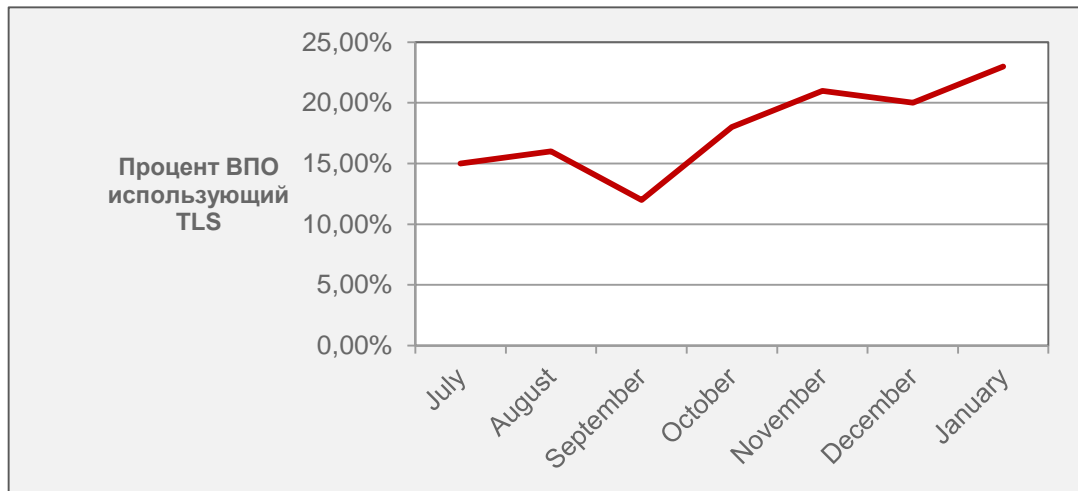


# Интересные факты:

- Зловреды наиболее часто используют:
  - DNS Suffix: org
  - DNS TTL: 3600
  - TLS\_RSA\_WITH\_RC4\_128\_SHA
  - HTTP Field: location
  - DNS Alexa: Not Found
  - HTTP Server: nginx
  - HTTP Code: 404
- Легитимный трафик наиболее часто использует:
  - TLS Ext: extended\_master\_secret
  - Content type: application/octet-stream
  - TLS\_DHE\_RSA\_WITH\_DES\_CBC\_SHA
  - HTTP Server: Microsoft-IIS/8.5
  - DNS Alexa: top-1,000,000
  - HTTP User-Agent: Microsoft-CryptoAPI/6.1

# Выводы

- Использование TLS вредоносным трафиком растёт;
- Данная работа хорошо дополняет ранее известные методы обнаружения ВПО;
- Рассмотренные методы машинного обучения и классификаторы могут применяться для пассивного обнаружения коммуникаций ВПО;
- <https://github.com/cisco/joy>
  - Напишите авторам, чтобы получить хорошо натренированные классификаторы ;)



# Исследователи



Blake Anderson – Technical Leader

PhD in Computer Science (Machine Learning)

Работает в Cisco с 2015 года



David McGrew – Cisco Fellow

PhD in Physics (Chaos Theory)

Работает в Cisco в с 1998 года

Препринт статьи, на основе которой сделана данная презентация, доступен по адресу:

<https://arxiv.org/abs/1607.01639>

Название:

«Deciphering Malware's use of TLS (without Decryption)»

Авторы:

[Blake Anderson](#), [Subharthi Paul](#), [David McGrew](#)

# Вопросы