

# Sandboxing

---

**CS-576 Systems Security**

Instructor: Georgios Portokalidis

Spring 2018

# Sandboxing Means Isolation

## Why?

Software has bugs

Defenses slip

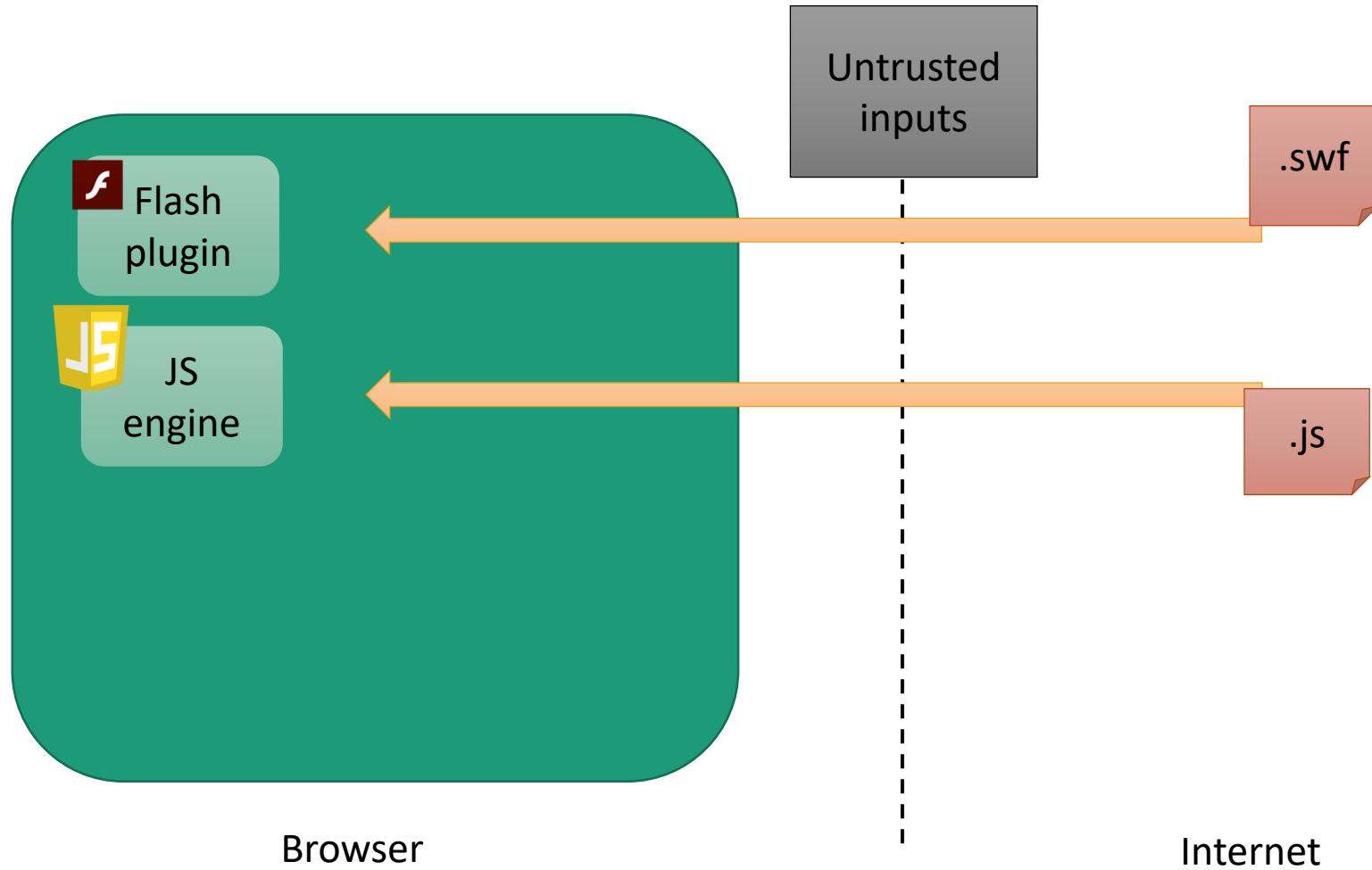
Untrusted code

Compartmentalization  
limits interference and  
damage!

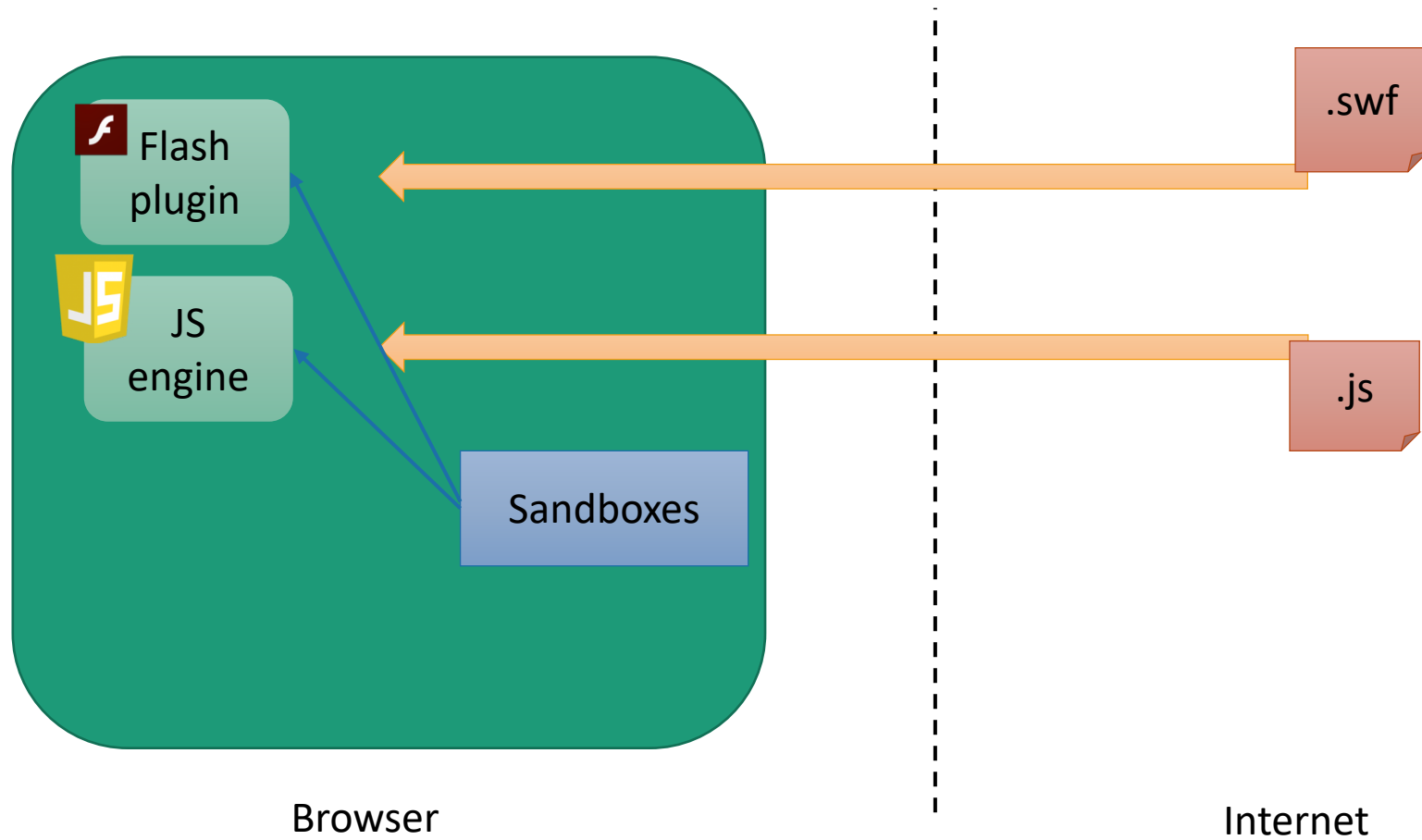


**“a sandbox is a security mechanism for separating running programs”**  
**-- wikipedia**

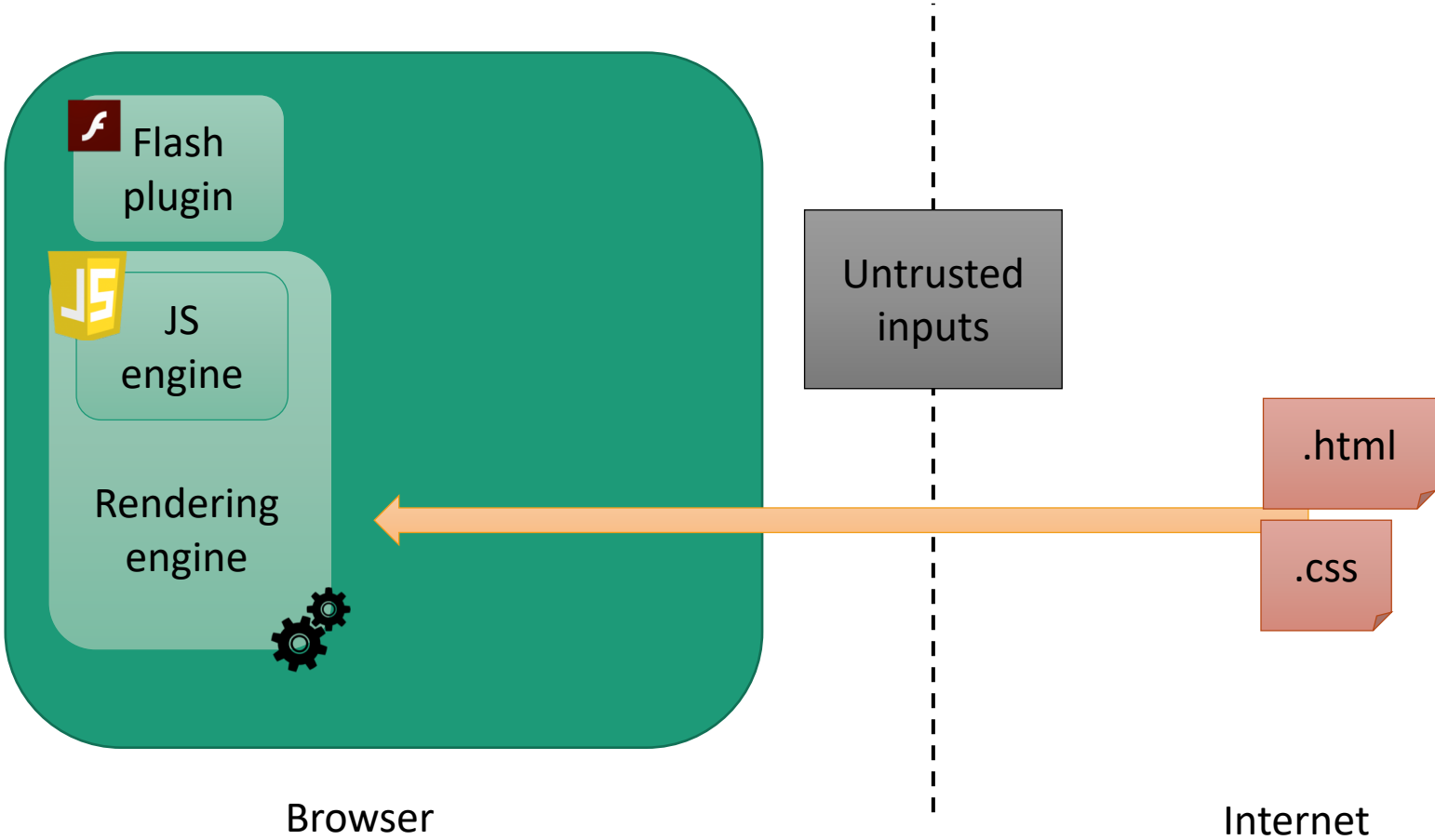
# Opportunities for Sandboxing: Browsers



