# CS 489 / 698: Software and Systems Security

**Module: Common Vulnerabilities**
Lecture: weird machine

Meng Xu *(University of Waterloo)*

Fall 2024

# Outline

**Introduction**
○●○○○○○

State machine
○○○○○○○○○○○○○

Security
○○○○○○○○○○○○○

## Based on paper

**Weird Machines, Exploitability, and Provable Unexploitability**

By *Thomas Dullien* published in 2017 when he was in Google
Project Zero.

Introduction
○○●○○○○

State machine
○○○○○○○○○○○○○

Security
○○○○○○○○○○○○○

# Why this paper?

**Introduction**
○○○○○○○

State machine
○○○○○○○○○○○○○

Security
○○○○○○○○○○○○○

## Why this paper?

It attempts to formalize a concept that has been intuitively known for quite a while in the community of security practitioners, i.e., both by the hackers and the researchers...

**Introduction**
○○○●○○○○

State machine
○○○○○○○○○○○○

Security
○○○○○○○○○○○○

# Why this paper?

It attempts to formalize a concept that has been intuitively known for quite a while in the community of security practitioners, i.e., both by the hackers and the researchers...

... and that concept is called "exploit".

# What is an exploit?

## What is an exploit?

- Magic
- Access (mostly unauthorized)
- Controls the instruction pointer (e.g., EIP/RIP register)
- A program does something it is not supposed to do
- I can recognize it when I see it

They are not technically wrong, but are clearly ill-defined for
academic research purposes.

Introduction
ooooooo

State machine
oooooooooooo

Security
oooooooooooo

## Why do we bother to define it?

## Why do we bother to define it?

We need to make justifications in the real-world that depends on the concept of "exploits":

- Mitigation strategies
  - e.g., difficulty of exploitation vs performance
  - e.g., difficulty of exploitation vs programmability
  - e.g., difficulty of exploitation vs complexity

- Exploitability of software/hardware defects
  - e.g., does the Rowhammer bug create a big security problem?
  - e.g., can the Spectre bug be used to launch general attacks? If yes, how?

**Introduction**
○○○○○●○

State machine
○○○○○○○○○○○○

Security
○○○○○○○○○○○○○

# The MitiGator

Raising the bar on exploitation until no more exploits can be seen
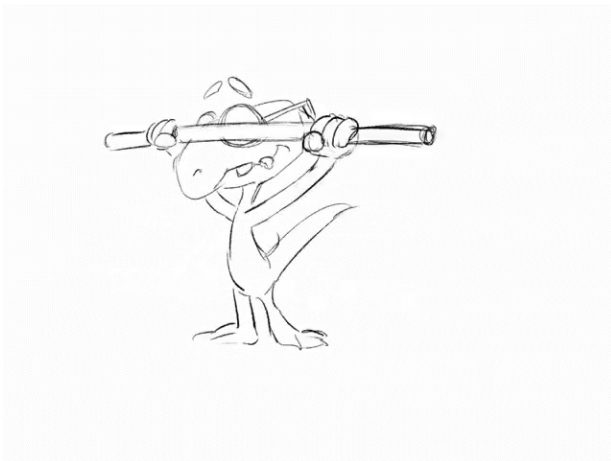


Copyright: The MitiGator animation

**Introduction**
ooooo●o

State machine
oooooooooooo

Security
ooooooooooooo

# The MitiGator

Raising the bar on exploitation until no more exploits can be seen



Copyright: The MitiGator animation

**Introduction**
○○○○○●○

State machine
○○○○○○○○○○○○

Security
○○○○○○○○○○○○○

## The MitiGator

Raising the bar on exploitation until no more exploits can be seen



Copyright: The MitiGator animation

**Introduction**
ooooooo●

State machine
oooooooooooooo

Security
ooooooooooooo

## Learn principles, not examples

**Introduction**
○○○○○○○●

State machine
○○○○○○○○○○○○○

Security
○○○○○○○○○○○○○

## Learn principles, not examples

An important message conveyed by this paper (which is also a message I want to share with you), is that exploitation IS NOT a "bag of tricks".

**Introduction**
ooooooo●

State machine
oooooooooooooo

Security
ooooooooooooo

## Learn principles, not examples

An important message conveyed by this paper (which is also a message I want to share with you), is that exploitation IS NOT a "bag of tricks".

In security courses (including this one), we teach
- Stack smashing, buffer overflows, heap exploitations
- SQL injection, XSS, etc
- ASLR, CFI, sandboxing, etc.

