

Turtle  
Sec

@PATI\_GALLARDO

@PATI\_GALLARDO

# LINUX SECURITY

## and the Chromium Sandbox

PATRICIA AAS, *TurtleSec*  
BlackHoodie 2018

Turtle  
Sec

PATRICIA AAS - CONSULTANT

@PATI\_GALLARDO

*C++ Programmer, Application Security*

Currently : **TurtleSec**

Previously : Vivaldi, Cisco Systems, Knowit, Opera Software

Master in Computer Science - main language Java

Pronouns: she/her

Turtle  
Sec

Remote Code Execution  
is what browsers do.

# SANDBOXING

@PATI\_GALLARDO

- **System External Threat :** Protect the system from the vulnerabilities in the browser
- **System Internal Threat :** Protect the browser from malware on the system





# OUTLINE

- Process Architecture
- Linux Security APIs
- The Initial Sandbox
- Shrinking the Sandbox

@PATI\_GALLARDO

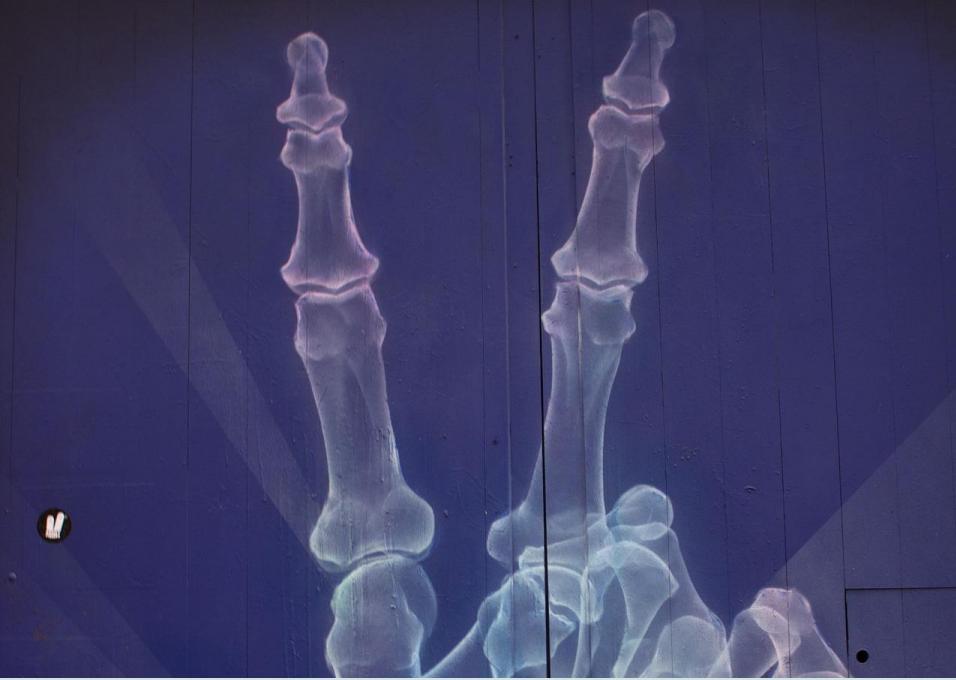
# Process Architecture

# THE EXECUTABLE FILES

- **vivaldi** : The bash script wrapper, launches vivaldi-bin. Sets up environment.
- **vivaldi-bin** : The browser binary
- **vivaldi-sandbox** : A setuid binary, not in much use today