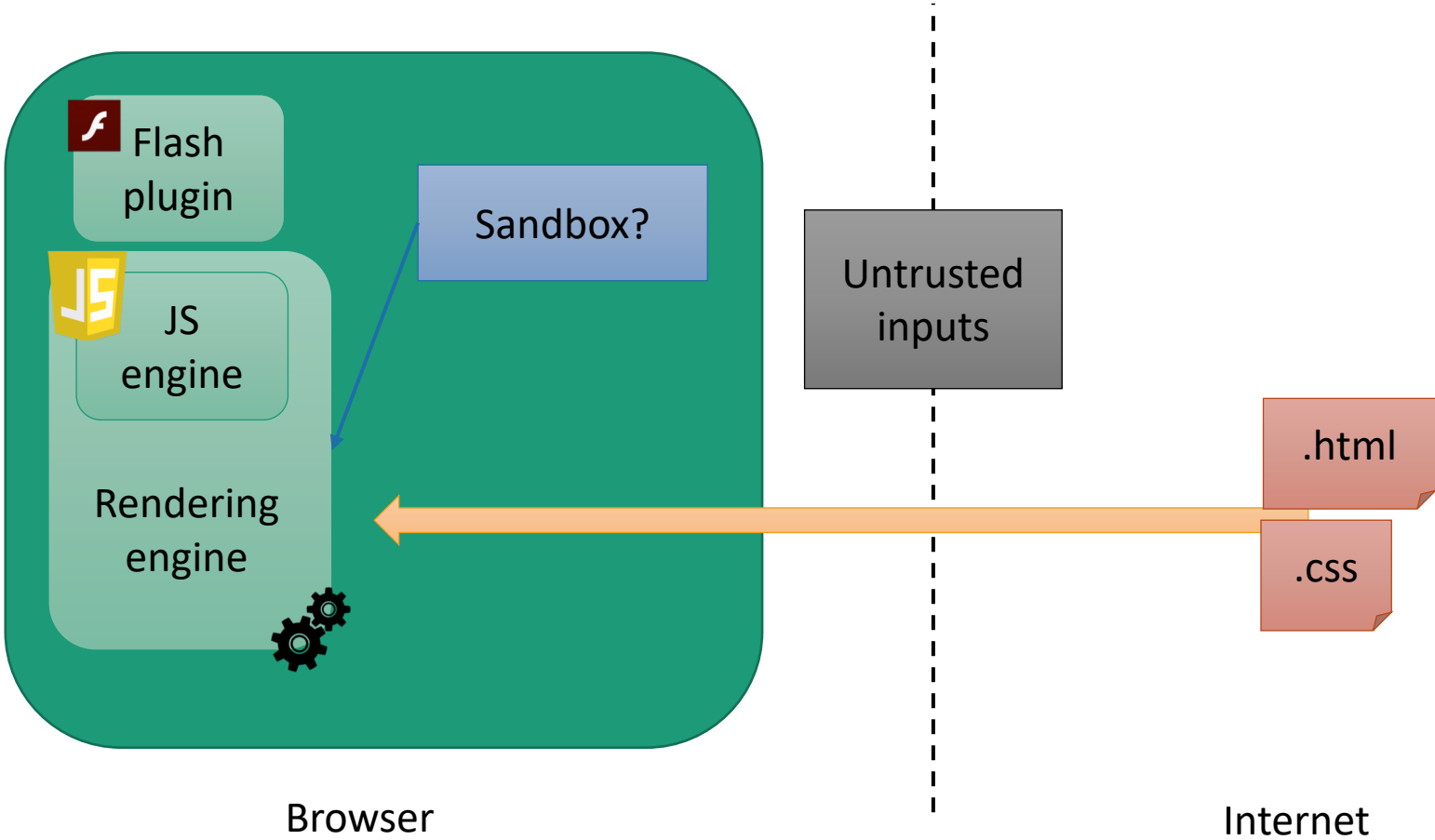
# Opportunities for Sandboxing: Browsers



# Untrusted Code in Browsers



# Untrusted Code in Browsers



# Sandboxing Methods

---

## VM-based

- Run entire OS in isolation

## OS-based

- Process-wide
- Available system calls and capabilities are restricted

## Language-based

- Language isolates components

## Inline reference monitor

- Integrated into untrusted code during compilation, code generation, or through emulation
- Security checks injected to enforce policy

# Sandboxing Methods

## VM-based

- Run entire OS in isolation

## OS-based

- Process-wide
- Available system calls and capabilities are restricted

## Language-based

- Language isolates components

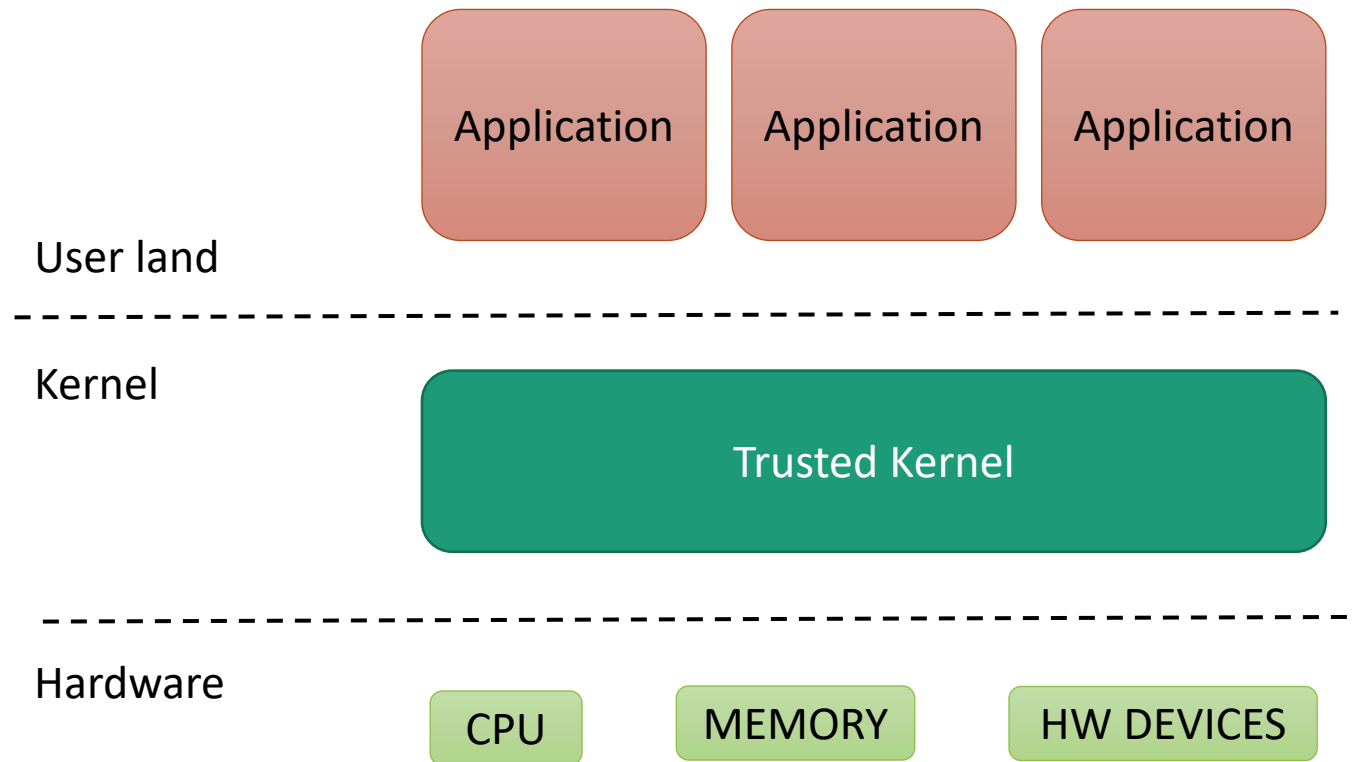
## Inline reference monitor

- Integrated into untrusted code during compilation, code generation, or through emulation
- Security checks injected to enforce policy

