

Создание доверенных систем. Модели долгосрочного развития

Арутюн Аветисян
директор ИСП РАН
академик РАН
arut@ispras.ru

15 ноября 2022 г.



С.А. Лебедев



В.А. Мельников



БЭСМ-6 в музее науки Лондона



В.П. Иванников



Л.Н. Королёв

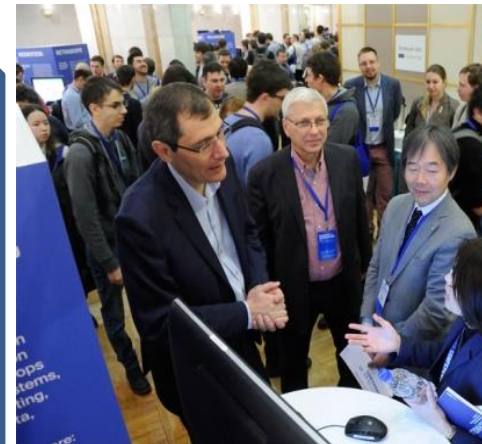


«Если мы глубже разберёмся в этом эпохальном советском суперкомпьютере, это позволит пересмотреть заявления времён холодной войны об отставании русской технологии, а также подтвердить или развеять мифы о технологическом совершенстве наших союзников».

Doron Swade, senior curator of computing and information technology

2018: 70 лет IT*
2019: 25 лет ИСП РАН
2023: 75 лет IT*

*** в России и странах
постсоветского пространства**



- ✓ **Индустриальные исследования и внедрение научных разработок (анализ данных, анализ программ)**
 - ✓ **Кафедры системного программирования на факультетах ВМК МГУ, ФУПМ МФТИ, ФКН ВШЭ**
 - ✓ **Три лаборатории за пределами ИСП РАН**

Кибербезопасность: вызовы и решения

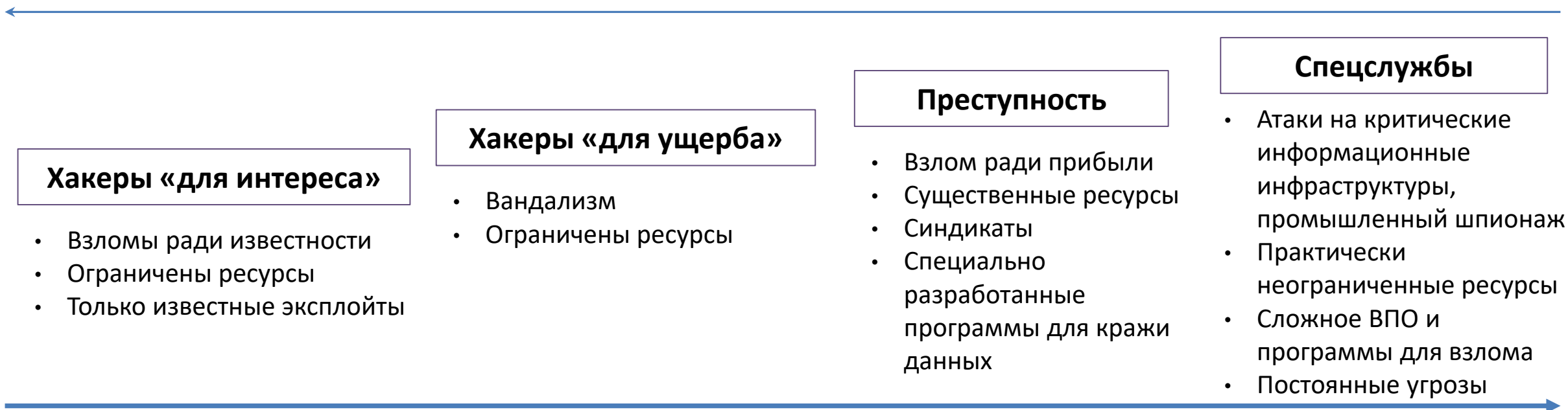
Ошибки в ПО – основная причина уязвимости систем

Принципиальное наличие **уязвимостей*** в ПО и аппаратуре:

функциональные, архитектурные, программного кода/микрокода.

**Границы между ошибками программиста, закладками и НДВ размыты*

Утечка информации о уязвимостях



Рост ресурсов и сложности атак

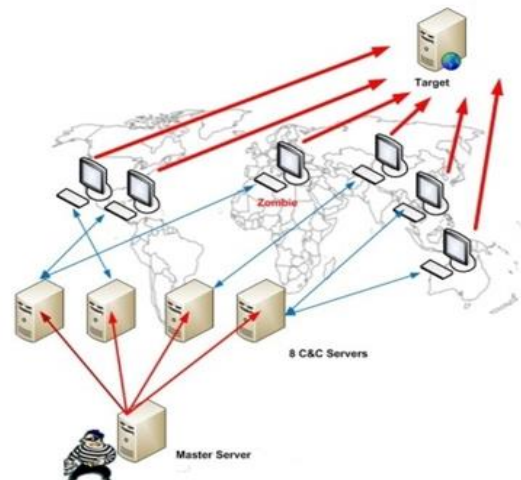
Эскалация размеров и сложности (миллиарды строк кода на GitHub), сложность среды разработки и сборки, облачные платформы и интернет вещей



Аварии на критических объектах



Перехват управления



Бот-нет



Кража паролей



Уязвимости



Кража информации о кредитных картах

Компьютерные атаки осуществляются путём эксплуатации дефектов в ПО и аппаратуре

```
void f(char * p)
{
    char s[6];
    strcpy(s, p);
}
```

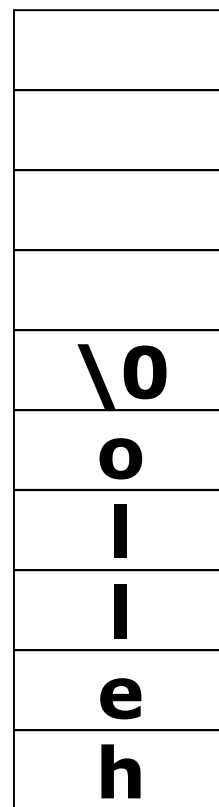
```
void main1 ()
{
    f("hello");
}

void main2 ()
{
    f("privet");
}
```

В случае main2 адрес возврата перезапишется, и управление будет передано не на main2, а на другой участок кода