@PATI\_GALLARDO

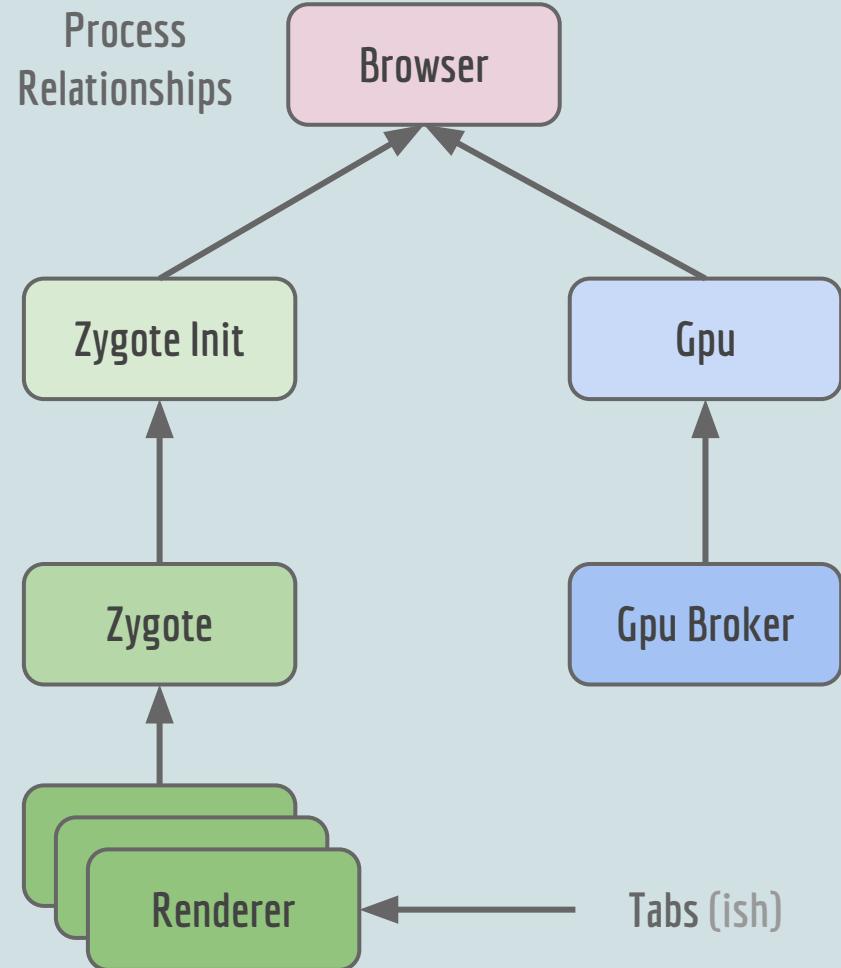




# PROCESS ARCHITECTURE

@PATI\_GALLARDO

9



# THE RUNNING PROCESSES

@PATI\_GALLARDO



1. Browser :  
vivaldi-bin
2. Zygote (Parent) :  
vivaldi-bin --type=zygote
3. Zygote (Child) :  
vivaldi-bin --type=zygote
4. Renderer (MANY) :  
vivaldi-bin --type=renderer
5. GPU (Parent) :  
vivaldi-bin --type=gpu-process
6. GPU Broker (Child) :  
vivaldi-bin --type=gpu-broker



# OUTLINE

- Process Architecture
- Linux Security APIs
- The Initial Sandbox
- Shrinking the Sandbox

@PATI\_GALLARDO

@PATI\_GALLARDO

# Linux Security APIs

# CHROME: //SANDBOX



@PATI\_GALLARDO

13

## Sandbox Status

SUID Sandbox	No
Namespace Sandbox	Yes
PID namespaces	Yes
Network namespaces	Yes
Seccomp-BPF sandbox	Yes
Seccomp-BPF sandbox supports TSYNC	Yes
Yama LSM Enforcing	Yes

You are adequately sandboxed.

Network

Web

NS(Net)

NS(Pid)

NS(User)

setrlimit

chroot

seccomp, capset

Resources

Filesystem

System Calls





# OUTLINE

- Process Architecture
- Linux Security APIs
- The Initial Sandbox
- Shrinking the Sandbox