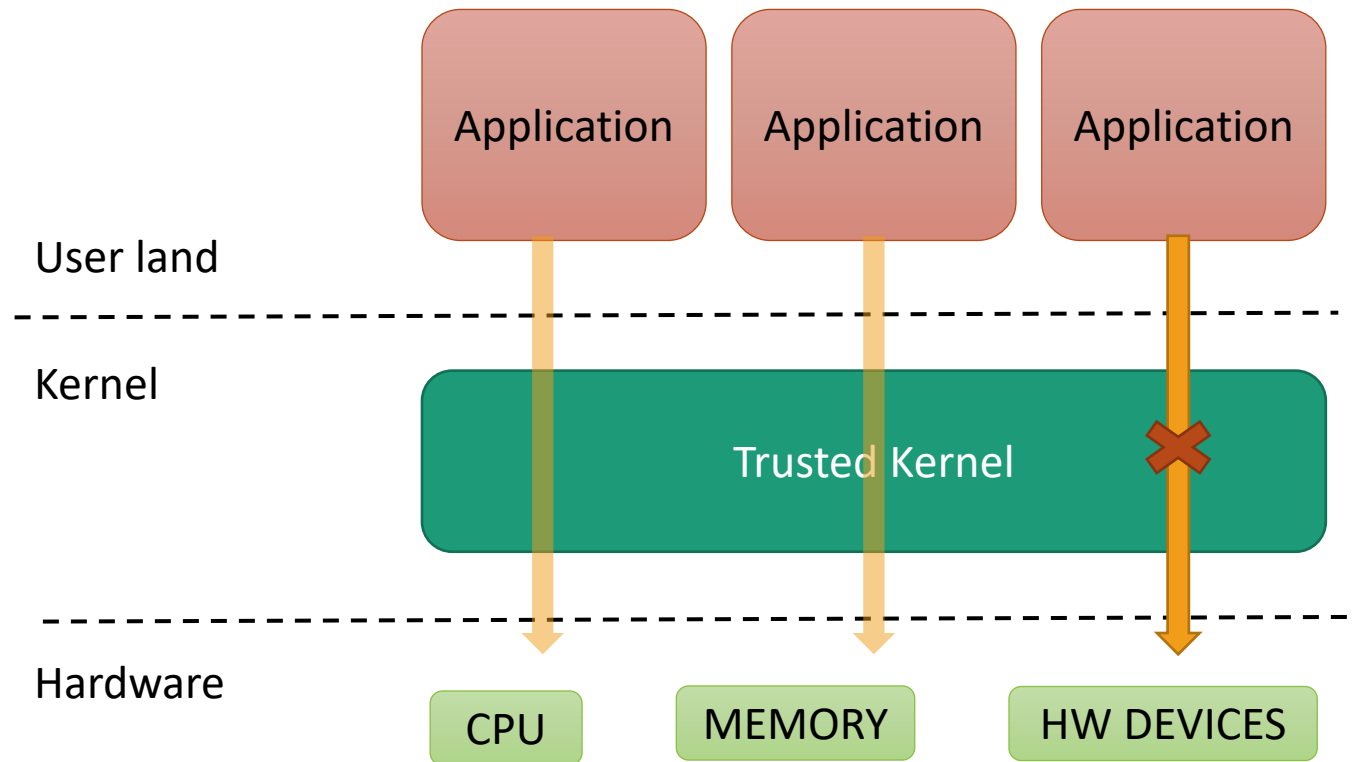


# Lets Refresh What We Know About OSes

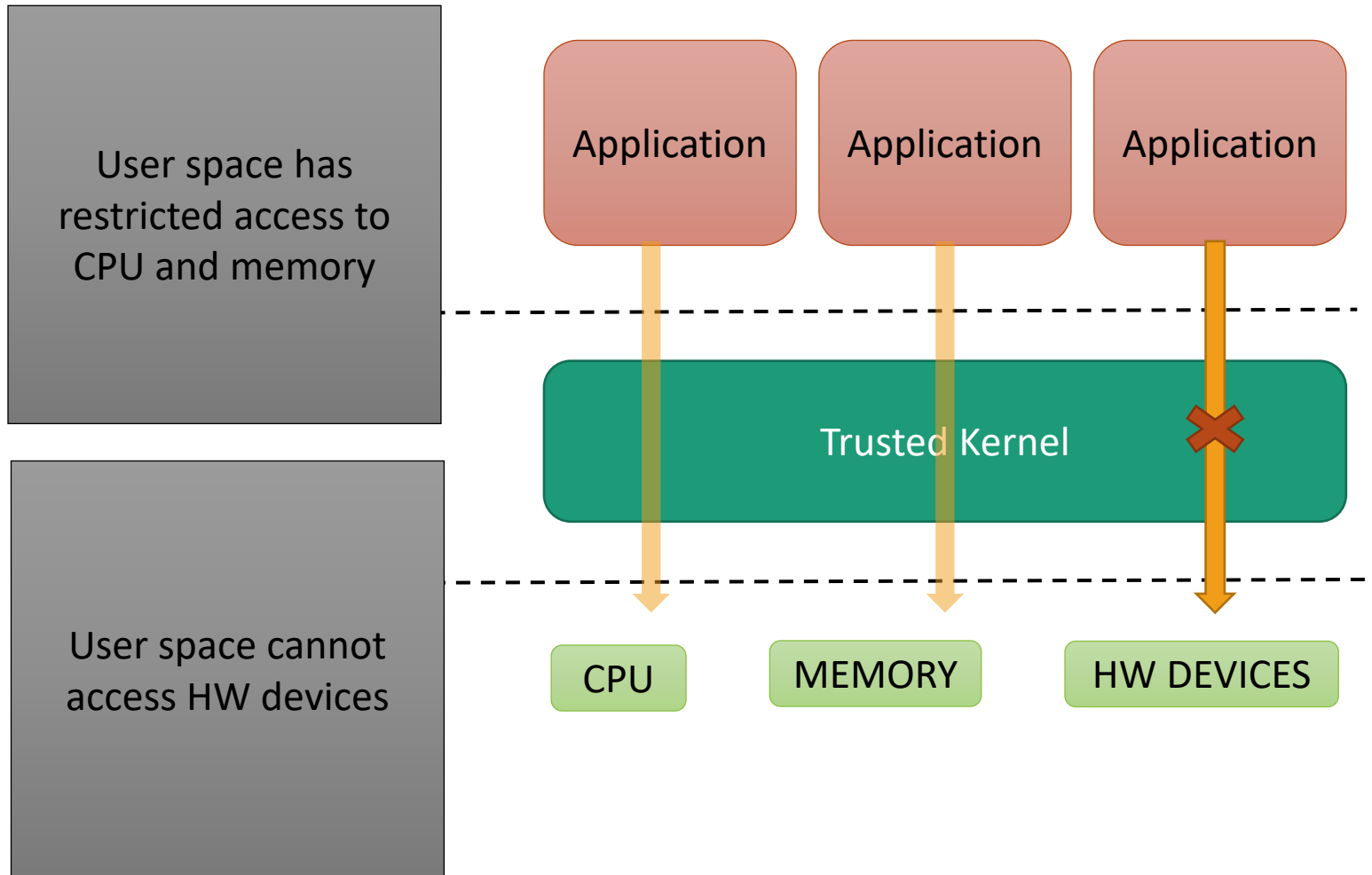
# OS Access Control of HW



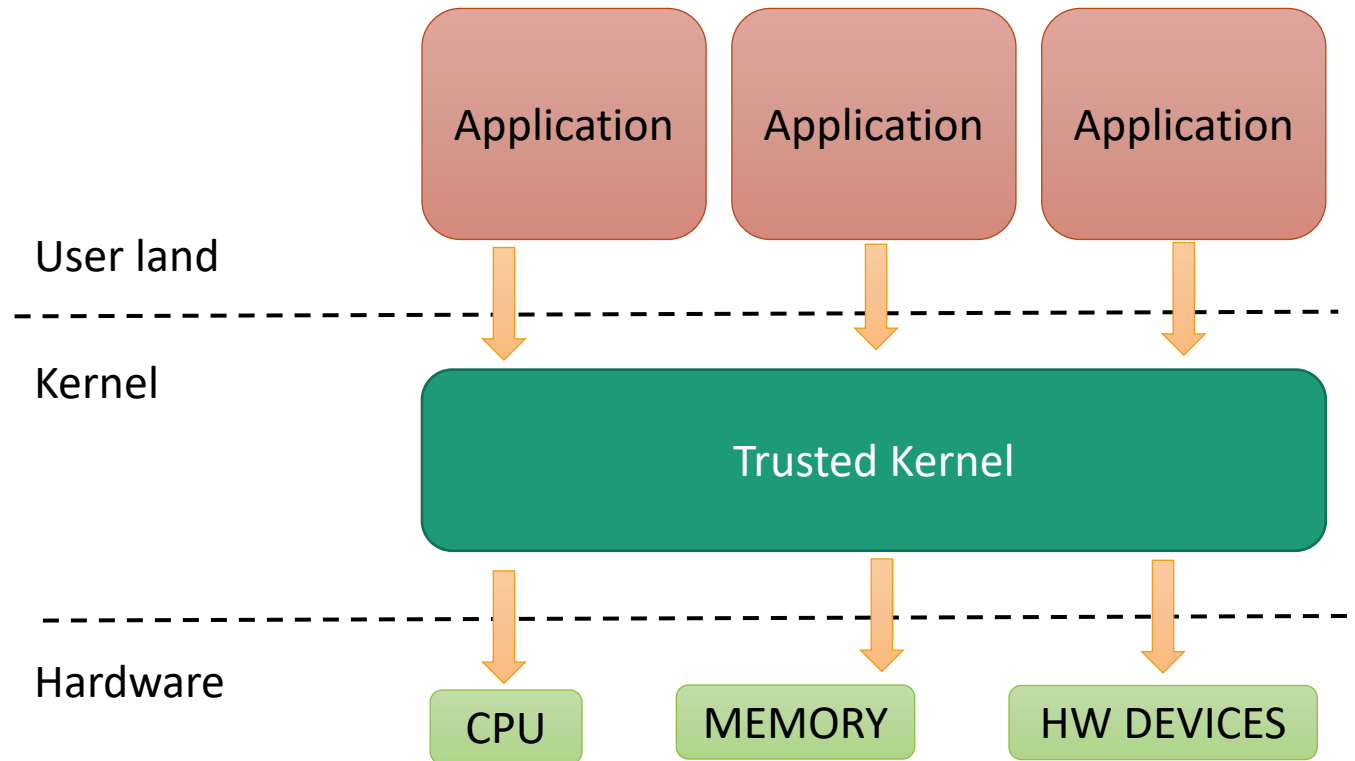
# OS Access Control of HW



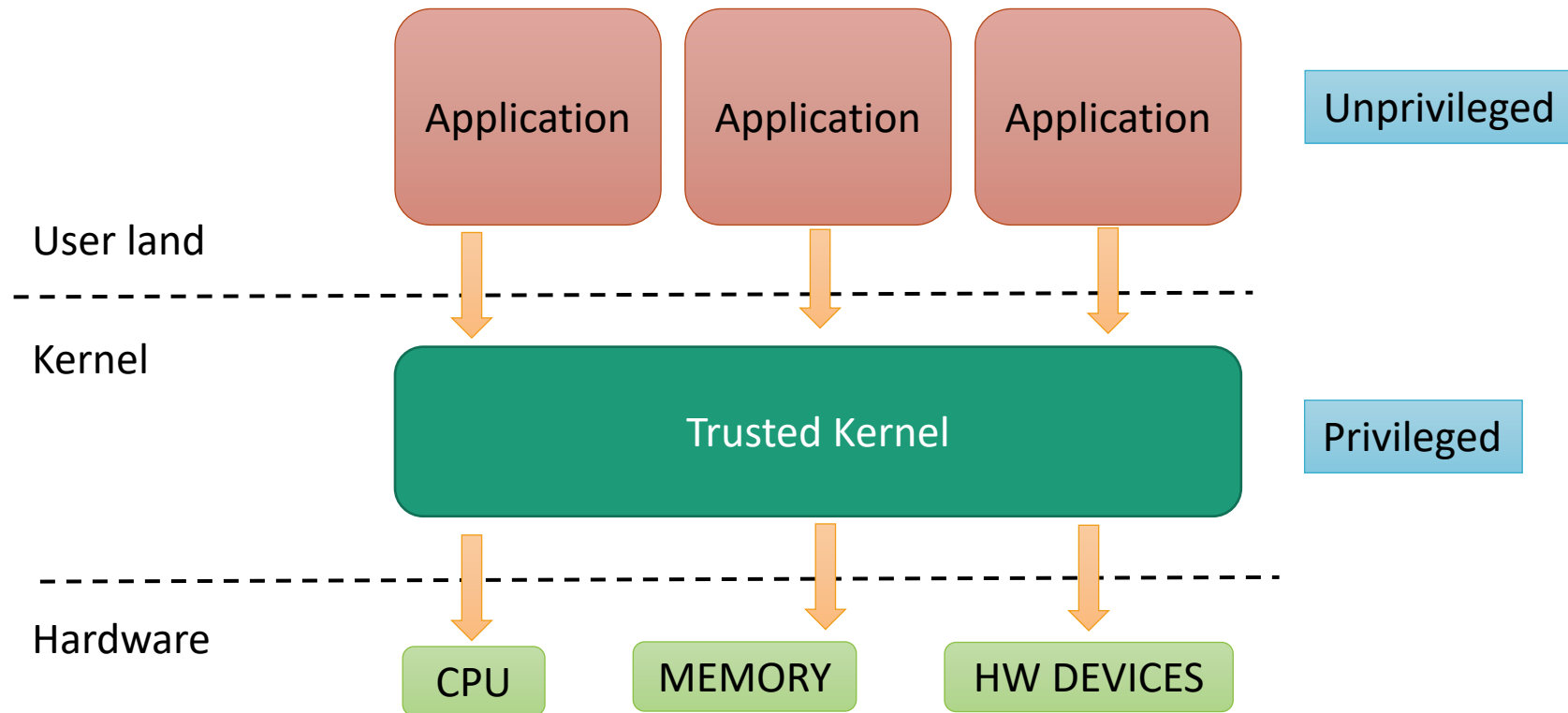
# OS Access Control of HW



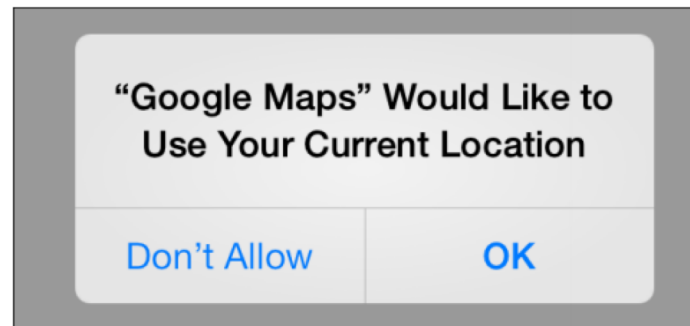
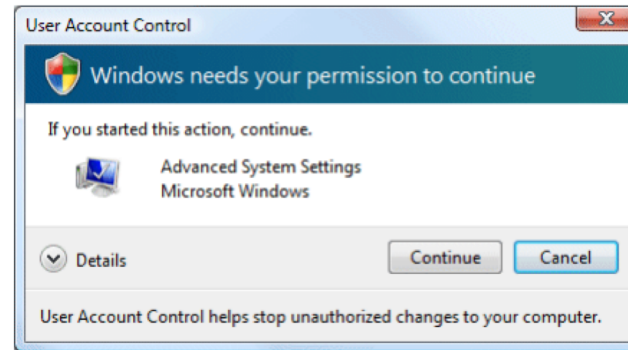
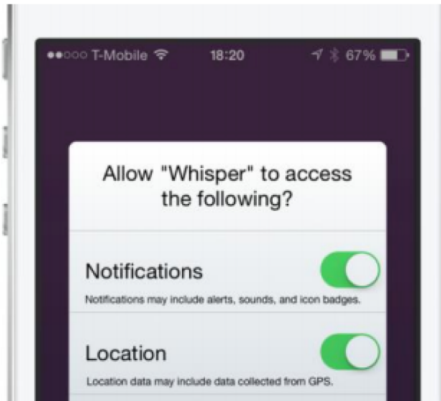
# OS Access Control of HW