Стек после выполнения функции f, вызванной из main1

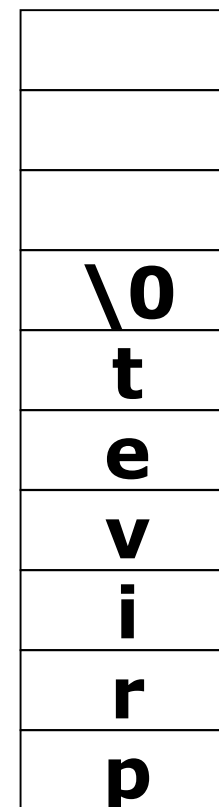
Адрес main1



адрес возврата


массив s

Стек после выполнения функции f, вызванной из main2



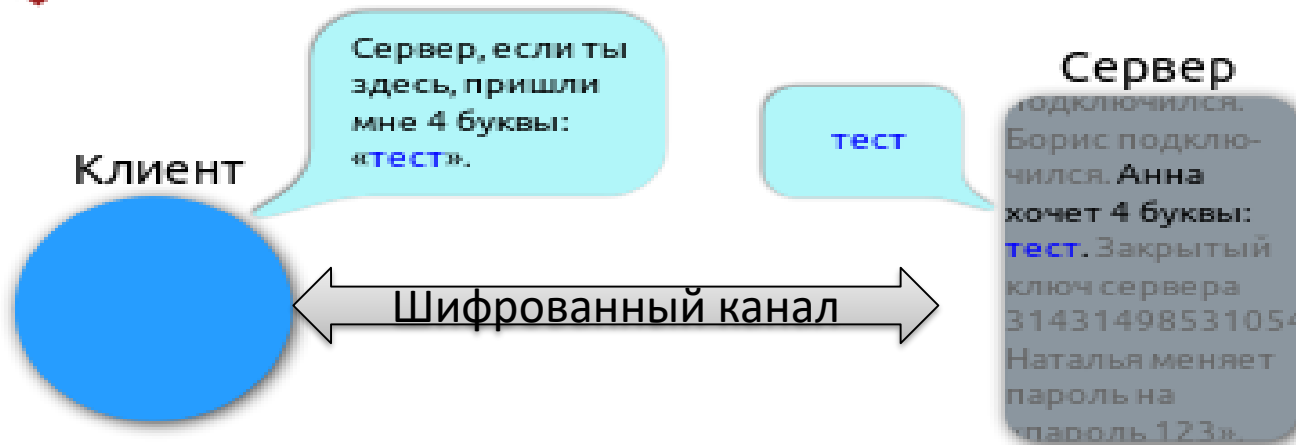
На место адреса main2 записали лишний байт

```
bool auth() {  
    char buf[N];  
    bool res;  
  
    read_password(buf, N);  
    res = check_password(buf);  
  
    memset(buf, 0, N);  
    return res;  
}
```

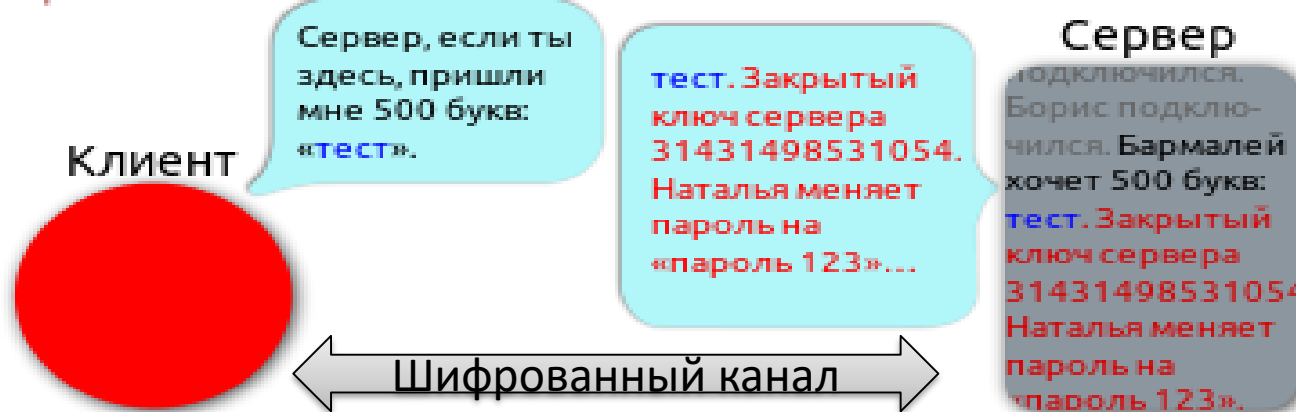


Компилятор удаляет
обнуление буфера с паролем,
т.к. с его точки зрения после
обнуления буфер не
используется.
При этом пароль останется на
стеке.

♥ Heartbeat — нормальная работа



♥ Heartbleed — эксплуатация ошибки



500000 сайтов заражено
\$500 млн потерь

- Ошибка чтения данных за границей буфера: злоумышленник контролирует длину посланного текста
- Происходит утечка пользовательских данных
- Весь обмен данными строго следует зашифрованному протоколу

Типы ошибки: слабость кодирования обработки входных данных, переполнение буфера



Прежде чем достичь места реализации ошибки, введённые извне данные «проходят» по многим функциям разных модулей

Модуль с функцией считывания файла-архива

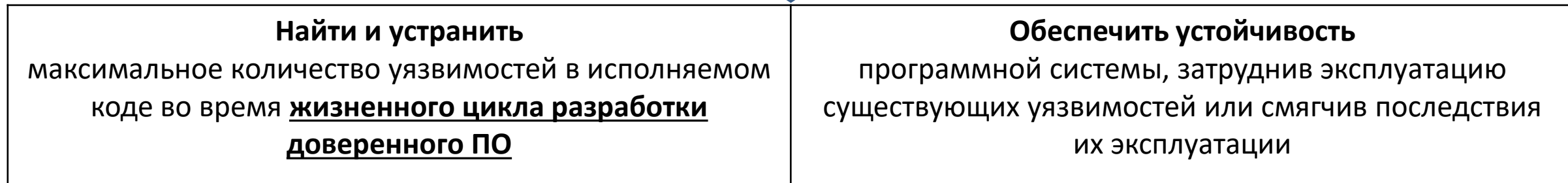
```
[tainted] Call of read
95     if(read(fd, file_hdr, sizeof_newlhd) != sizeof_newlhd) {
96         free(file_hdr);
97         return NULL;
98     }
99     file_hdr->flags = unrar_endian_convert_16(file_hdr->flags);
100    file_hdr->head_size = unrar_endian_convert_16(file_hdr->head_size);
101    file_hdr->pack_size = unrar_endian_convert_32(file_hdr->pack_size);
102    file_hdr->unpack_size = unrar_endian_convert_32(file_hdr->unpack_size);
103    file_hdr->file_crc = unrar_endian_convert_32(file_hdr->file_crc);
104    Composite 'file_hdr' taints element 'file_hdr->name_size'
105    file_hdr->name_size = unrar_endian_convert_16(file_hdr->name_size);
106    if(file_hdr->flags & 0x100) {
116    return file_hdr;
```

Функция в другом модуле. Ранее считанные извне данные определяют размер копируемой памяти

```
1484    /* Enter response type, length and copy payload */
1485    *bp++ = TLS1_HB_RESPONSE;
1486    s2n(payload, bp);
1487    Tainted data from /home/shimnik/openssl/ssl/s3_pkt.c+239 reached a sink.
1488    8. [SINK] *(s->s3->rrec.data + @) reaches the sink
1489    memcpy(bp, pl, payload);
1490    bp += payload;
1491    /* Random padding */
1492    RAND_pseudo_byte(paddi, 3 + payload + paddi);
1493    r = dtls1_write(ssl, buffer, 3 + payload + paddi);
```

От обучения до реагирования

- Недостаточно использовать классические методы защиты (защита по периметру, проверка доступа, антивирусы и др.)
- Необходима разработка новых моделей, методов и технологий в области анализа и трансформации программ



США (NIST и другие)

NIST: Национальный институт стандартов и технологий

С 1999: развитие госстандартов Common Criteria

С 2004: появление Microsoft Security Development Lifecycle

Инструменты анализа выбираются
сертификационными лабораториями

Разработчики ПО обязаны пользоваться
инструментами требуемого уровня. Заложено
непрерывное обновление стандартов по мере
развития технологий

Принятие нормативных документов привело к
взрывному росту рынка технологий анализа

Россия (ФСТЭК и другие)