It is important to remember that there is a more fundamental principle behind these examples — *exploitation is all about entering and programming a weird machine*.

Introduction
0000000

State machine
●000000000000

Security
00000000000000

## Outline

## Behind an "exploit"

By just saying "I exploited something", you are conveying at least two messages:

- There exists some software running on top of some hardware
- There are "defects" in either the software or hardware (or both).

Introduction
0000000

State machine
0000000000000

Security
00000000000

What is software?

Introduction
○○○○○○○

State machine
○○●○○○○○○○○○○

Security
○○○○○○○○○○○○○

## What is software?

A software is an emulator for a finite-state machine (FSM) we would like to have but we don't.

Introduction
0000000

State machine
0000000000000

Security
00000000000000

## What is software?

A software is an emulator for a finite-state machine (FSM) we would like to have but we don't.

Instead, we only have a general-purpose CPU which is designed to model a huge spectrum of FSMs.

Introduction
0000000

State machine
000●000000000

Security
00000000000

## What is software?

A software is an emulator for a finite-state machine (FSM) we would like to have but we don't.

Instead, we only have a general-purpose CPU which is designed to model a huge spectrum of FSMs.

Hence, the reason we develop software is to confine the CPU to follow and only follow the FSM we intend to have.

# The intended finite-state machine (IFSM)

The state machine we want to have is called the "intended finite-state machine" (IFSM).

- It is usually not explicitly specified
- It is "perfect" by design — fully implements our intentions
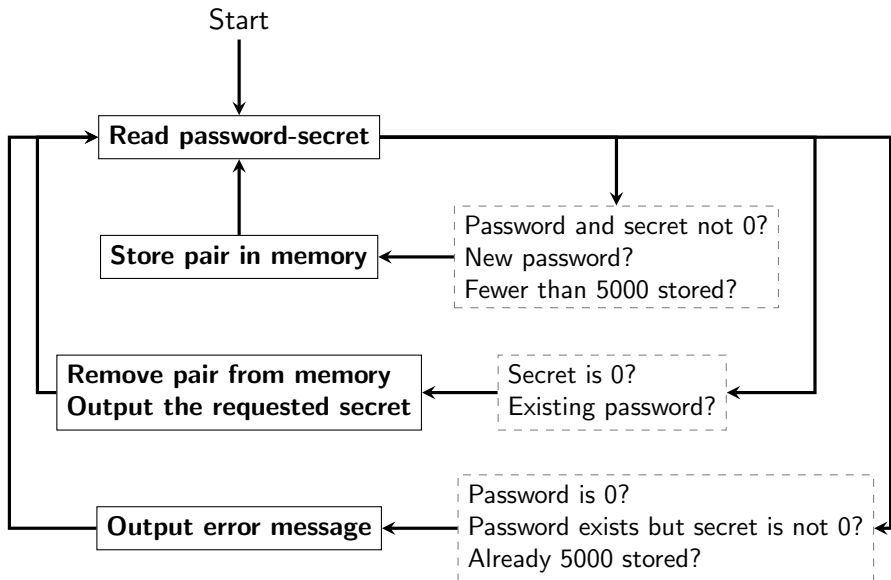- It cannot, by definition, have security problems.

Introduction
0000000

State machine
0000●00000000

Security
00000000000

# A concrete example: a secret-keeping machine

The machine has the following functionalities:

- Reads a password / secret $(p, s)$ from a user and remembers it.
  - NOTE: neither $p$ nor $s$ can be 0 (0 is reserved as an error code)
- Given a password $(p)$ that exists in the memory, the machine returns a previously-stored secret $(s)$ and forget both.
- The machine will not need to store more than 5000 such pairs.

Introduction
0000000

State machine
000000●0000000

Security
00000000000

## IFSM diagram

## IFSM diagram

Introduction
State machine
Security
ooooooo
oooooooo●ooooo
ooooooooooooo

## IFSM formalization

The set of all *Memory*, denoted as $\mathcal{M}$, can be formally defined as

$$
\mathcal{M} = \left\{ \begin{array}{l} \emptyset \\ \{(p_1, s_1)\} \\ ... \\ \{(p_1, s_1), ..., (p_{5000}, s_{5000})\} \end{array} \middle| \begin{array}{l} p_i, s_i \in \{0, 1\}^{32} - \{0\} \\ p_i \neq p_j \end{array} \right\}
$$