# OS Access Control of HW

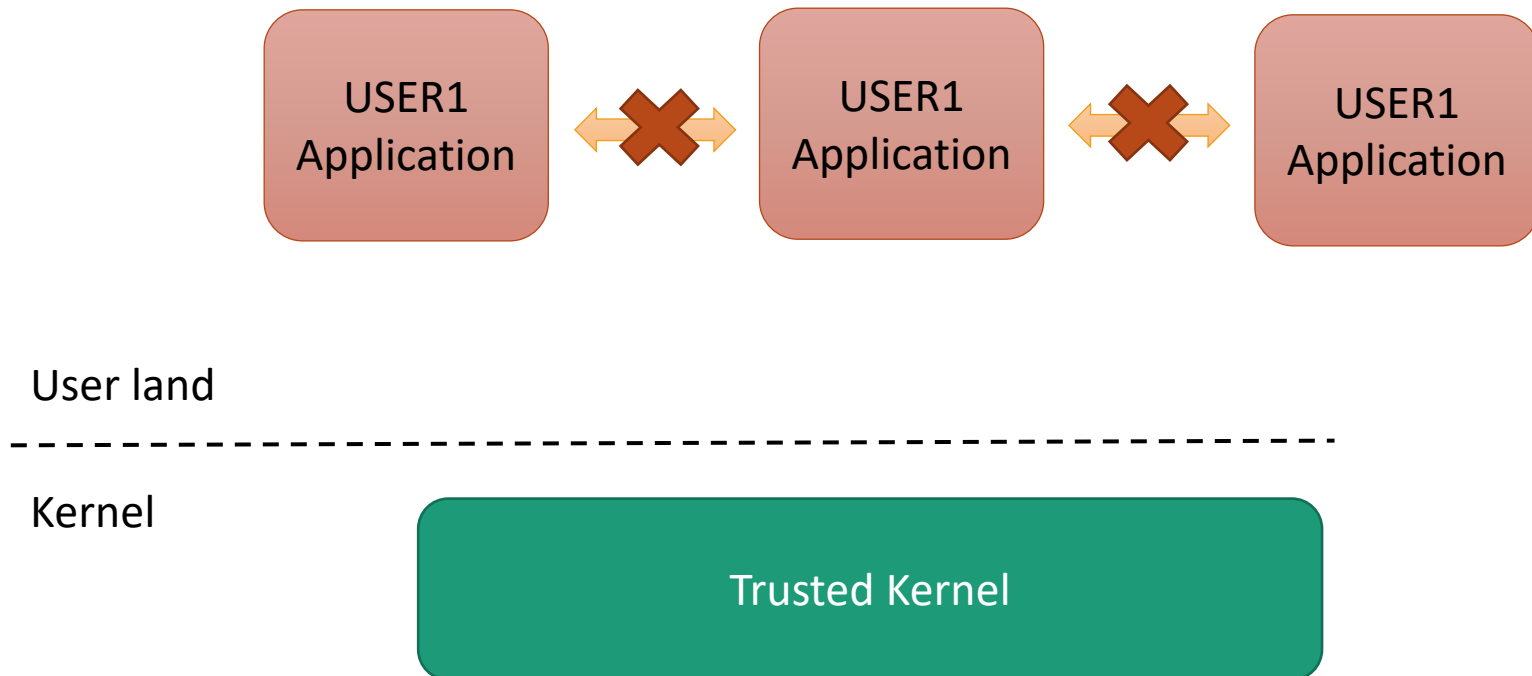


# Example of OS-Level Access Control to HW



# Process-level Isolation

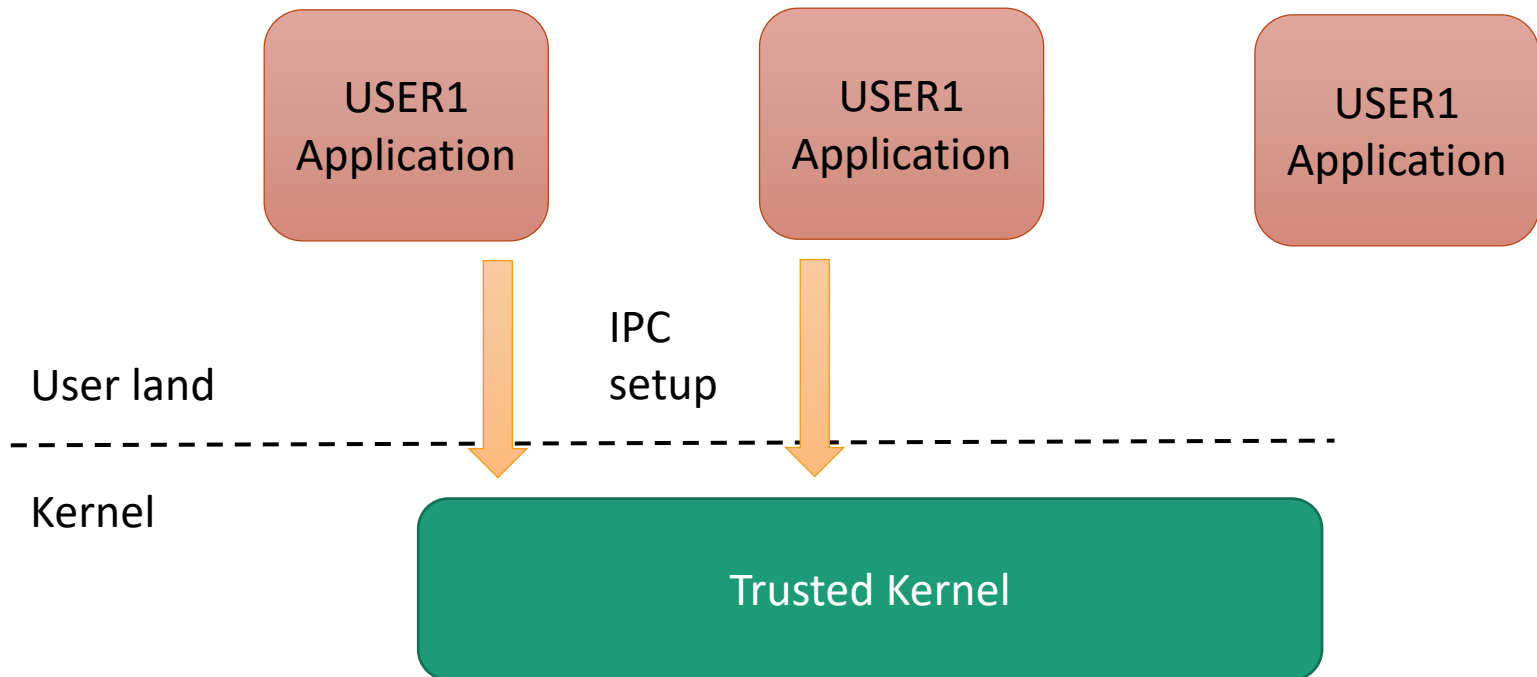
Processes cannot directly access each other's state