**ГОСТ Р 56939-2016 «Защита информации. Разработка
безопасного программного обеспечения. Общие
требования»**

**Методика выявления уязвимостей и недеklarированных
возможностей в программном обеспечении**

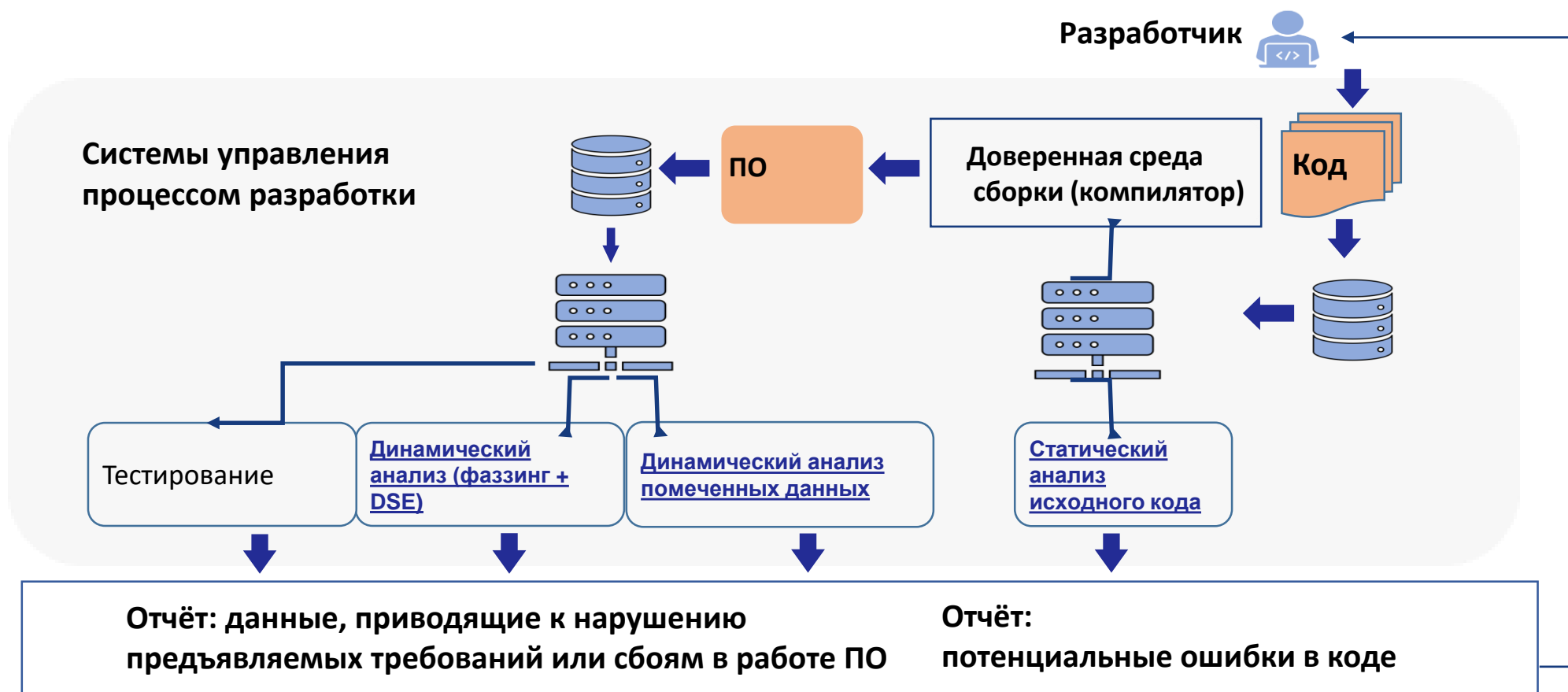
В разработке:

ГОСТ Р «Защита информации. Разработка безопасного
программного обеспечения. **Руководство по проведению
статического анализа. Общие требования**»

ГОСТ Р «Защита информации. Разработка безопасного
программного обеспечения. **Руководство по проведению
динамического анализа. Общие требования**»


ГОСТ Р «Защита информации. Разработка безопасного
программного обеспечения. **Безопасный компилятор
языков Си/Си++. Общие требования**»

...



!! Проблемы: сложные инструменты, траты ресурсов на анализ миллионов строк кода

<p>Статический анализ исходного кода</p>	<p><u>Svace (ИСП РАН, Россия)</u></p> 	<p>Аналоги: Klocwork (Perforce, США), Coverity (Synopsys, США), Fortify (MicroFocus (ранее HP), США, Великобритания). Открытые: Clang Static Analyzer, SpotBugs</p>
<p>Фаззинг + DSE (динамический анализ)</p>	<p><u>ИСП Crusher (ИСП РАН, Россия)</u></p> 	<p>Аналоги, недоступные в РФ: Peach Fuzzer (США, Peach Tech), Synopsys Defensics (Synopsys, США), MAYHEM (ForAllSecure, США). Открытые: angr (США), American Fuzzy Lop (США), Driller (США)</p>
<p>Динамический анализ помеченных данных</p>	<p><u>Блесна (ИСП РАН, Россия)</u></p> <p>Предназначена для испытательных лабораторий, органов по сертификации и др.</p>	<p>Аналоги, недоступные в РФ: MAYHEM (ForAllSecure, США), TETRANE's Timeless Analysis (Tetrane, США), APAC (2013-2015), VET (2013), CGC (2016), CHESS (2019) (проекты DARPA)</p>

Безопасная компиляция (не добавляющая ошибки в бинарный код)	<p><u>SafeC (ИСП РАН, Россия)</u></p> 	<p>Аналоги: отдельные части в GCC и Clang (флаги, санитайзеры), исследовательские или ограниченно применимые компиляторы и диверсификаторы (CompCert C, kcc, Selfrando, Oхumoron, Shuffler и др.)</p>
Определение поверхности атаки	<p><u>Natch (ИСП РАН, Россия)</u></p> <p>Предназначен для разработчиков безопасного ПО и для испытательных лабораторий</p> <p>Распространяется в сообществе</p>	<p>Аналоги: фреймворки Panda, Decaf (нет средств интроспекции), инструмент Panorama (исследовательский проект для Windows)</p>

Технологии ИСП РАН используют более 100 компаний в России и за рубежом:
Samsung, Huawei, «Лаборатория Касперского», «РусБИТех», Postgres Professional и др.

Статический анализ: аннотации служебных функций ОС (выделение и освобождение ресурсов, ввод-вывод и пр.), анализ вызовов по указателю, восстановление графа вызовов

Динамический анализ: использование инструментов, предназначенных для ядра ОС (фаззер syzkaller)

- Специальная сборка ядра с поддержкой инструментации и санитайзеров
- Подготовка списка системных вызовов, для фаззинга которых будет генерироваться тестовая программа

Возможно использование гипервизора (создание снимка состояния ОС – и запуска ОС на новых данных через гипервизор, начиная с этого снимка)

Статический анализ: использование контролируемой сборки через кросс-компиляторы на серверной хост-системе, возможность совместного анализа кода всей системы вместе

Динамический анализ: использование подходов частичной эмуляции (выполнение тестируемой программы через QEMU) или удаленной инструментации (фаззинг на хост-системе, легковесная инструментация на целевой встраиваемой системе)

- Требуется снятие дампа памяти тестируемой программы в эмуляторе
- Требуется эмуляция системных вызовов и периферии

Требуется анализ зависимостей для поиска уязвимых компонент в многокомпонентной программе

Динамический анализ: фаззинг входных интерфейсов (базы данных: JSON-формат, выполнение и обработка SQL-запросов)

- Нашли зависание БД (20+ минут) на определенном запросе
- Нашли переполнение буфера при фаззинге XML-парсера
- Нашли целочисленное переполнение в библиотеке FreeImage, приводящее к отказу в обслуживании

Эффективный фаззинг внутренних интерфейсов требует написания функций-оберток, которые непосредственно передают данные от фаззера внутрь программы

- Такой обертке нужно передавать сложный контекст выполнения, который нужно инициализировать
- Для написания обертки нужен квалифицированный разработчик

Статический анализ: написание аннотаций библиотек, настройка критичности предупреждения (например, для языка Java некоторые предупреждения имеют меньшую критичность)

Динамический анализ: инструментация интерпретируемого или компилируемого байткода и перехват исключений интерпретатора

- Фаззинг программ, выполняемых JIT-компилятором или виртуальной машиной, может потребовать их модификации

2018: принято решение Президиума РАН о новом научном направлении «Анализ, трансформация программ и кибербезопасность»*

2021: новая специальность ВАК «Кибербезопасность» утверждена Приказом Минобрнауки №118

Направления исследований:

- Анализ и систематизация уязвимостей
- Моделирование политик информационной безопасности, угроз и атак
- Методы, алгоритмы и средства пост-релизного глубокого анализа защищенности ПО
- Методы интеграции средств защиты на уровне аппаратуры и на уровне ПО
- Интеллектуальный масштабируемый мониторинг инцидентов безопасности в распределенных программно-аппаратных системах
- Масштабируемые средства интеллектуального анализа данных и процессов в распределенных системах

И другие

*По предложению ИСП РАН; поддержано ФСТЭК России, АО «Лаборатория Касперского», ОАО «РусБИТех».