# FSM quick recap

An FSM can be defined by a 7-tuple: $(Q, i, F, \Sigma, \Delta, \delta, \sigma)$

- $Q$: Set of states
- $i$: The initial state
- $F$: The set of final states
- $\Sigma$: The input alphabet
- $\Delta$: The output alphabet
- $\delta$: State transition function $\delta : Q \times \Sigma \to Q$
- $\sigma$: Output mapping function $\sigma : Q \times \Sigma \to \Delta$

# IFSM formalization — what we intend to have

The IFSM of our secret-keeping program can be defined as:

- $Q$: $\{A_M, M \in \mathcal{M}\}$
- $i$: $A_\emptyset$
- $F$: $\emptyset$
- $\Sigma$: $\{(p, s) \mid p, s \in \{0, 1\}^{32}\}$
- $\Delta$: $\{0, 1\}^{32}$
- $\delta$: $A_M \times (p, s) \to A_M \mid A_{M \cup (p,s)} \mid A_{M-(p,s)}$
- $\sigma$: $A_M \times (p, s) \to s' \mid 0$

## What we actually have: a realistic CPU

The Cook-and-Reckhow RAM machine

- $2^{16}$ memory cells each holding a 32-bit value
- 7 CPU registers ($r_0$ to $r_6$)
- A small set of instructions
  - Constant: LOAD($C$, $r_d$)
  - Register operations: ADD($r_{s1}$, $r_{s2}$, $r_d$)
  - Register operations: SUB($r_{s1}$, $r_{s2}$, $r_d$)
  - Memory read: ICOPY($r_p$, $r_d$)
  - Memory write: DCOPY($r_d$, $r_s$)
  - Control flow: JNZ/JZ($r$, $l_z$)
  - Environment IO: READ($r_d$)
  - Environment IO: PRINT($r_s$)
- Harvard architecture (program is provided and external to RAM)

# CPU FSM formalization — what we actually have

The FSM of a general-purpose CPU can be defined as:

- $Q$: $(q_1, ..., q_{2^{16}}) \times (r_0, ..., r_6) \times p_i$ where $q_i, r_i \in \{0,1\}^{32}$, $p_i \in P$
- $i$: $q_i = 0, r_i = 0, p_i = P_0$
- $F$: $\emptyset$
- $\Sigma$: CPU Instruction Set $\{I\}$
- $\Delta$: $\{0,1\}^{32}$
- $\delta$: $Q \times I \to Q'$
- $\sigma$: $Q \times I \to (e \in \Delta)$

# From spec to execution: a series of refinement

We want to translate our IFSM $S_{spec}$ into our CPU FSM $S_{execution}$.

Introduction
0000000

State machine
00000000000●

Security
00000000000

## From spec to execution: a series of refinement

We want to translate our IFSM $S_{spec}$ into our CPU FSM $S_{execution}$.

It is actually a multi-stage process, involving (non-exhaustively)

$$S_{spec} \ \sqsupseteq \ S_{language} \ \sqsupseteq \ S_{machine} \ \sqsupseteq \ S_{execution}$$

Introduction
0000000

State machine
00000000000●

Security
00000000000

# From spec to execution: a series of refinement

We want to translate our IFSM $S_{spec}$ into our CPU FSM $S_{execution}$.

It is actually a multi-stage process, involving (non-exhaustively)

$$S_{spec} \supseteq S_{language} \supseteq S_{machine} \supseteq S_{execution}$$

- $S_{spec} \not\supseteq S_{language}$: software bug, blame the developer
- $S_{language} \not\supseteq S_{machine}$: compiler bug, blame the compiler
- $S_{machine} \not\supseteq S_{execution}$: hardware bug, blame the machine

# Outline

1 Introduction

2 A tale of two state machines

3 Defining security

Introduction
0000000

State machine
000000000000

Security
00●000000000

23 / 33

## Bug $\implies$ exploits?

Does having a bug in the refinement chain always imply a security issue (a.k.a., an exploit)?

Introduction
0000000

State machine
0000000000000

Security
0000000000000

# What is security?

Introduction
0000000

State machine
000000000000

Security
000●000000000

## What is security?

Security is properties of the IFSM that we want to hold in the presence of an adversary with a specific attack model.

Introduction
0000000

State machine
0000000000000

Security
0000●00000000

# Security of our secret-keeper

Introduction
0000000

State machine
000000000000

Security
0000000000000

## Security of our secret-keeper

Informally, we want to ensure that anyone who interact with our program needs to know (or guess) the right password in order to obtain the stored secret.