@PATI\_GALLARDO

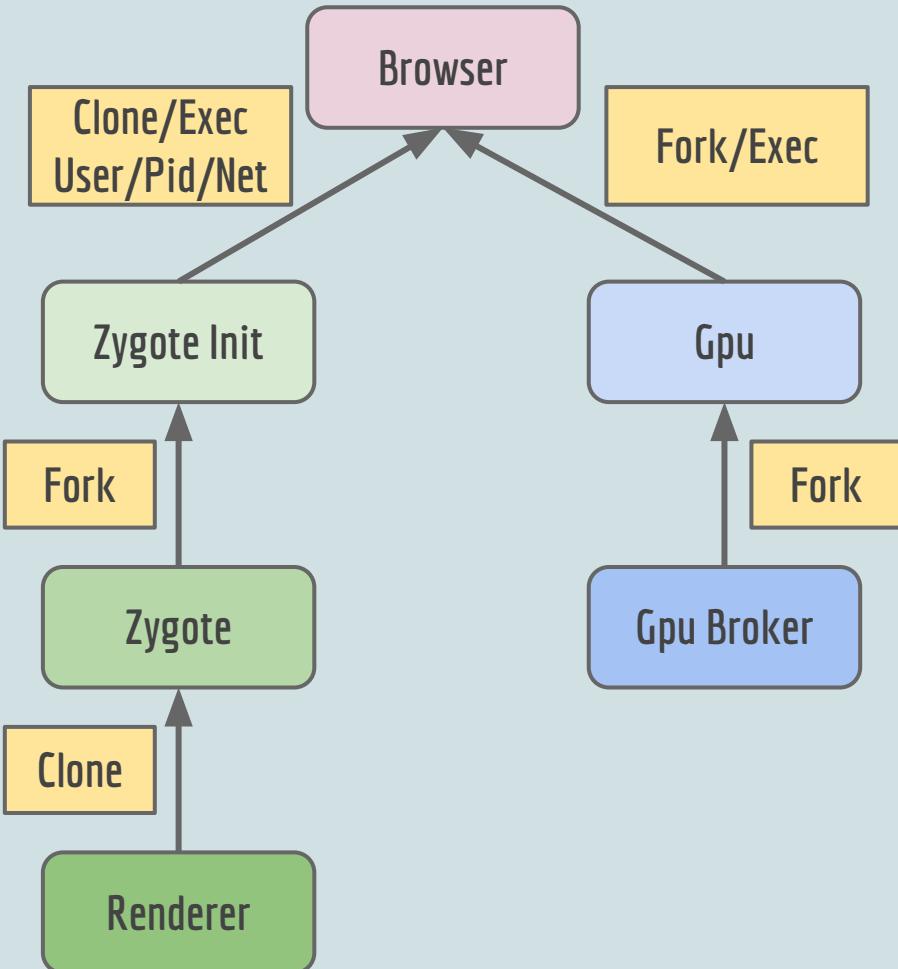
# The Initial Sandbox

# ONE BINARY



@PATI\_GALLARDO

17

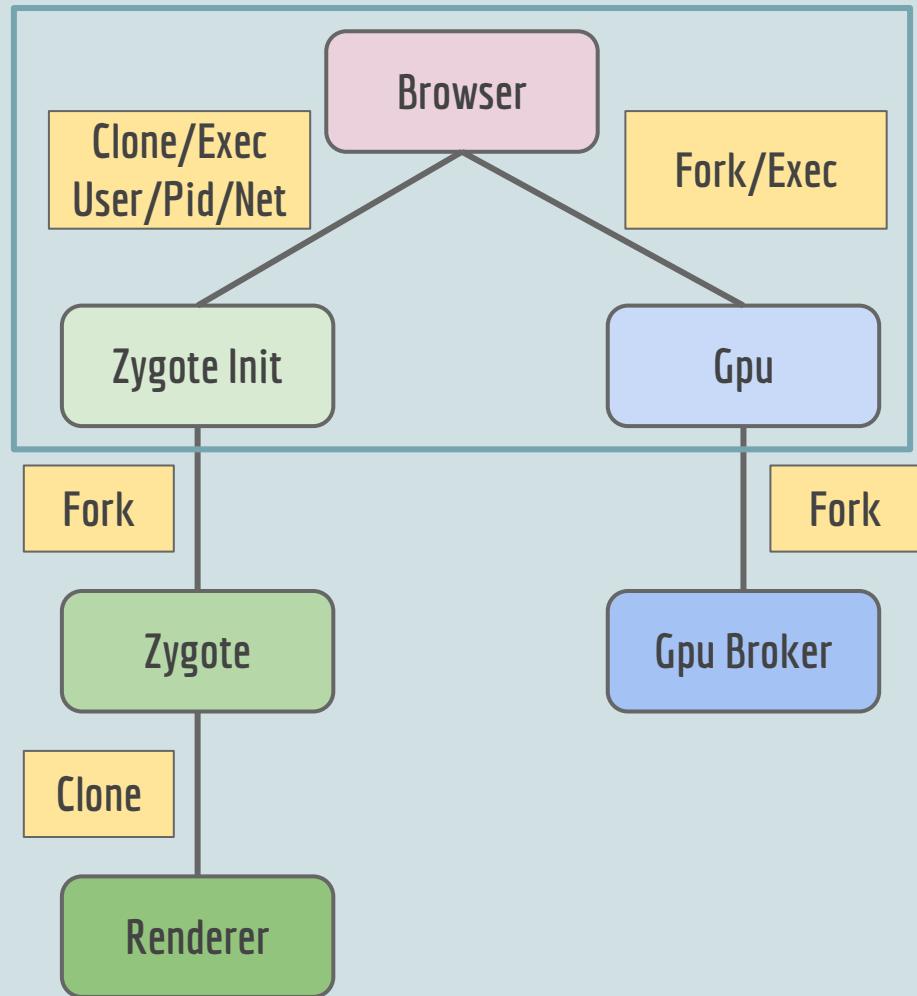


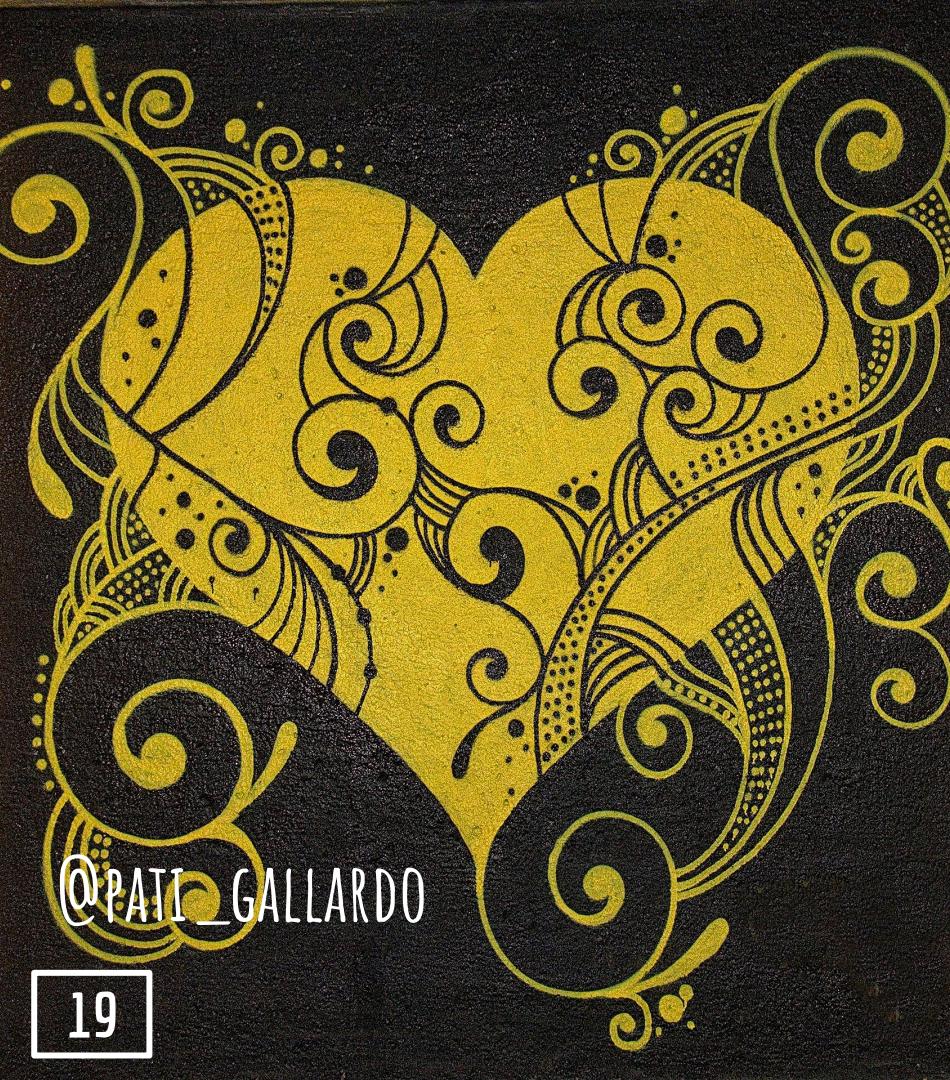
# INITIAL SANDBOX



@PATI\_GALLARDO

18





# FORK / EXEC

- A typical start of a new process on Linux is a “fork/exec”
- “Forking” is creating a new process
- “Exec” is executing a binary in a process
- “Clone” is a type of fork which can restrict a process at creation