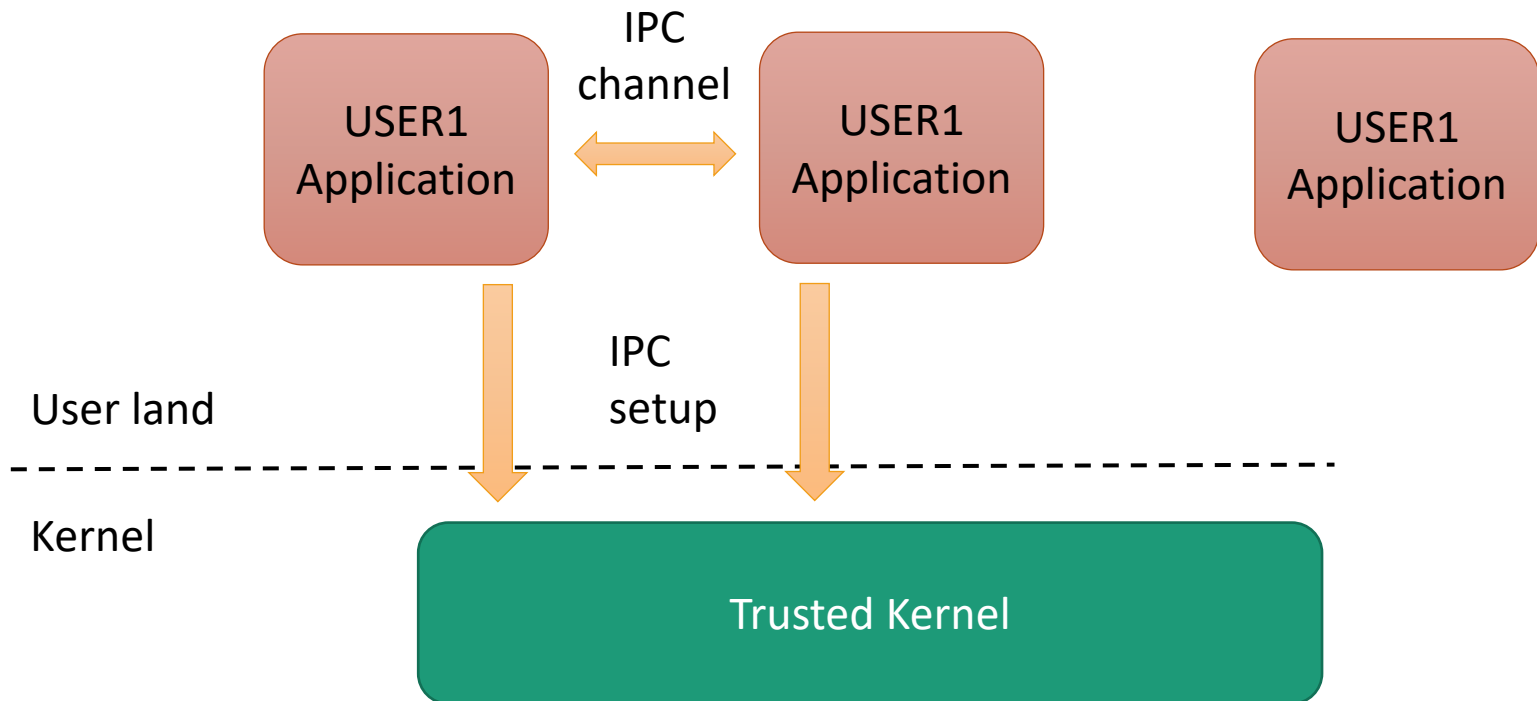
# Process-level Isolation

The kernel can setup inter-process communication



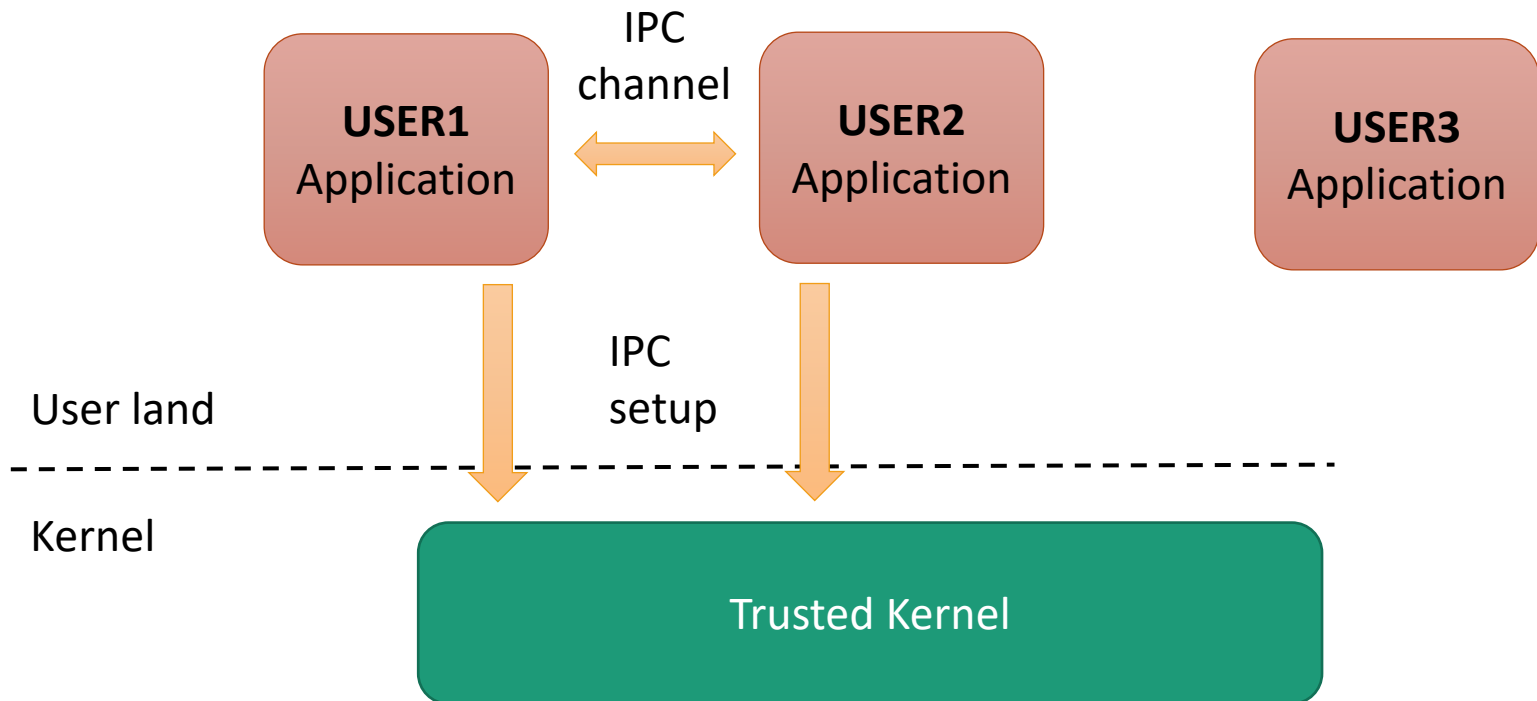
# Process-level Isolation

The kernel can setup inter-process communication



# Process-level Isolation

Same for processes owned by different users



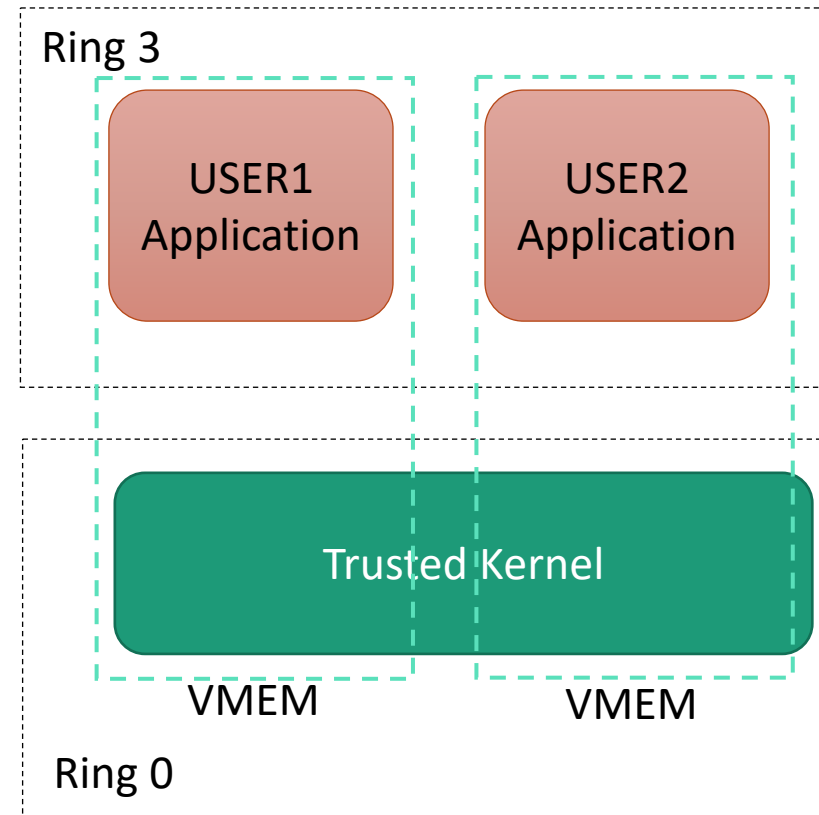
# Hardware-based Enforcement

The memory-management unit (MMU) provides virtual memory

Execution rings separate user and kernel space

- Indicated by bits in CPU status register

Processes are isolated into different virtual memory address spaces



# Back to Sandboxing

---

# Sandboxing Methods

## VM-based

- Run entire OS in isolation

## OS-based

- Process-wide
- Available system calls and capabilities are restricted

## Language-based

- Language isolates components

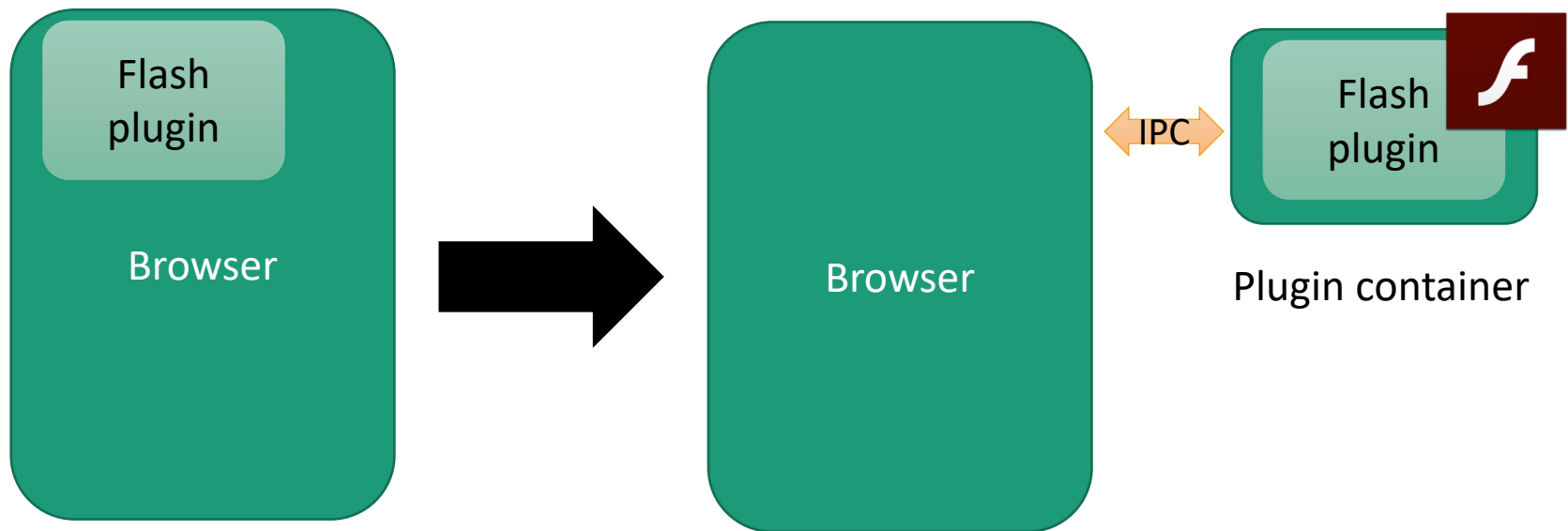
## Inline reference monitor

- Integrated into untrusted code during compilation, code generation, or through emulation
- Security checks injected to enforce policy

# Building on Process Isolation

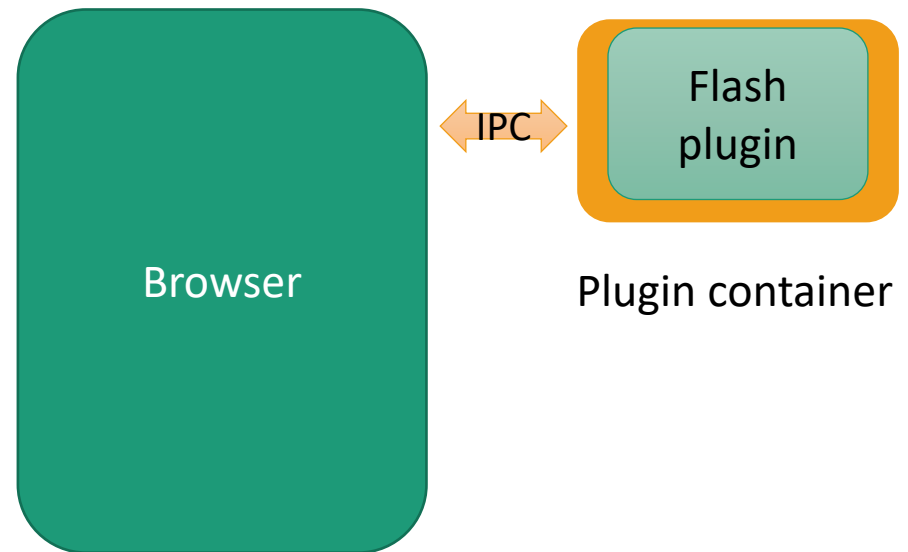
Run code in its own process space to isolate it from browser process

Congratulations you have just executed untrusted code from the Internet!



# Building on Process Isolation

Container must have limited privileges





# Chromium Sandboxing in Linux

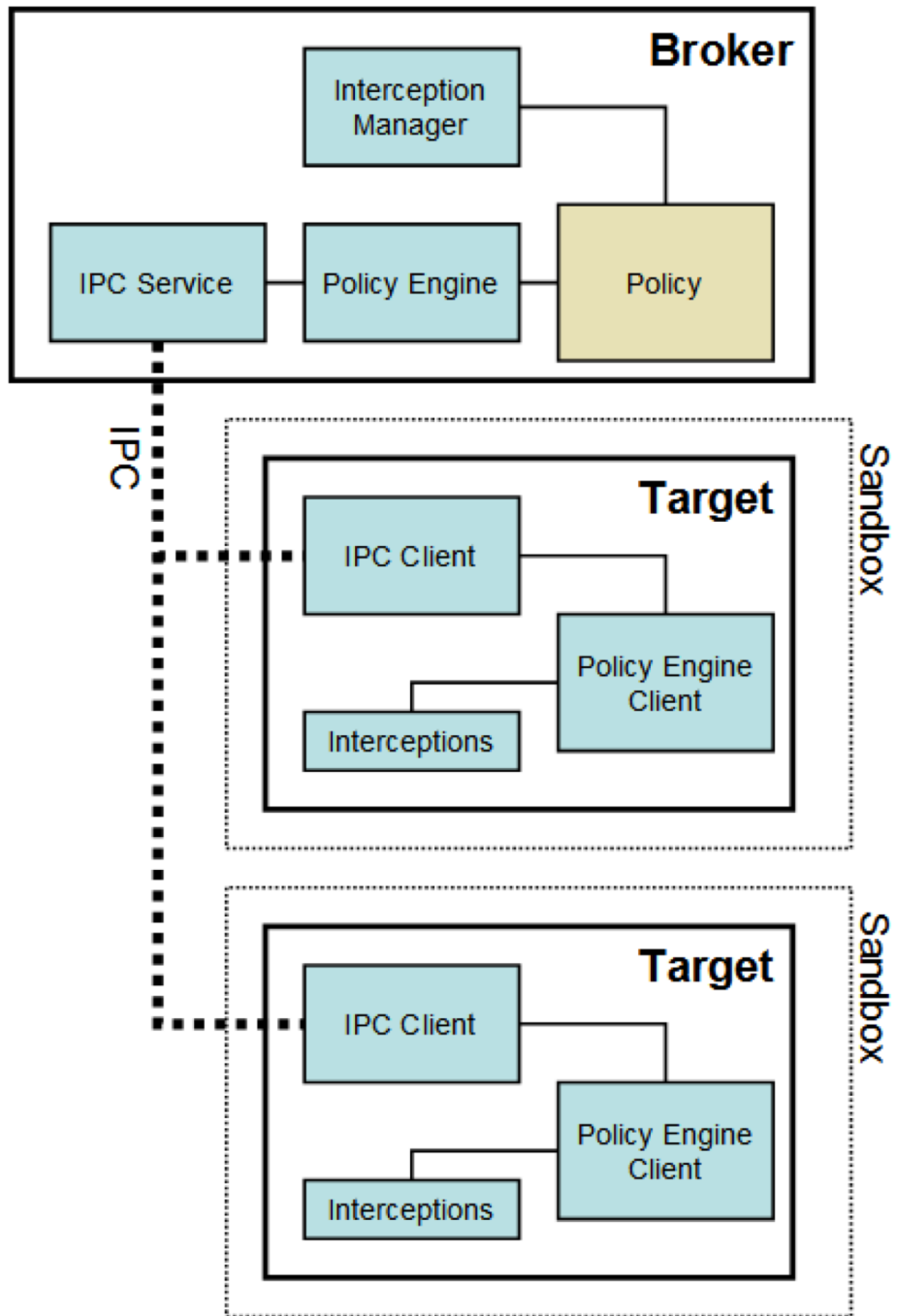
Chromium runs plugins and the rendering engine for each tab in a separate process