Искусственный интеллект: НОВЫЕ ВЫЗОВЫ В КИБЕРБЕЗОПАСНОСТИ



Компьютер выполняет «интеллектуальные» действия, которые до этого были доступны только человеку:

- Играет в шахматы
- Ведет простой диалог

Решения основанные на правилах

Искусственный интеллект

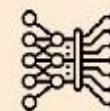


Машинное обучение

Решения основанные на данных:

- **Этап обучения:**
алгоритм получает примеры вопросов и правильных ответов
- **Этап применения:**
алгоритм получает вопрос и выдает ответ

Глубокое обучение



Иерархия алгоритмов машинного обучения: ответы одних попадают в вход другим

Модель нейронной сети

1950

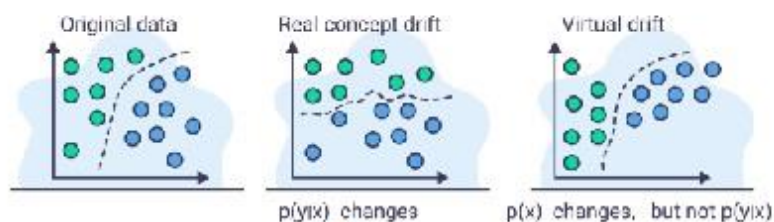
1980

2010

Переход от построения **модели** объекта автоматизации к решению проблемы **по аналогии**

Проблемы машинного обучения

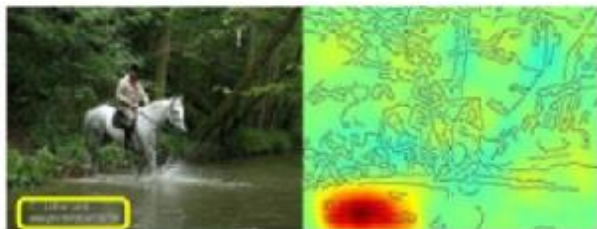
Проблемы разработки и эксплуатации
(переобучение, сдвиги в данных в течение
эксплуатации)



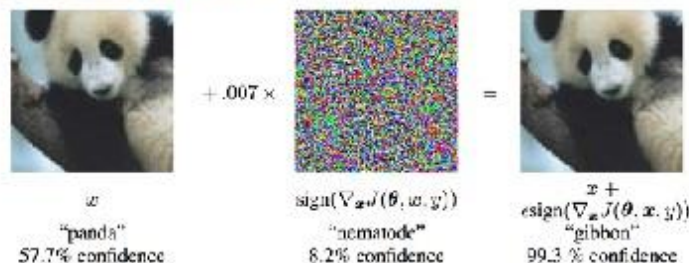
Предвзятость моделей



Интерпретация результатов



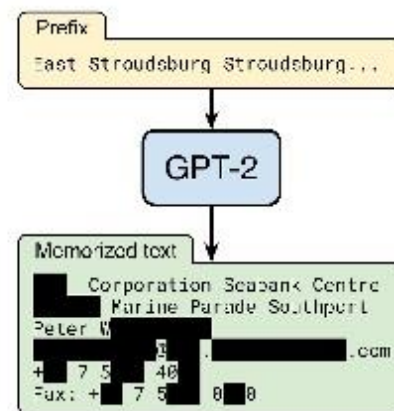
Состязательные атаки для манипуляции
поведением системы



Встраивание закладок на этапе обучения



Извлечение конфиденциальных данных из
обученных



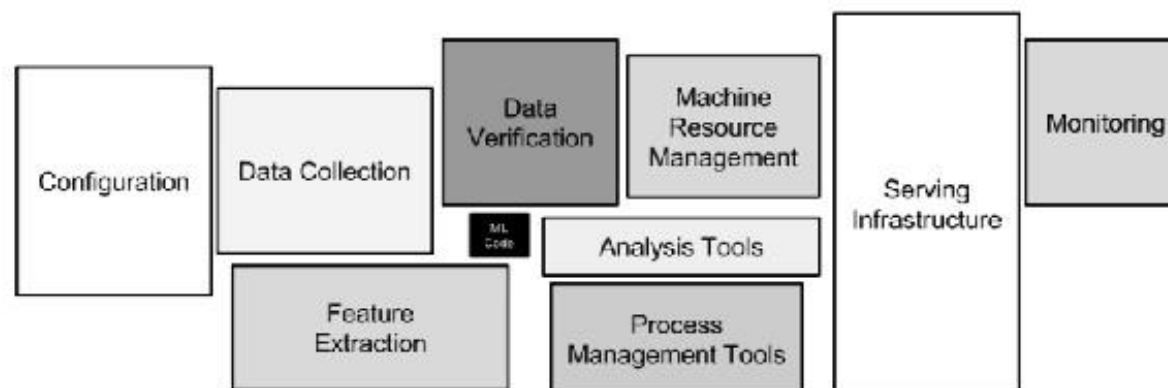
«Кража» самих моделей

Проблемы создания доверенных систем, использующих ТИИ, уже активно рассматриваются:

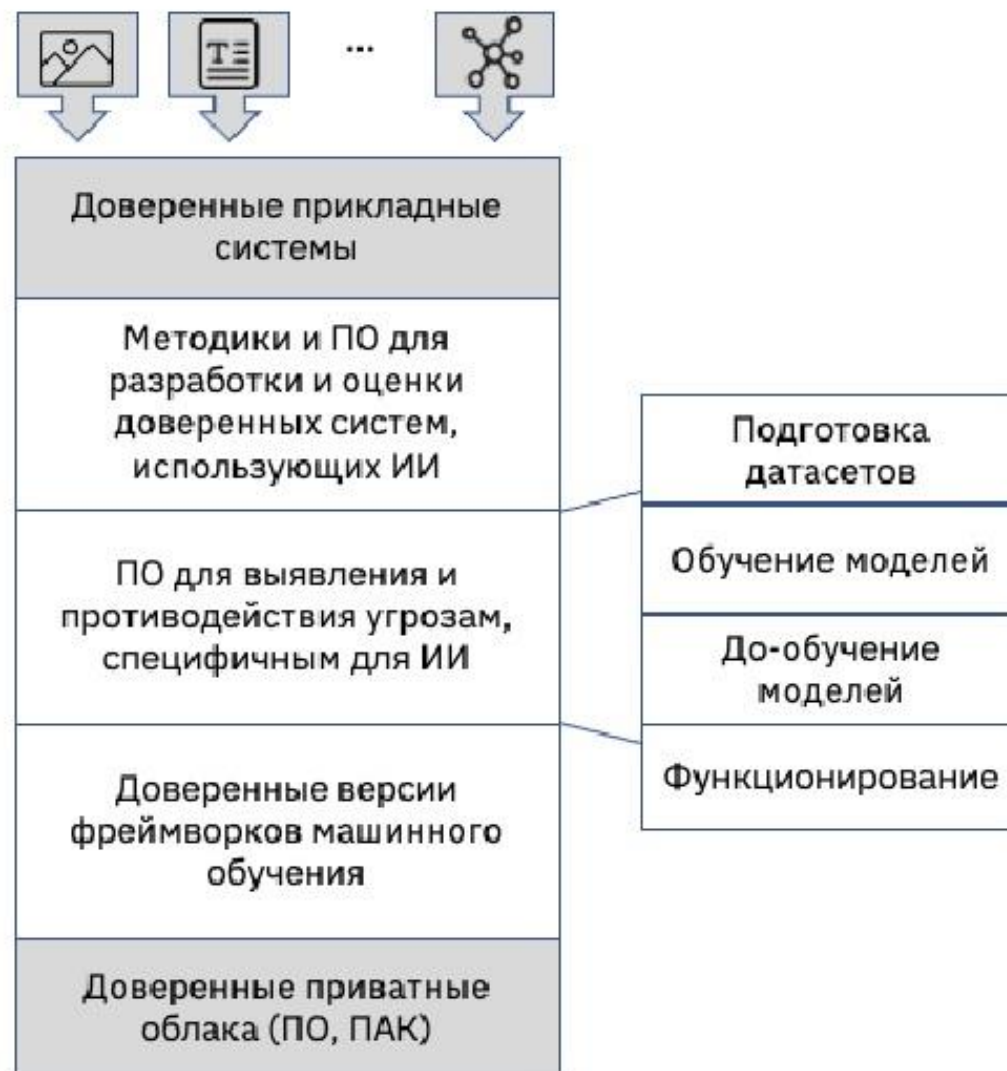
- NIST AI Risk Management Framework (США)
 - DIN DKE German Standardization Map on Artificial Intelligence (Германия)
 - MITRE ATLAS, Adversarial Threat Landscape for Artificial-Intelligence Systems
 - Google Responsible AI practices
-