# SANDBOXING OPPORTUNITIES: FORK

1. BEFORE FORK
2. AFTER FORK

# SANDBOXING OPPORTUNITIES: "FORK"/EXEC

1. BEFORE CLONE/FORK
2. (AT CLONE)
3. BEFORE EXEC
4. AT STARTUP (INPUT : ARGV, ENV)



# Namespaces

@PATI\_GALLARDO

# NAMESPACES (ZYGOTES/ RENDERERS)

Limits what a process can see

Api : clone/unshare



ZYGOTE + RENDERER



@PATI\_GALLARDO

24

## NAMESPACES IN USE

### CLONE\_NEWUSER

Inside a USER NS we can create a PID NS.

### CLONE\_NEWPID

Same PID number can represent different processes in different PID namespaces. One init (PID 1) process per PID NS

### CLONE\_NEWWNET

Isolate a process from network

ZYGOTE + RENDERER



@PATI\_GALLARDO

25

## AT CLONE : CREATE NAMESPACES

Clone flags define the process\* created and will create namespaces (NS) for it

1. Test which NS are available
2. Fail if not sufficient
3. Construct the biggest supported and applicable set

Emulates fork with longjmp

\* Also used to create threads

```
int flags = CLONE_NEWUSER | CLONE_NEWPID | CLONE_NEWNET;
jmp_buf env;
if (setjmp(env) == 0) {
    return CloneAndLongjmpInChild(flags, ptid, ctid, &env);
}
```

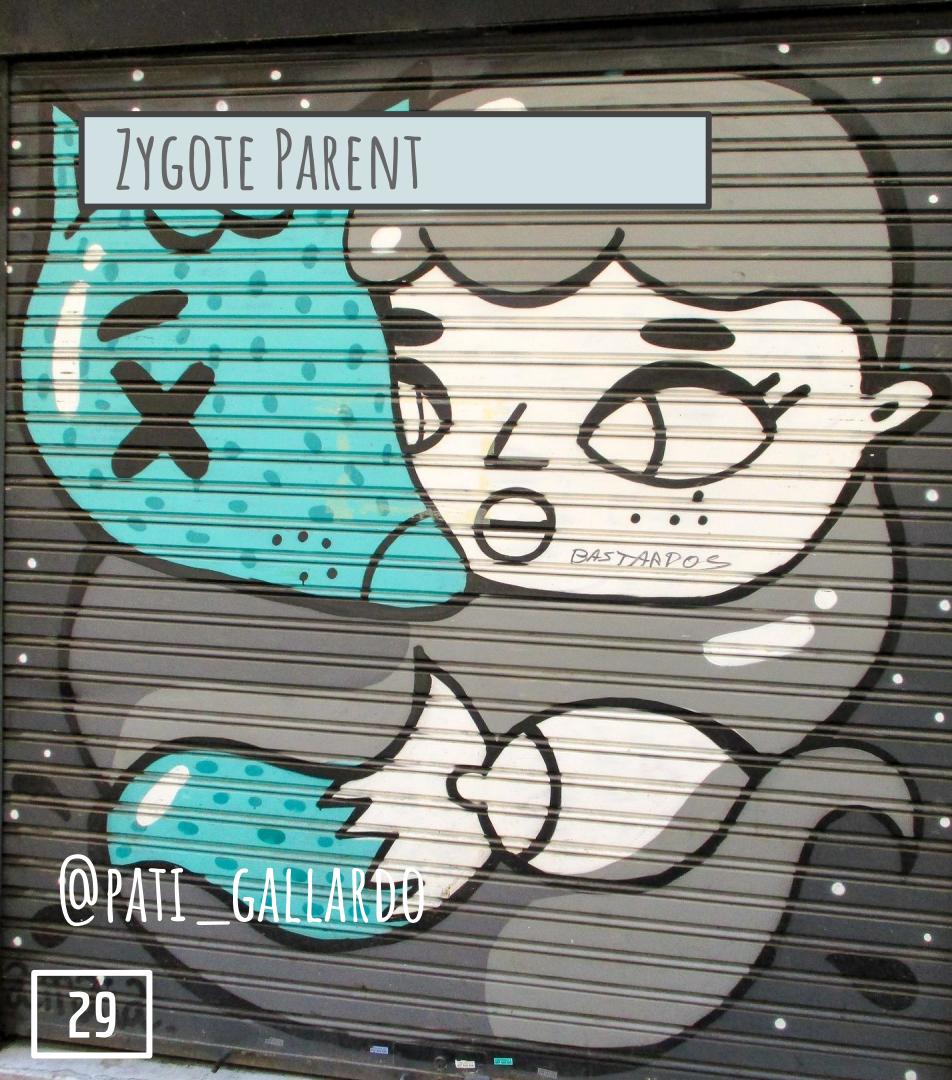
chromium/base/process/launch\_posix.cc

```
pid_t CloneAndLongjmpInChild(unsigned long flags,  
                           pid_t* ptid,  
                           pid_t* ctid,  
                           jmp_buf* env) {  
    char stack_buf[PTHREAD_STACK_MIN];  
    void* stack = stack_buf + sizeof(stack_buf);  
    return clone(&CloneHelper,  
                stack, flags, env, ptid, nullptr, ctid);  
}
```