Rendering processes are sandboxed

Sandboxed processes are managed by a broker process over IPC



[https://chromium.googlesource.com/chromium/src/+master/docs/linux\\_sandboxing.md](https://chromium.googlesource.com/chromium/src/+master/docs/linux_sandboxing.md)



# Process Sandbox: SUID

A helper binary with the setuid bit set is used

The SUID bit causes the execution of the process as root

- Enables access to privileged kernel APIs

chroot() is used to change the process' root directory

- Take away file system access from the process

Process is placed in new PID namespace

- Process cannot terminate or signal processes outside the namespace

Process is placed in new network namespace

- Restrict network access of process

Finally drop super-user privileges

# Process Sandbox: User Namespaces

---

User namespaces are an unprivileged API

Used as an alternative to SUID sandbox

A process is placed a new namespace

Isolates:

- Filesystem
- Network
- PID
- IPC

# User Namespaces

A newly launched process can be put in a new namespace

- Through the clone() system call

## Available namespaces

| Namespace | Constant        | Isolates                             |
|-----------|-----------------|--------------------------------------|
| Cgroup    | CLONE_NEWCGROUP | Cgroup root directory                |
| IPC       | CLONE_NEWIPC    | System V IPC, POSIX message queues   |
| Network   | CLONE_NEWNET    | Network devices, stacks, ports, etc. |
| Mount     | CLONE_NEWNS     | Mount points                         |
| PID       | CLONE_NEWPID    | Process IDs                          |
| User      | CLONE_NEWUSER   | User and group IDs                   |
| UTS       | CLONE_NEWUTS    | Hostname and NIS domain name         |

Reading material: <https://lwn.net/Articles/531114/>

# Process Sandbox: SECCOMP BPF

Filters the kernel APIs available to a process

Used together with previous sandboxes

Aims to protect the kernel from a malicious process

Available system calls are defined using **Berkeley packet filters**

- Filters are compiled to a program that enforces policy

```

static int install_syscall_filter(void)
{
    struct sock_filter filter[] = {
        /* Validate architecture. */
        VALIDATE_ARCHITECTURE,
        /* Grab the system call number. */
        EXAMINE_SYSCALL,
        /* List allowed syscalls. */
        ALLOW_SYSCALL(rt_sigreturn),
#ifdef __NR_sigreturn
        ALLOW_SYSCALL(sigreturn),
#endif
        ALLOW_SYSCALL(exit_group),
        ALLOW_SYSCALL(exit),
        ALLOW_SYSCALL(read),
        ALLOW_SYSCALL(write),
        KILL_PROCESS,
    };
    struct sock_fprog prog = {
        .len = (unsigned short)(sizeof(filter)/sizeof(filter[0])),
        .filter = filter,
    };
}

```

```
if (prctl(PR_SET_NO_NEW_PRIVS, 1, 0, 0, 0)) {
    perror("prctl(NO_NEW_PRIVS)");
    goto failed;
}
if (prctl(PR_SET_SECCOMP, SECCOMP_MODE_FILTER, &prog)) {
    perror("prctl(SECCOMP)");
    goto failed;
}
return 0;
```

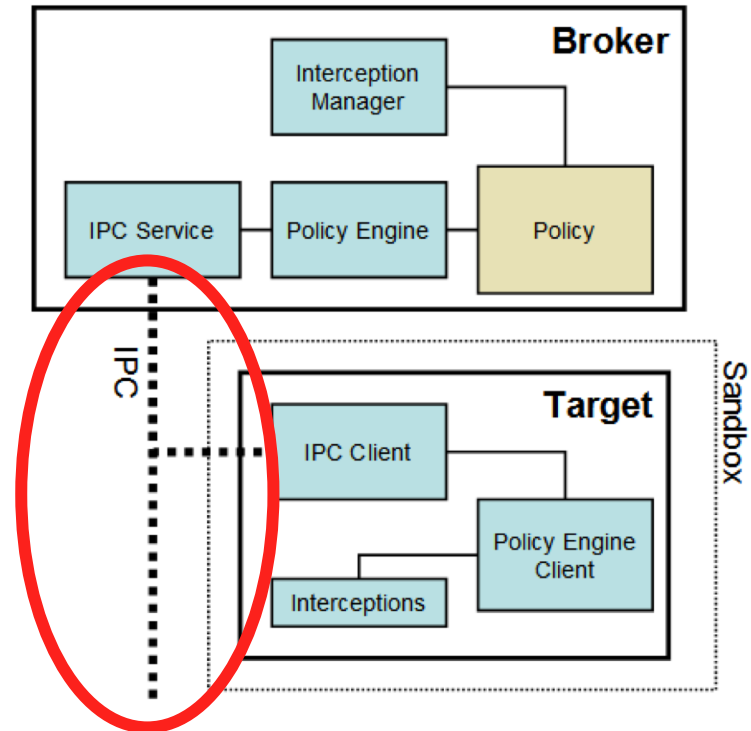
failed:

```
if (errno == EINVAL)
    fprintf(stderr, "SECCOMP_FILTER is not available. :(\n");
return 1;
}
```



# Limitations of OS and VM-based Sandboxing

Context switches between broker and sandboxed processes can be expensive



# Sandboxing Methods

## VM-based

- Run entire OS in isolation

## OS-based

- Process-wide
- Available system calls and capabilities are restricted

## Language-based

- Language isolates components

## Inline reference monitor

- Integrated into untrusted code during compilation, code generation, or through emulation
- Security checks injected to enforce policy

# Example: JS/Java

The language and the runtime environment/VM is enforcing security

- Memory safe languages
- Memory corruption or leakage is not possible (at least in theory)

Access control done at the API level, for example:

- Which files can be loaded
- Which frames are accessible through the DOM
- Where can code be loaded from
- **The VM acts as a reference monitor**

# Sandboxing Methods

## VM-based

- Run entire OS in isolation

## OS-based

- Process-wide
- Available system calls and capabilities are restricted

## Language-based

- Language isolates components

## Inline reference monitor

- Integrated into untrusted code during compilation, code generation, or through emulation
- Security checks injected to enforce policy

# Sandboxing Unsafe Languages

---

Pointers can be used to potential read/write arbitrary memory

Memory accesses need to be isolated first

- Can rarely rely on HW to contain memory operations
- Software checks are introduced in application code

# Software-fault Isolation

Run multiple programs in the same address space that run in isolation

Each program runs in a different logical fault domain

Programs can access memory within their domain

- Ensures memory secrecy and integrity

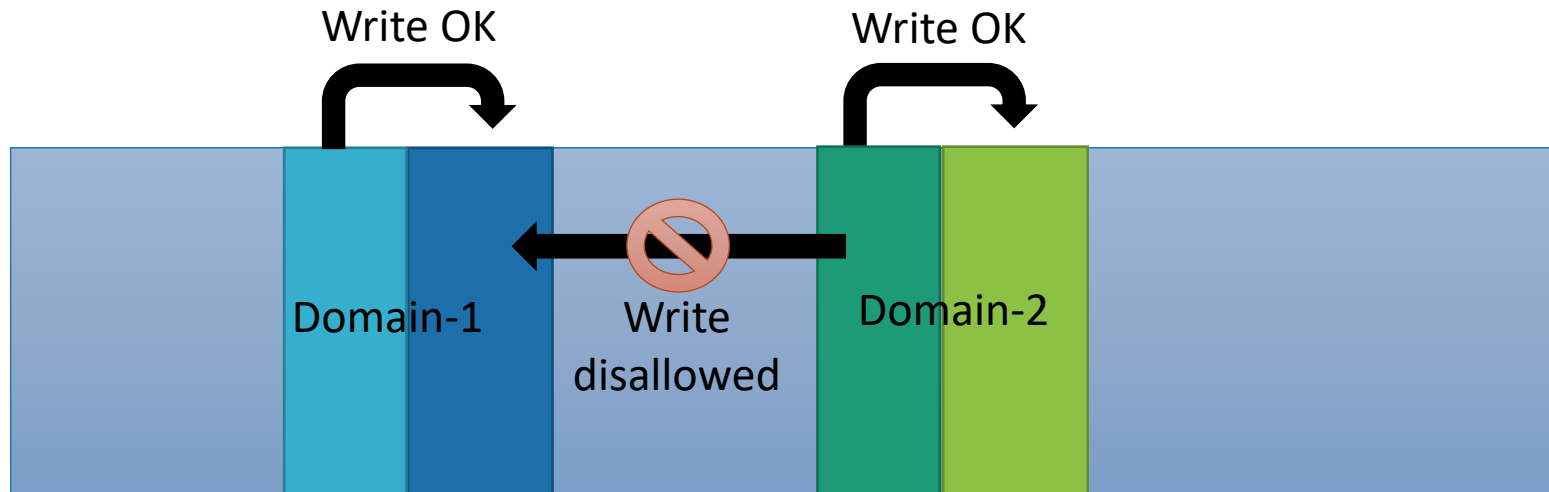
Code within a domain cannot call/jump to code in other domains

- Unless through secure interfaces

# Software-fault Isolation

Programs can only access memory within their domain

- Ensures memory secrecy and integrity



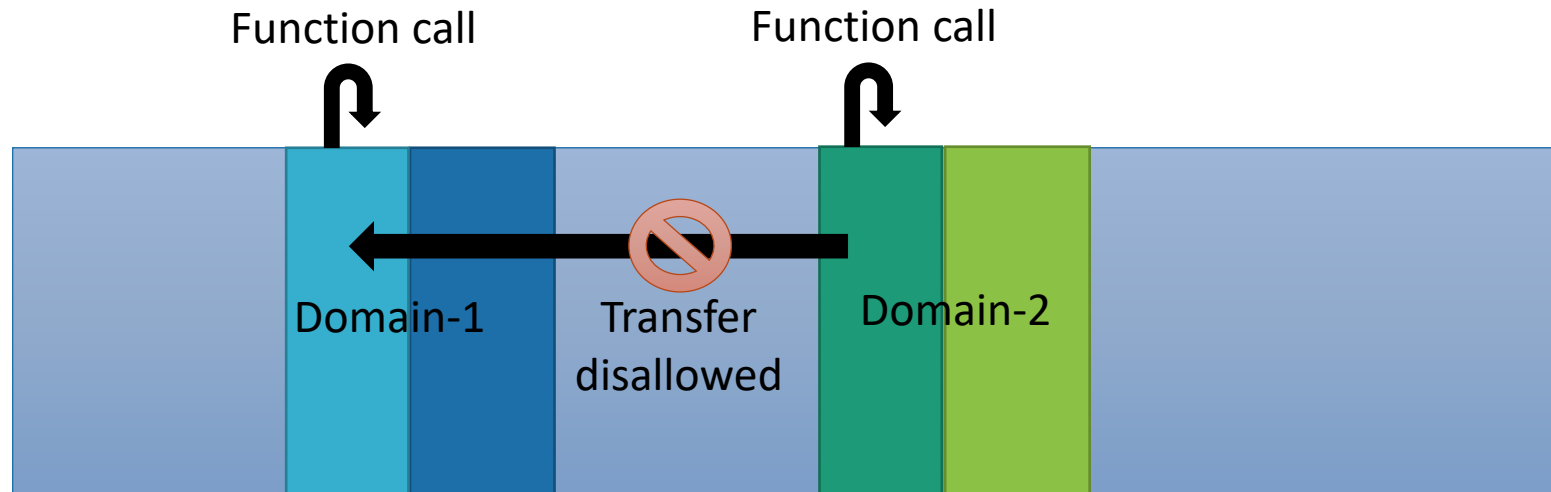
# Software-fault Isolation

Programs can only access memory within their domain

- Ensures memory secrecy and integrity

**Code within a domain cannot call/jump to code in other domains**

- Unless through secure interfaces





# Software-fault Isolation

---

Programs can only access memory within their domain

- Ensures memory secrecy and integrity

Code within a domain cannot call/jump to code in other domains

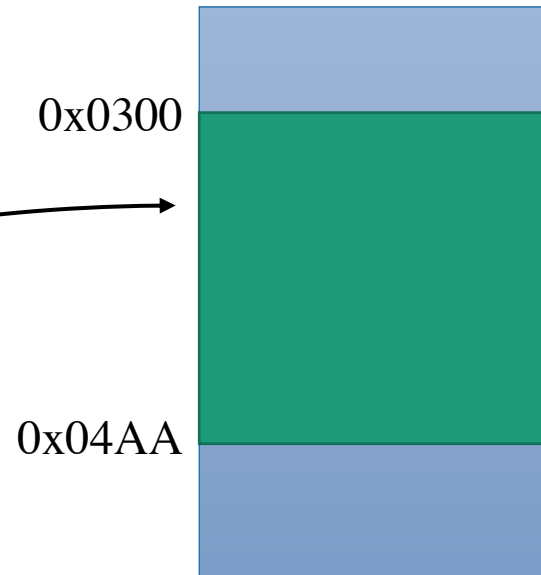
- Unless through secure interfaces

**Modify programs during compilation or by rewriting to enforce these properties**

# Constraining Memory Accesses

Through boundary checking

```
cmp 0x0300  
if less Error  
cmp 0x04AA  
if greater Error  
write x
```