Put in a different way, the best way to attack our program to extract some secret is to guess the password.

Introduction
0000000

State machine
0000000000000

Security
00000●000000

## Security of our secret-keeper

Formally, we want the security property to hold at our IFSM:

$$Pr[s \in O_{IFSM}] \leq \frac{|I_{attempt}|}{2^{32}}$$

Introduction
0000000

State machine
000000000000

Security
00000●000000

## Security of our secret-keeper

Formally, we want the security property to hold at our IFSM:

$$Pr[s \in O_{IFSM}] \leq \frac{|I_{attempt}|}{2^{32}}$$

As well as at the final execution stage, after the refinement chain

$$Pr[s \in O_{execution}] \leq \frac{|I_{attempt}|}{2^{32}}$$

## Security of our secret-keeper

Formally, we want the security property to hold at our IFSM:

$$Pr[s \in O_{IFSM}] \leq \frac{|I_{attempt}|}{2^{32}}$$

As well as at the final execution stage, after the refinement chain

$$Pr[s \in O_{execution}] \leq \frac{|I_{attempt}|}{2^{32}}$$

Even in the presence of an attacker with the assumed power of performing single chosen bit-flip.

Introduction
0000000

State machine
0000000000000

Security
00000●000000

# The security property depends on the implementation

- Naive implementation: Simulate the *Memory* set as a flat linear array with sequential scanning

Introduction
0000000

State machine
000000000000

Security
00000●000000

# The security property depends on the implementation

- Naive implementation: Simulate the *Memory* set as a flat linear array with sequential scanning

- Clever implementation: Simulate the *Memory* set with two singly-linked lists.

# The security property depends on the implementation

- Naive implementation: Simulate the *Memory* set as a flat linear array with sequential scanning

- Clever implementation: Simulate the *Memory* set with two singly-linked lists.

**Conclusion**: the clever implementation is actually vulnerable.

## An attack on the clever implementation

1. Attacker sends $(p_0, s_0)$, $(p_1, s_1)$, $(p_2, s_2)$
2. Victim sends $(p_d, s_d)$
3. Attacker sends $(p_2, 0)$, $(p_1, 0)$, $(p_3, s_3)$, $(p_4, s_4)$
4. Attacker gets to corrupt a single bit: flip the least significant bit for memory cell content at b'0101 (i.e., cell 0x5)
5. Attacker sends $(s_4, 0)$
6. Attacker sends $(12, 0)$ and obtains $s_d$

Introduction
0000000

State machine
0000000000000

Security
00000000●0000

# The naive implementation is secure

Please refer to the paper for the details of the proof.

# Programming the weird machine

Introduction
○○○○○○○
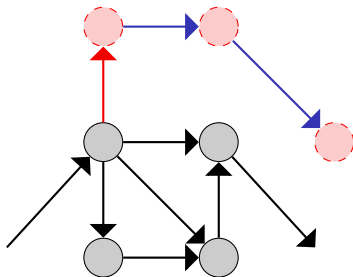
State machine
○○○○○○○○○○○○○○

Security
○○○○○○○○○○●○○○

# Programming the weird machine

Introduction
0000000

State machine
0000000000000

Security
00000000●0000

# Programming the weird machine

Introduction
0000000

State machine
0000000000000

Security
00000000000●000

# Programming the weird machine

Introduction
0000000

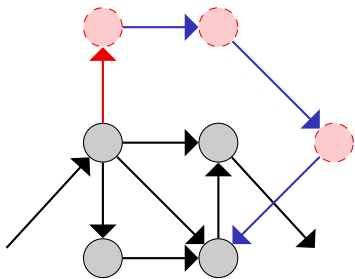State machine
000000000000

Security
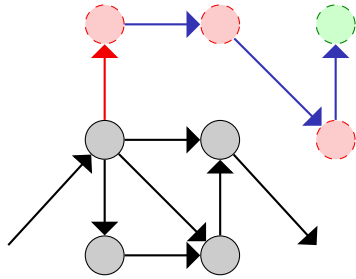00000000000●00

## An emergent instruction set

This weird machine creates an emergent instruction set that is constrained by:

- The IFSM
- The program that is refined from the IFSM
- The CPU FSM

Introduction
0000000

State machine
000000000000

Security
00000000000●0

# Outcomes of weird machine programming



Reverted back to the IFSM

Reached the target state

⟨ **End** ⟩