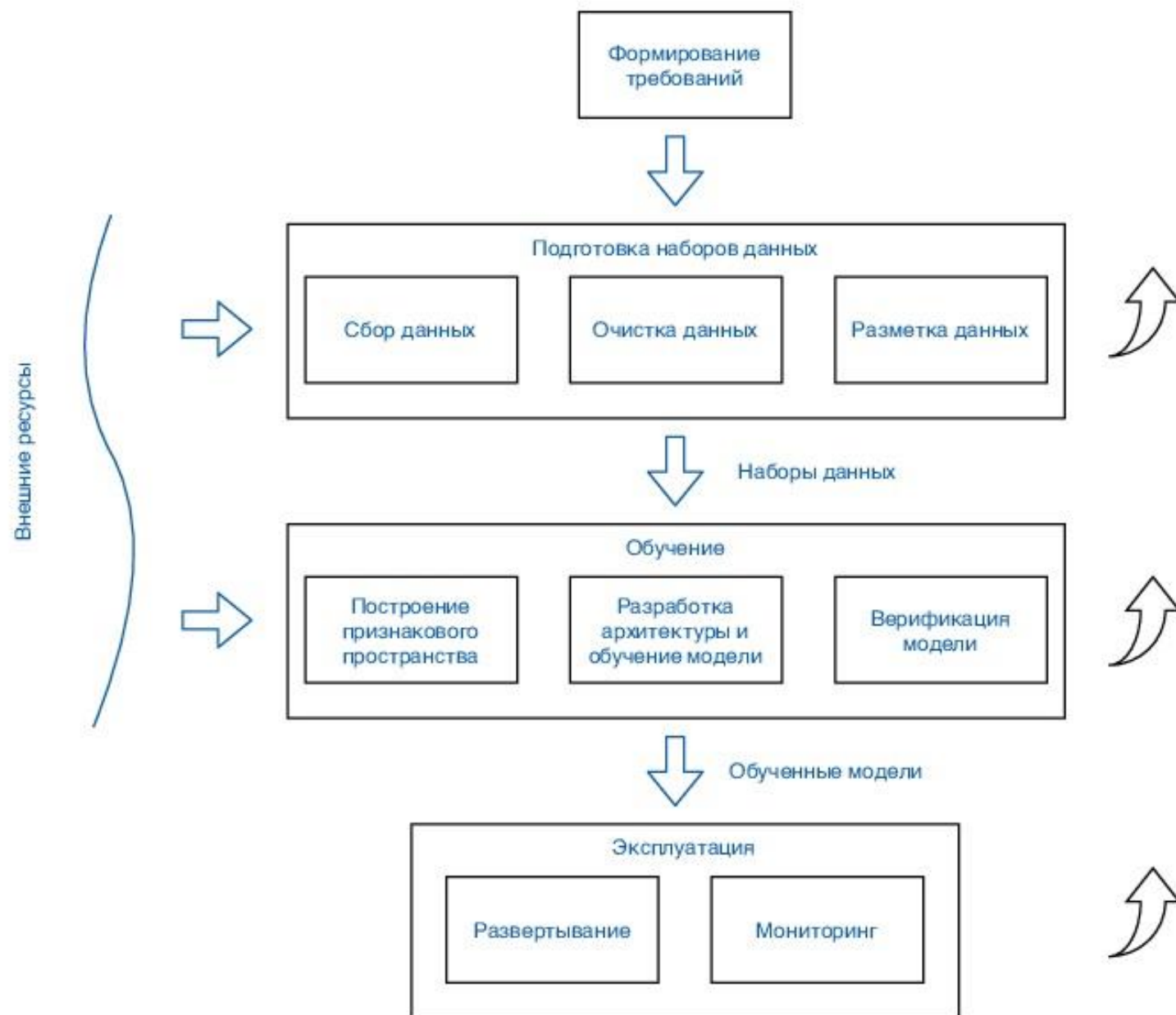
Историческая аналогия

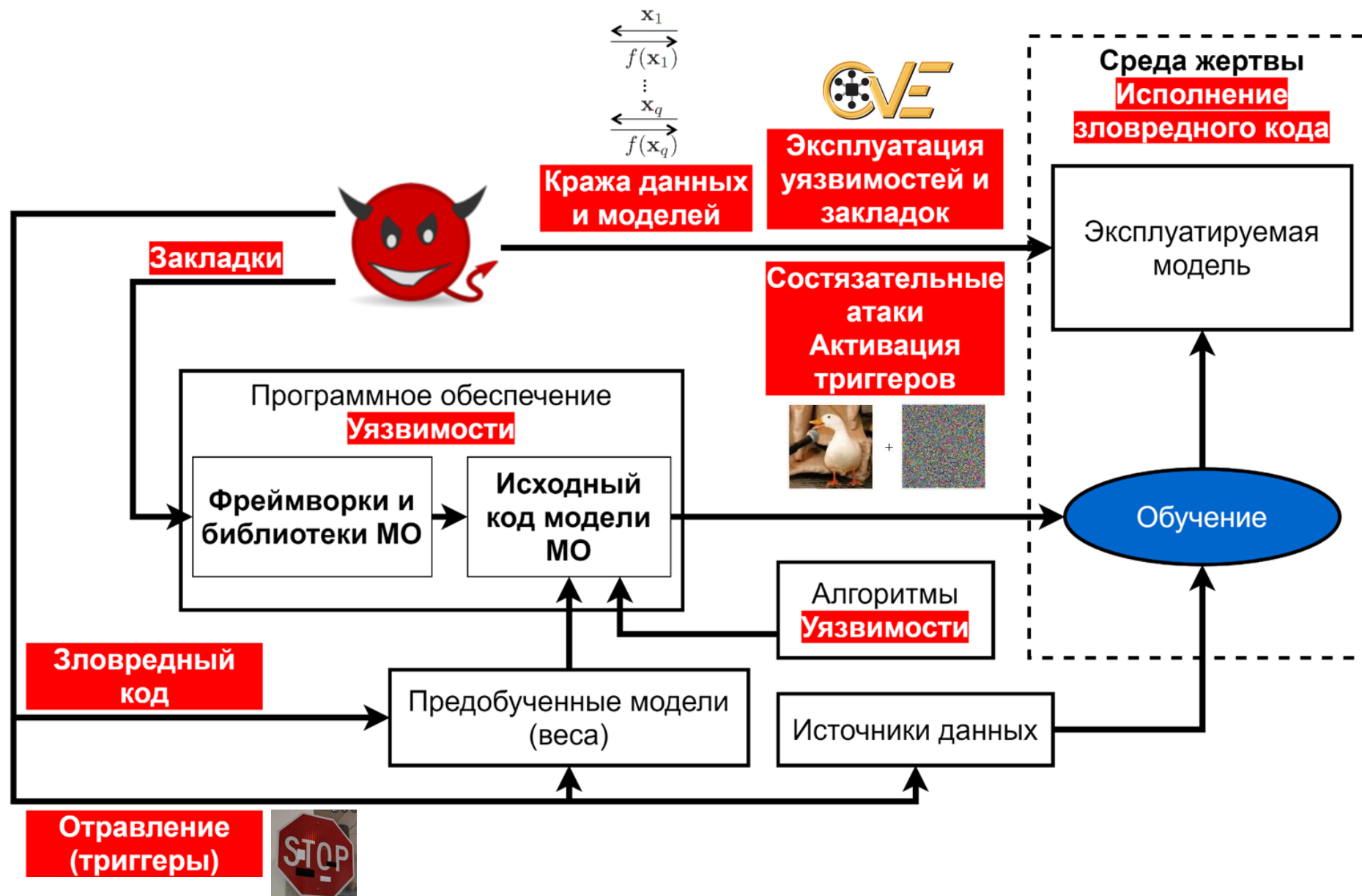
- 1999: развитие госстандартов Common Criteria
- 2004: появление Microsoft Security Development Lifecycle
- Продукты не соответствующие требованиям не допускаются к внедрению в госсекторе США и в крупных международных компаниях