# Constraining Memory Accesses

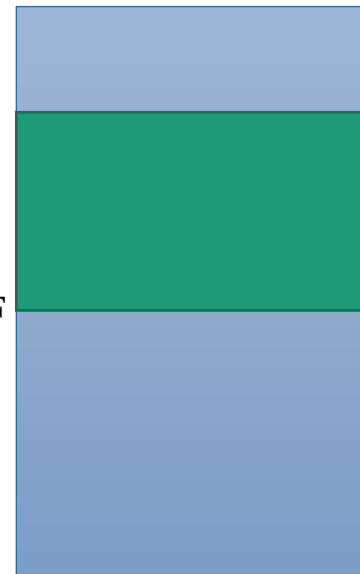
We can improve the boundary checks

- By allocating domains in aligned memory ranges
- Using bit masking to help with checking

```
tmp := x & FF00  
cmp tmp 0300  
if not equal Error  
write x
```

0x0300

0x03FF



# Constraining Memory Accesses

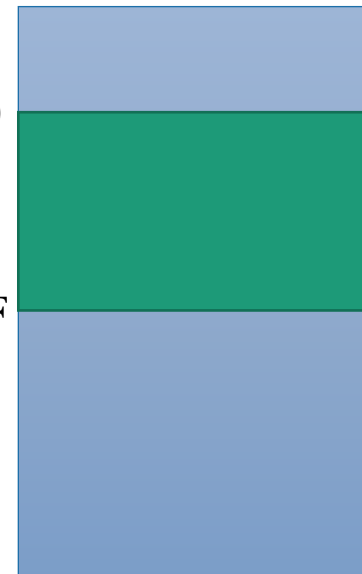
## Further improvements

- Do not detect error
- Constrain memory access to domain

```
tmp := x & 00FF  
tmp := tmp | 0300  
write tmp
```

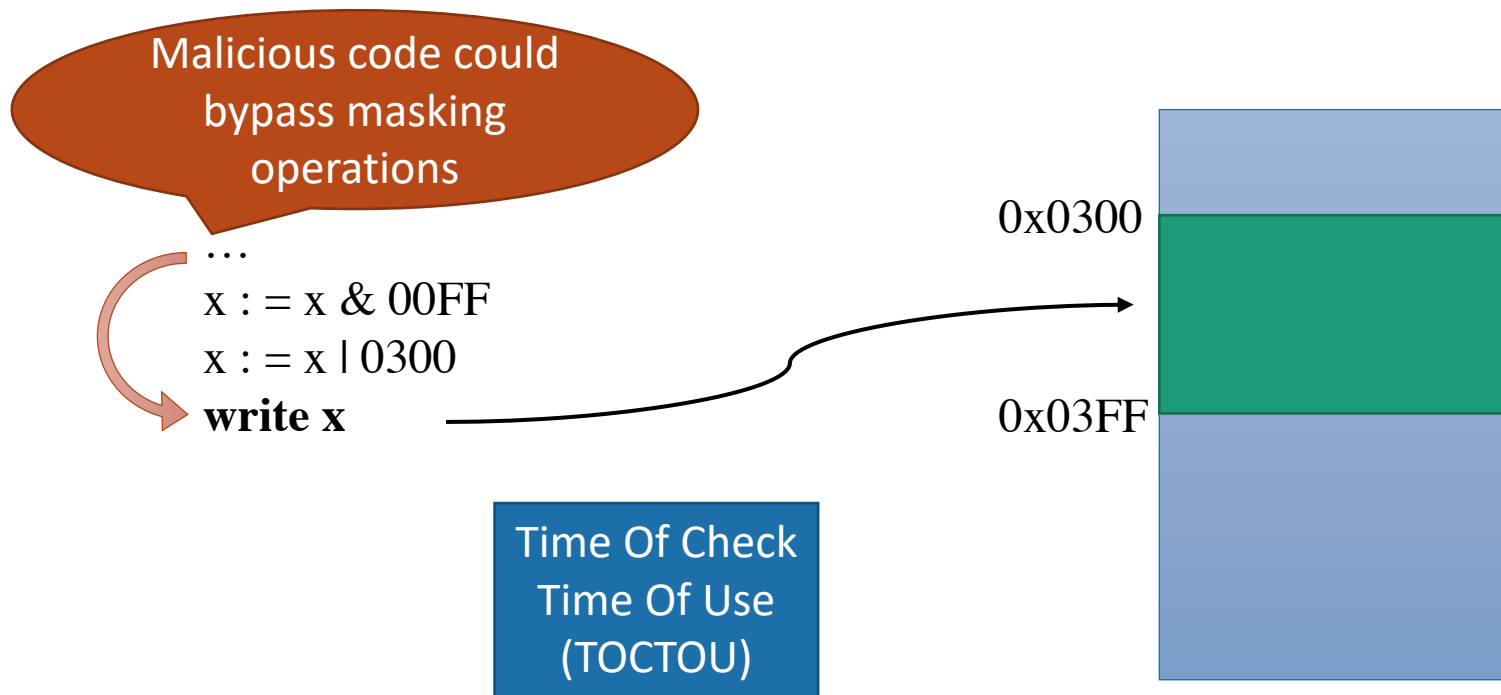
0x0300

0x03FF



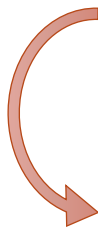
# Constraining Memory Accesses

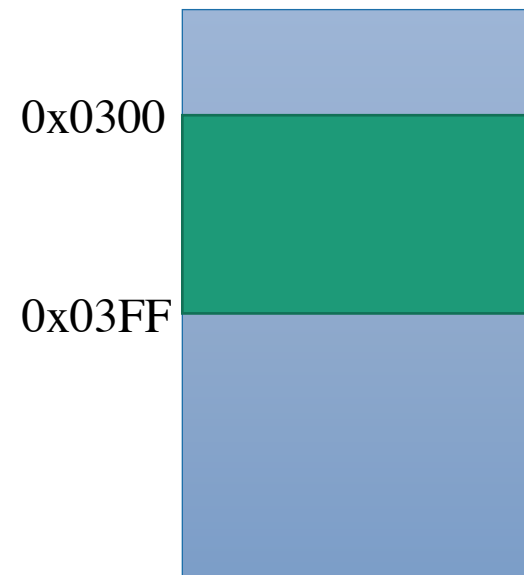
Eliminating temporary registers is not always a good idea



# Constraining Memory Accesses

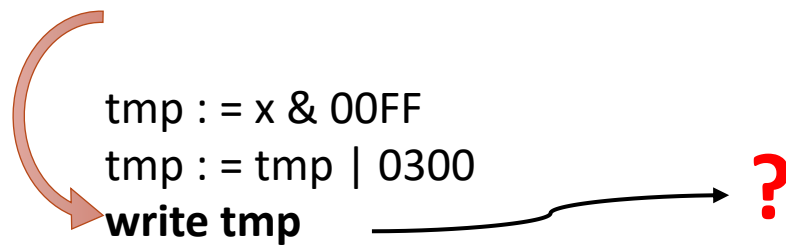
Can malicious code bypass checks with temporary registers?

 `tmp := x & 00FF`  
`tmp := tmp | 0300`  
**write tmp**

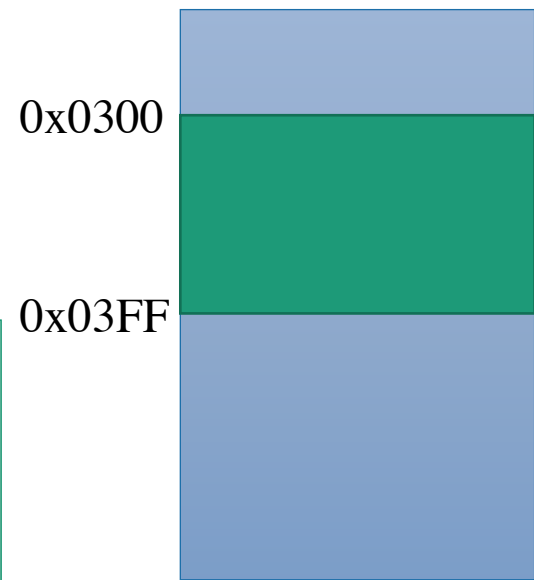


# Constraining Memory Accesses

Can malicious code bypass checks with temporary registers?

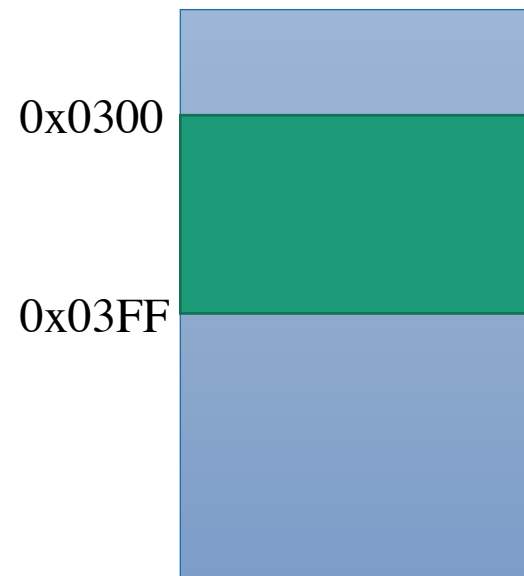
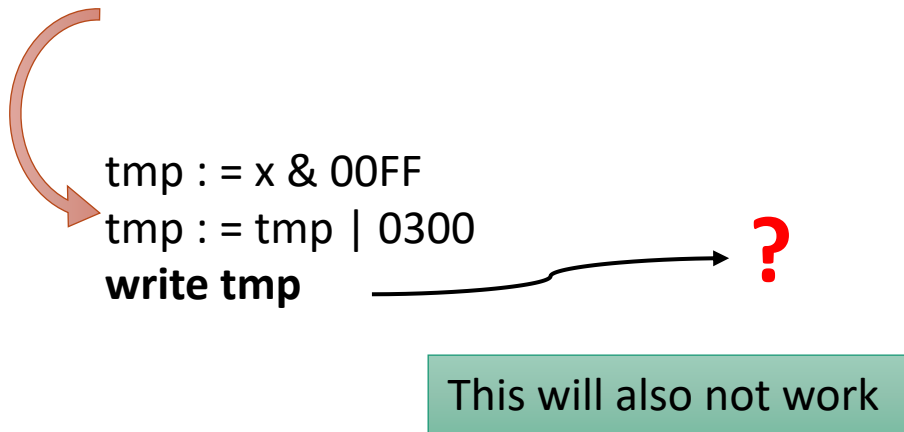


- tmp has not been initialized and will probably cause the program to crash.
- Can be forced to crash by setting tmp to bad address (e.g., 0xffffffff) after write



# Constraining Memory Accesses

Can malicious code bypass checks with temporary registers?