chromium/base/process/launch\_posix.cc

```
int CloneHelper(void* arg) {
    jmp_buf* env_ptr = reinterpret_cast<jmp_buf*>(arg);
    longjmp(*env_ptr, 1);
    // Should not be reached
    assert(false);
    return 1;
}
```

chromium/base/process/launch\_posix.cc



## UNSHARE - USER NAMESPACE

### **unshare(CLONE\_NEWUSER)**

Alternative to `clone(CLONE_NEWUSER)`, will not create new process, but move caller.

Credentials::CanCreateProcessInNewUserNS

# SANDBOXING OPPORTUNITIES: "FORK"/EXEC

1. BEFORE CLONE/FORK

2. AT CLONE

3. BEFORE EXEC

4. AT STARTUP (INPUT : ARGV, ENV)



## BEFORE EXEC : LAUNCH OPTIONS



@PATI\_GALLARDO

31

# PREPARE FOR EXEC

1. Fix the environment
2. Fix file descriptors
3. Fix signal handling
4. Set up process group
5. Maximize resource limits
6. Set PR\_SET\_NO\_NEW\_PRIVS
7. Change current dir
8. Select executable path
9. Setup command-line



# OUTLINE

- Process Architecture
- Linux Security APIs
- The Initial Sandbox
- Shrinking the Sandbox