- Модели машинного обучения — центральная, но не самая большая часть интеллектуальных систем
- Понятие «доверие» к программным системам определяется национальными стандартами (ГОСТ Р 56939-2016, приказ ФСТЭК №76 от 2 июня 2020 г.)
- Принципиальное отличие «интеллектуальных систем» - информация содержится не в программном коде, а в данных
- Для обеспечения доверия к интеллектуальным системам **отсутствует научно-технологическая база**

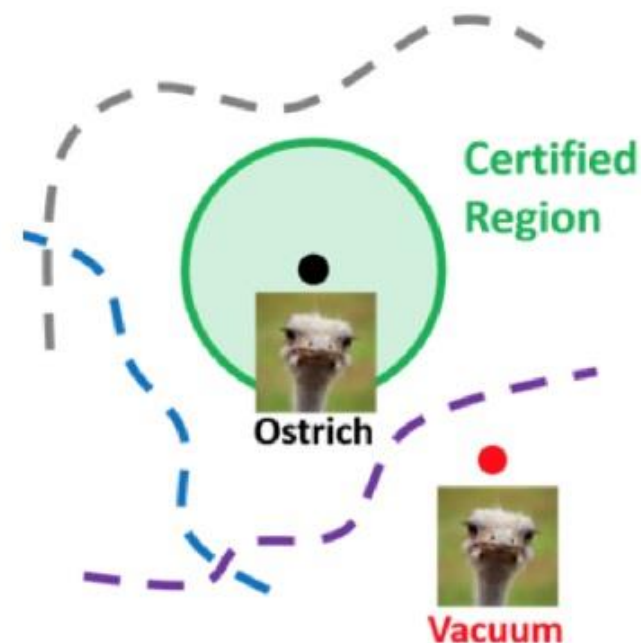


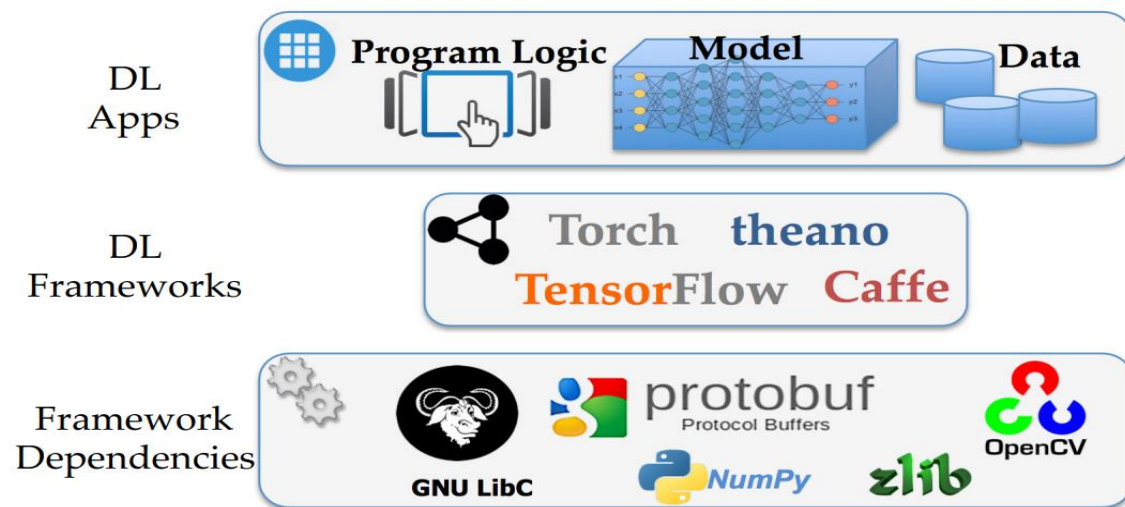




- Модели, устойчивые к заданному типу атак
 - Сетифицированная робастность
 - Составительное обучение
- Выбор алгоритма обучения, для поиска устойчивых решений

CNN-Cert finds a
certified region of
robustness





Xiao, Q., Li, K., Zhang, D., & Xu, W. (2018). Security Risks in Deep Learning Implementations. 2018 IEEE Security and Privacy Workshops (SPW), 123-128.

- Классические уязвимости (CVE) в исходном коде фреймворков и библиотек расширяют поверхность атаки на эксплуатируемые модели
- Пример: переполнение буфера в библиотеке OpenCV, атака с помощью специально подготовленного BMP-изображения
- Кроме того, в исходном коде возможны **закладки**
- Необходимы **статический анализ** исходного кода фреймворков и создание их **доверенных версий**

7 миллионов строк кода  PyTorch


TensorFlow

10 миллионов строк кода

! ЭТО ЛИШЬ ЧАСТЬ СТЕКА ПО АНАЛИЗУ ДАННЫХ, НА КОТОРОМ ОСНОВАНЫ ТЕХНОЛОГИИ ИИ !

• WARNING

`torch.load()` uses `pickle` module implicitly, which is known to be insecure. It is possible to construct malicious pickle data which will execute arbitrary code during unpickling. Never load data that could have come from an untrusted source, or that could have been tampered with. **Only load data you trust.**

- Загрузка моделей из недоверенных источников (Hugging Face)
- PyTorch использует Pickle
- Pickle позволяет сохранять и загружать **любые** объекты Python
- То есть возможно дописать к сохраненной модели произвольную функцию, и вызвать ее через `eval` при загрузке

<https://youtu.be/2ethDz9KnLk>

- ✓ Внедрение технологий ИИ без технологий обеспечения доверия невозможно
 - ✓ В одиночку обеспечить доверие нельзя

Модели развития доверенного системного ПО: что есть и что нужно?