# Constraining Control Flow

Sandboxes are mainly to used to constrain untrusted code so obviously this is a general problem

...  
`jmp ptr` 

# Constraining Control Flow

Similar tricks can be applied

...  
**jmp ptr**



...  
tptr := ptr & 00FF  
tptr := tptr | 0300  
**jmp tptr**

...  
**call ptr**



...  
tptr := ptr & 00FF  
tptr := tptr | 0300  
**call tptr**

...  
**ret**



**?**

# Constraining Control Flow

Naive approach

**ret**



```
pop tptr  
tptr := tptr & 00FF  
tptr := tptr | 0300  
jmp ptr
```

# CISC Trouble

Constraining within the domain is not enough

- Instructions may be hidden within instructions in CISC programs



# Pseudo Fixed-size Instructions

Align every “pseudo” instruction on a 32-byte boundary

- 0x1F bits are always zero

Force pointer so it can only point to a pseudo instruction

```
pop tptr  
tptr := tptr & 00E0  
tptr := tptr | 0300  
jmp ptr
```

# Benefits of SFI

---

No context switches

Faster if run-time checks are faster than context switching

# Google Native Client (NaCL)

A sandboxing technology for running a subset of Intel x86, ARM, or MIPS **native** code in a sandbox

<https://developer.chrome.com/native-client>

NaCL programs are compiled with modified compiler

Supports subset of language

Produces sandboxed programs

# Escaping Sandboxes

Exploitation of a sandboxed component grants limited control

But sandboxes may have bugs

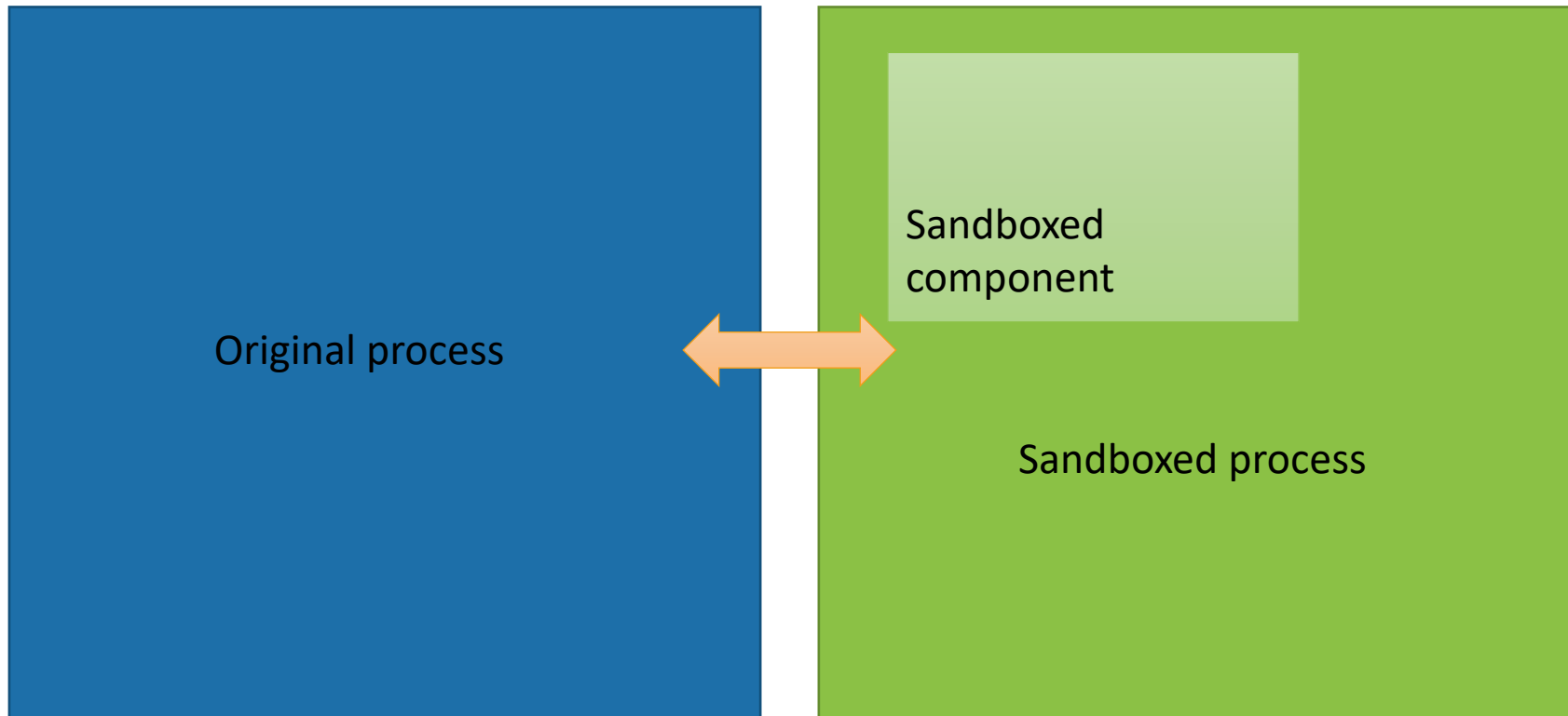
Multiple exploits in different components are usually required

In 2012's pwnium competition 14 bugs were needed to take down chrome

- <http://blog.chromium.org/2012/05/tale-of-two-pwnies-part-1.html>



# Multiple Layers of Sandboxes



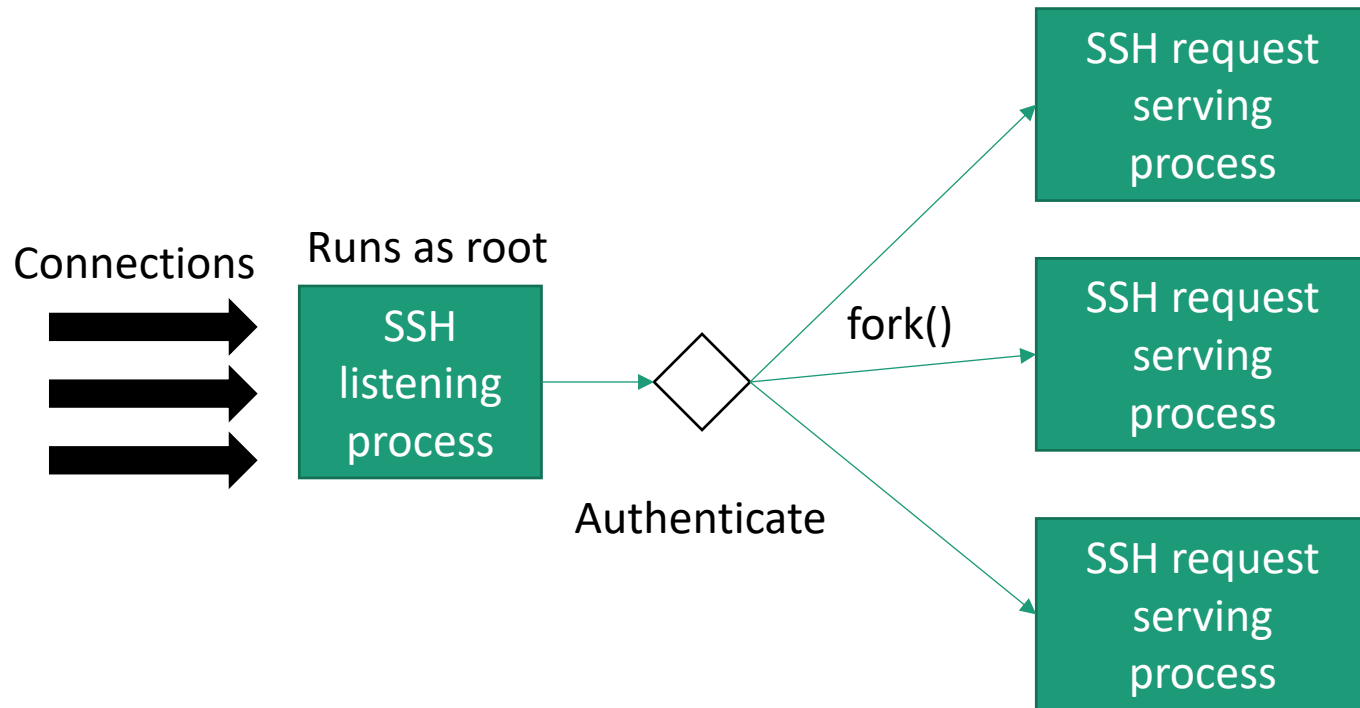
# Other Use Cases for Isolation

---

Process-level Isolation from the OS is frequently used to realize the principle of least privilege in servers

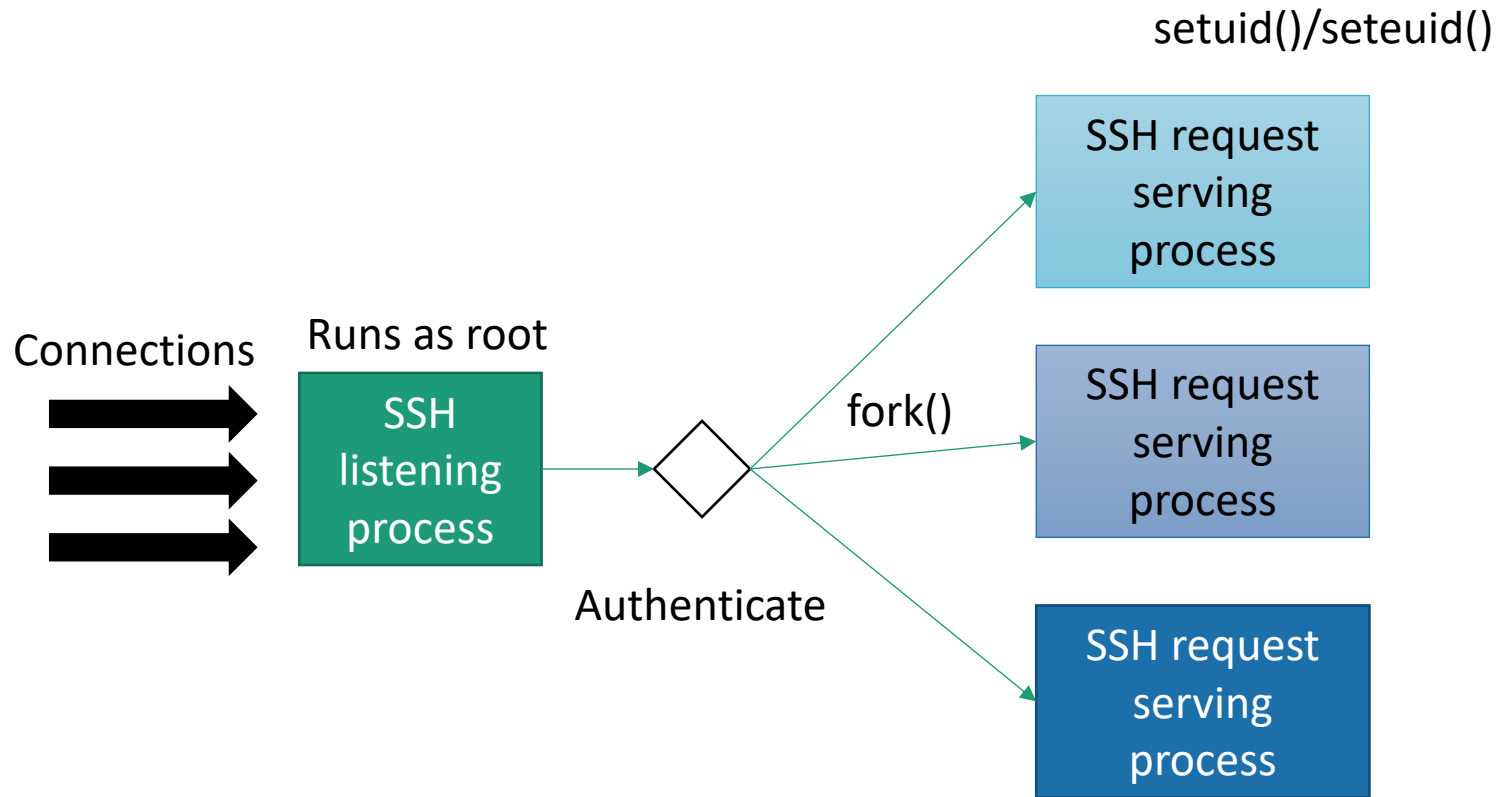
Examples: SSH, Web servers

# SSH



How is access control done here?

# SSH



**Process drop privileges and run as the authenticated user**