@PATI\_GALLARDO

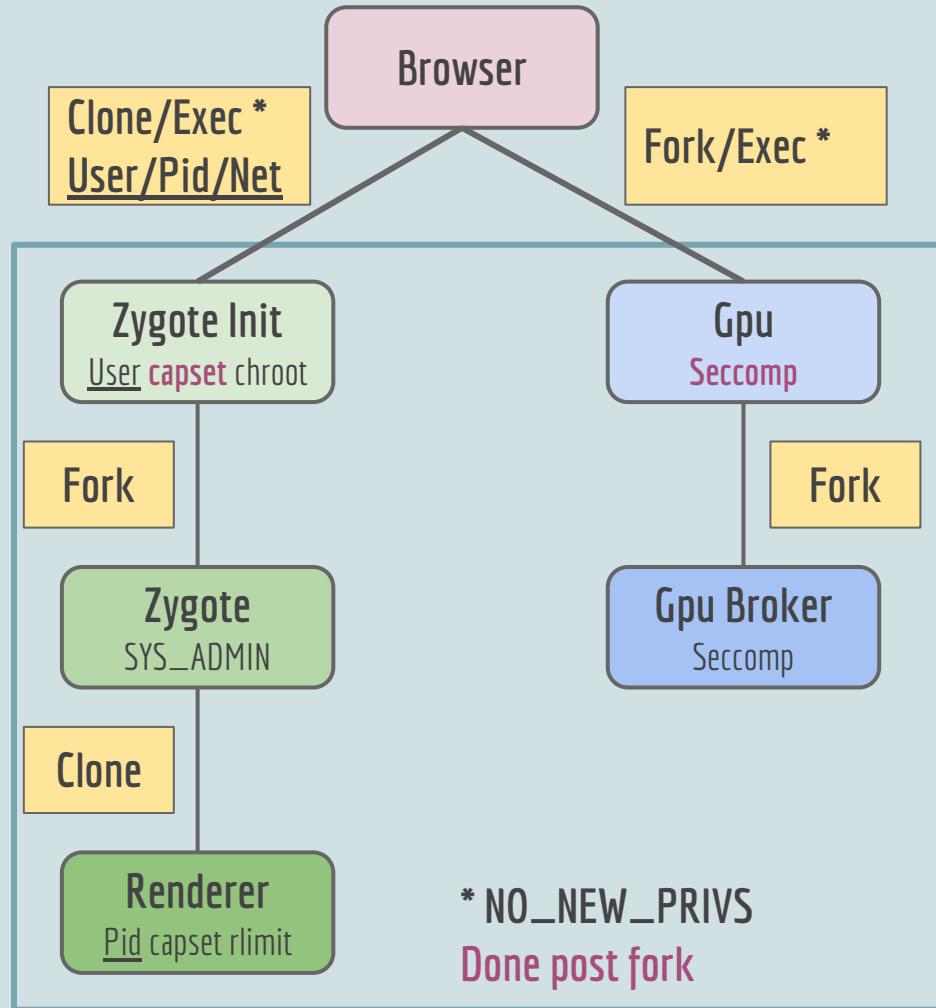
# Shrinking the Sandbox



## SHRINKING THE SANDBOX



34



ALL EXCEPT BROWSER



PR\_SET\_NO\_NEW\_PRIVS

prctl(PR\_SET\_NO\_NEW\_PRIVS)

Preserved across fork, clone and execve

"If no\_new\_privs is set, then operations that grant new privileges (i.e. execve) will either fail or not grant them. This affects suid/sgid, file capabilities, and LSMs."

/usr/include/linux/prctl.h



# Seccomp

@PATI\_GALLARDO

36

S  
@PATI\_GALLARDO

# SECCOMP

(RENDERERS/GPU/BROKER)

Limits which syscalls a process  
can call and how these calls are  
handled

Api : seccomp



@PATI\_GALLARDO



## SECCOMP BPF PROGRAM

Program written in an assembly-like language to filter system-calls.

Runs in a simple VM in kernel space. All syscalls will be filtered by this program

**TSYNC**: Once a Seccomp Program is installed it applies to all threads in a process



# SECCOMP : BPF POLICIES

## BPF Program defined in a Policy

Fundamentally a whitelist, allows a set of syscalls and has custom handling of others.