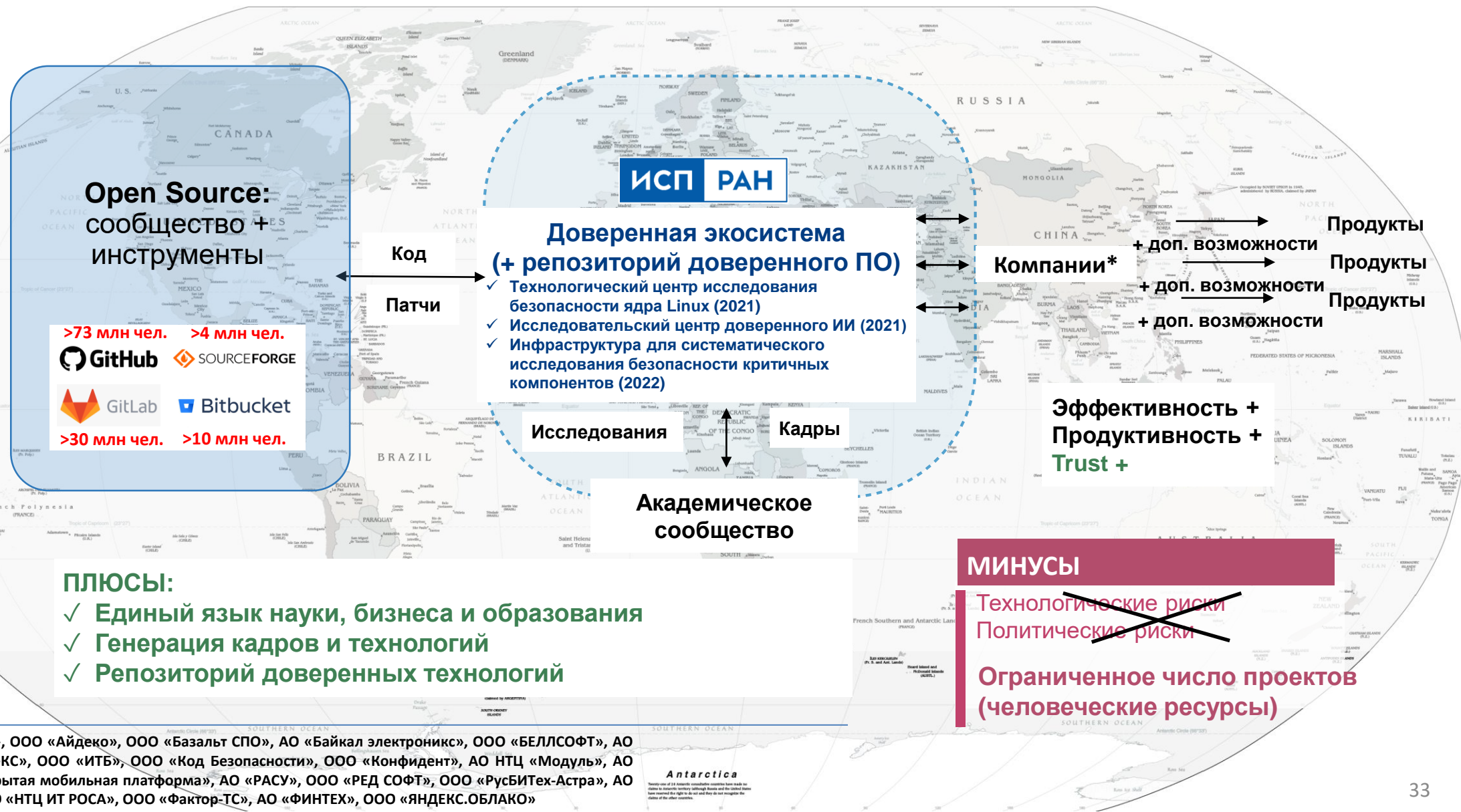
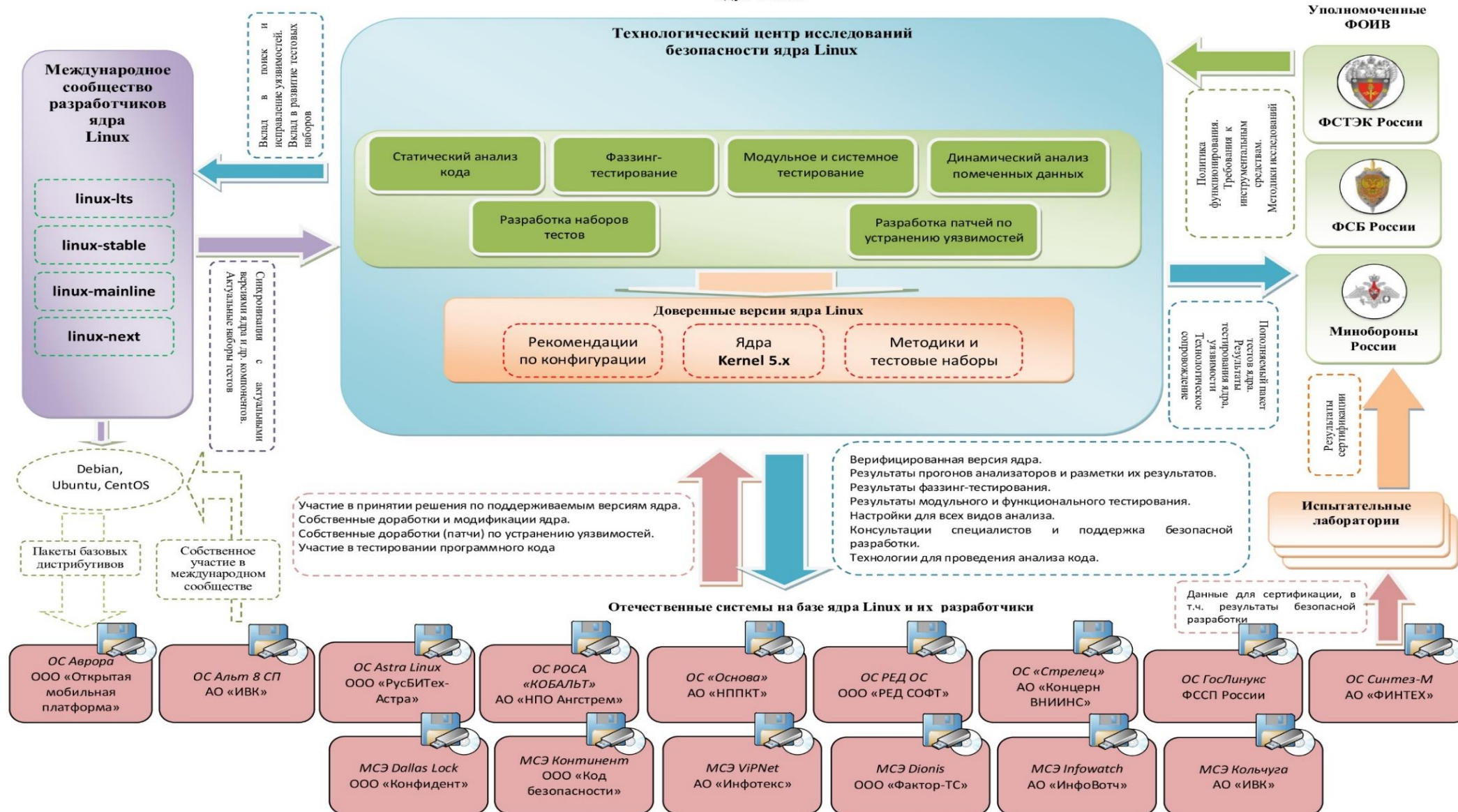


Схема
организации работ по исследованиям безопасности операционных систем на базе
ядра Linux



- Началось сопровождение ветки ядра Linux, основанной на стабильной версии 5.10.
- Подготовлены методики проведения исследования ядра, включая:
 - ✓ статический анализ при помощи инструмента Svasc;
 - ✓ системное и модульное тестирование (наполнение тестовыми наборами продолжается);
 - ✓ фаззинг-тестирование при помощи инструмента syzkaller;
 - ✓ проведение архитектурного анализа с целью определения поверхности атаки;
 - ✓ проведение анализа помеченных (чувствительных) данных.
- Создана экспертная группа из представителей 21 компании, в рамках которой:
 - ✓ формируются принципы функционирования Технологического центра;
 - ✓ проведена разметка более 6000 предупреждений инструмента Svasc;
 - ✓ подготовлено ~100 исправлений ошибок в ядре; 53 патча уже приняты в основную ветку ядра.
- Ведётся подготовка рекомендаций по конфигурированию ядра с целью повышения его безопасности.
- Ведётся разработка улучшений в ядре, нацеленных на повышение безопасности его работы на стадии развёртывания и инициализации.

ПАРТНЕРЫ:

АО «Аладдин Р.Д.»
ООО «Айдеко»
ООО «Базальт СПО»
АО «Байкал электроникс»
ООО «БЕЛЛСОФТ»
АО «ИВК»
АО «ИнфоТеКС»
ООО «ИТБ»
ООО «Код Безопасности»
ООО «Конфидент»
АО НТЦ «Модуль»
АО «МЦСТ»
ООО «ОМП»
АО «РАСУ»
ООО «РЕД СОФТ»
ООО «РусБИТех-Астра»
АО МВП «Свемел»
ООО «НТЦ ИТ РОСА»
ООО «Фактор-ТС»
АО «ФИНТЕХ»
ООО «ЯНДЕКС.ОБЛАКО»

❑ Создана аппаратно-программная инфраструктура для обеспечения доверия к базовым фреймворкам машинного обучения (TensorFlow, PyTorch)

- ✓ Первичный анализ исходного кода
- ✓ Постоянная проверка новых изменений кода
- ✓ Синхронизация с оригинальными открытыми версиями

❑ Разработаны образцы доверенных фреймворков машинного обучения. Тестирование осуществляется индустриальными партнерами центра

ПОДРОБНОСТИ

ЗАДАЧА

Создать доверенные версии фреймворков Tensorflow и PyTorch

РЕШЕНИЕ


Проведён анализ поверхности атаки. Перспективным вектором атаки был выбран компонент загрузки моделей. Кроме того, для фреймворка TensorFlow были восстановлены фаззинг-цели в проекте OSS-Fuzz. Фаззинг фреймворков производился с помощью инструмента **Sydr** (ИСП РАН). Исходный код был проанализирован инструментом **Svace** (ИСП РАН).


РЕЗУЛЬТАТ

В результате фаззинга обнаружено 7 ошибок для PyTorch и одна ошибка для TensorFlow.

В результате анализа исходного кода обнаружено 13 ошибок в PyTorch и 8 ошибок в TensorFlow.


← → ↻ github.com/tensorflow/tensorflow/pull/57892


 Open Fix bugs found by static analysis #57892
apach301 wants to merge 2 commits into tensorflow:master from apach301:static-analysis-bugs


 mihairuseac commented on Sep 28 Collaborator ...

CC @pak-laura @learning-to-play @izuk

@apach301 can you describe the static analysis framework you have used? Thank you very much for the PR

 gbaned added this to Assigned Reviewer in PR Queue via automation on Sep 29

 gbaned removed the awaiting review label on Sep 29

 apach301 commented on Sep 30 Author ...

CC @pak-laura @learning-to-play @izuk

@apach301 can you describe the static analysis framework you have used? Thank you very much for the PR

@mihairuseac We used [Svace](#) static analysis tool. It integrates into the project build system, creates IR and performs a number of checks. More details in [paper](#) and [slides](#).

Ошибки, найденные Svace

PyTorch:

5 поданы и приняты

6 поданы и рассматриваются

TensorFlow:

6 поданы и рассматриваются

Ошибки, найденные фаззингом:

PyTorch:


5 поданы и приняты

3 поданы и рассматриваются


TensorFlow:

1 подана и принята

Fix endless loop in DecodeLin16WaveAsFloatVector #56455

 Merged copybara-service merged 1 commit into tensorflow:master from anfedotoff:master on Jun 28

 Conversation 2  Commits 1  Checks 4  Files changed 2

 anfedotoff commented on Jun 14 Contributor ...

Hi!

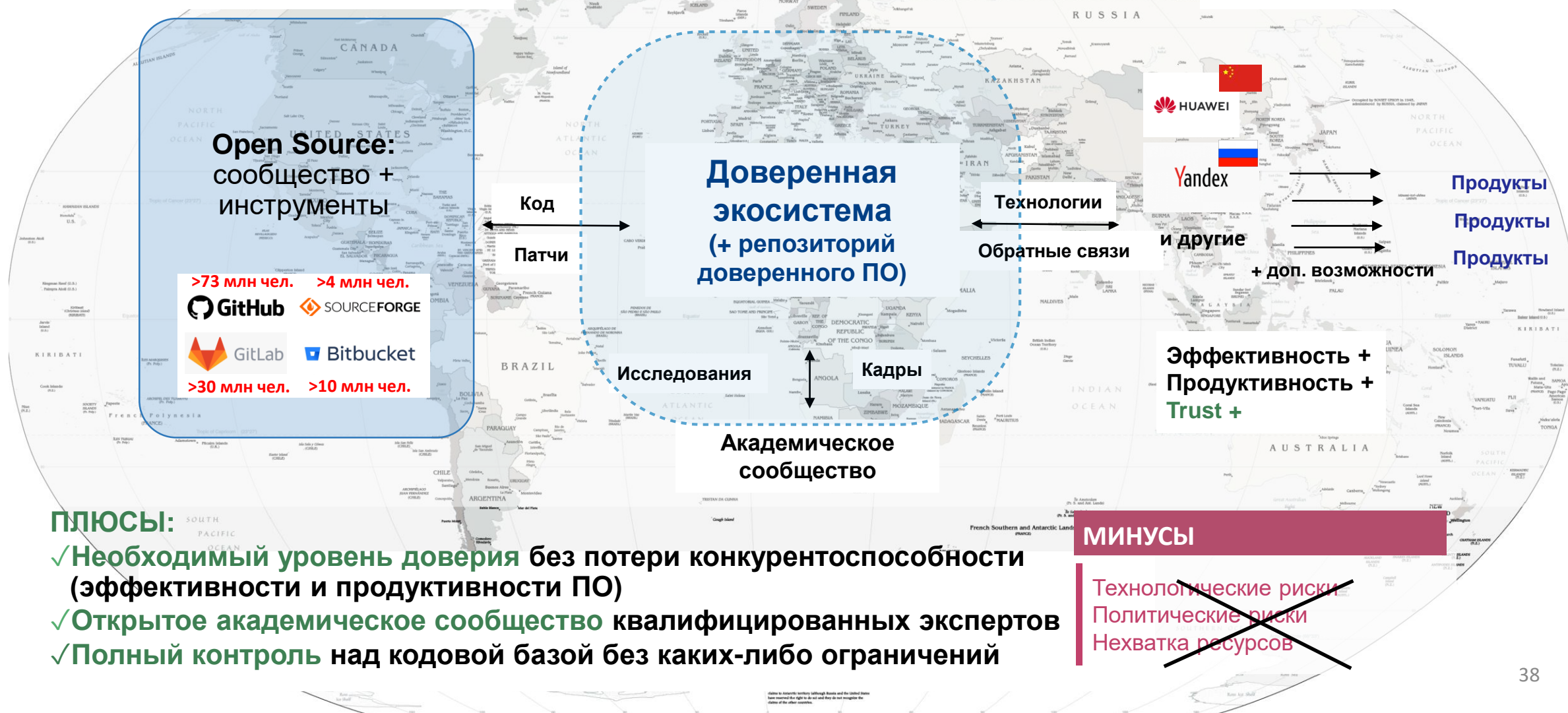
We were doing some fuzzing using libFuzzer and symbolic execution tool Sydr for fuzz targets from this [PR](#). And we found an interesting issue. Loop at [wav_io.cc:233](#) could be endless under certain conditions. I'll try to explain: if variable `found_text` is equal for one of keywords from [if at wav_io.cc:240](#) (for example `found_text = "bext"`) and the next value of `size_of_chunk` is equal to -8 in sign representation there will be an endless loop.

Глобальный вызов:

Долгосрочное устойчивое развитие
доверенного системного ПО

Глобальная цель:

Технологическая
независимость для всех





Спасибо за внимание!

Арутюн Аветисян

arut@ispras.ru