An extended Policy is generally more permissive

1. BaselinePolicy

- 1.1 GpuProcessPolicy

- 1.1.1 GpuBrokerProcessPolicy

- 1.2 RendererProcessPolicy

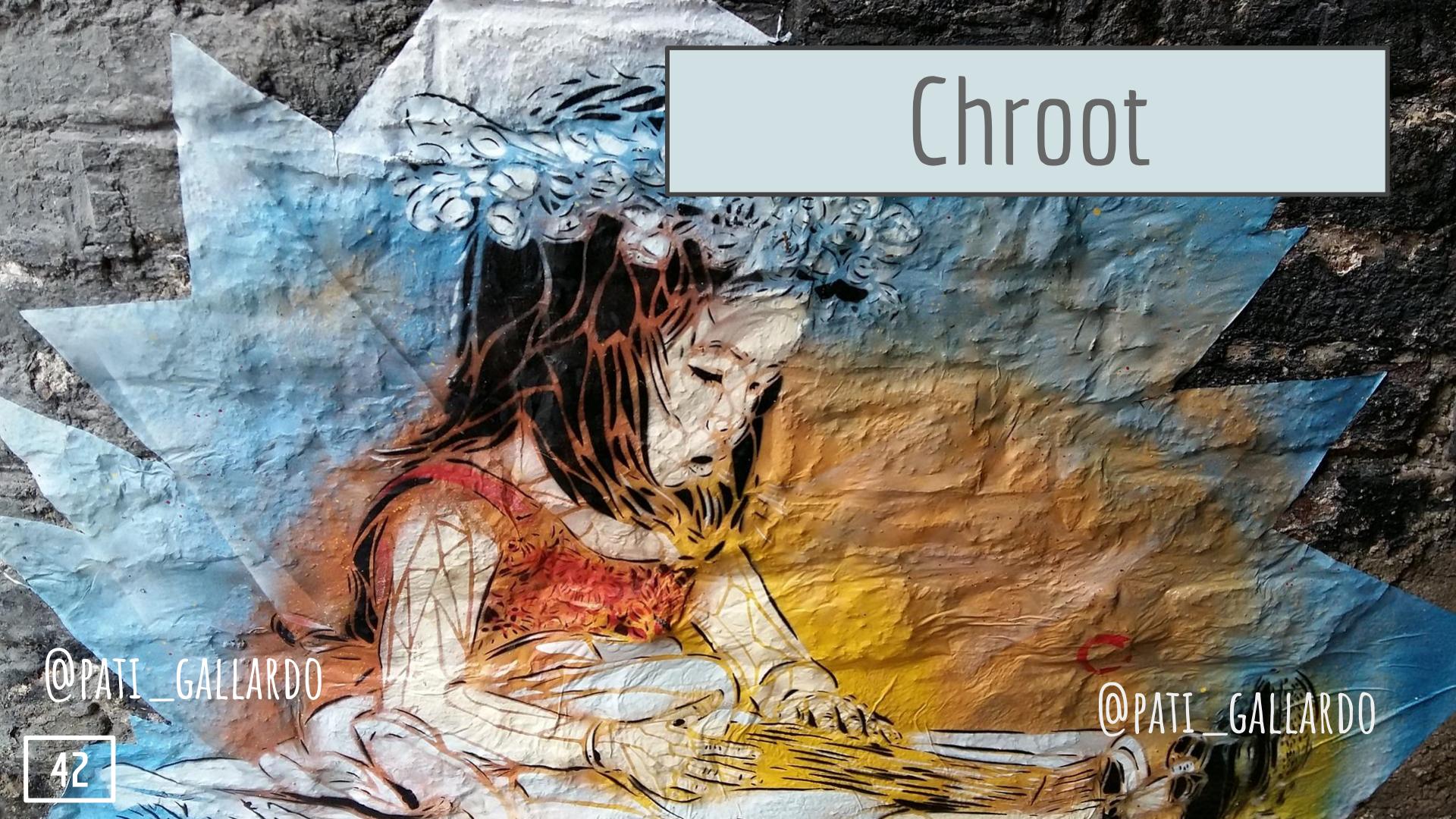
```
void StartSandboxWithPolicy(sandbox::bpf_dsl::Policy* policy)
{
    SandboxBPF sandbox(policy);
    assert(sandbox.StartSandbox());
}

bool SandboxBPF::StartSandbox() {
    InstallFilter();
    return true;
}
```

chromium/services/service\_manager/sandbox/linux/sandbox\_seccomp\_bpf\_linux.cc  
chromium/sandbox/linux/seccomp-bpf/sandbox\_bpf.cc

```
void SandboxBPF::InstallFilter() {
    CodeGen::Program program = AssembleFilter();
    struct sock_filter bpf[program.size()];
    const struct sock_fprog prog =
        { static_cast<unsigned short>(program.size()), bpf };
    memcpy(bpf, &program[0], sizeof(bpf));
    assert(prctl(PR_SET_NO_NEW_PRIVS, 1, 0, 0, 0) == 0);
    assert(seccomp(SECCOMP_SET_MODE_FILTER,
                   SECCOMP_FILTER_FLAG_TSYNC, &prog) == 0);
}
```

chromium/sandbox/linux/seccomp-bpf/sandbox\_bpf.cc



# Chroot

@PATI\_GALLARDO

42

S  
©PATI\_GALLARDO

# CHROOT

## (ZYGOTES/ RENDERERS)

Limits what a process can see of  
the filesystem

Api : chroot





ZYGOTES + RENDERER

@PATI\_GALLARDO

44

## CHROOT : DROP ACCESS TO FS

- clone(CLONE\_FS) a child process
- Child chroot("/proc/self/fdinfo/")
- Child immediately does a chdir("/")
- Child does \_exit(kExitSuccess)

You can see this by looking at

ls -l /proc/<pid>/root

Of the Zygote or any ancestor

Credentials::DropFileSystemAccess

```
bool ChrootToSafeEmptyDir() {
    pid_t pid = -1;
    char stack_buf[PTHREAD_STACK_MIN];
    void* stack = stack_buf + sizeof(stack_buf);
    int clone_flags = CLONE_FS | LINUX_SIGCHLD;
    pid = clone(ChrootToSelfFdinfo, stack, clone_flags, nullptr);
    int status = -1;
    assert(HANDLE_EINTR(waitpid(pid, &status, 0)) == pid);
    return WIFEXITED(status) && WEXITSTATUS(status) == kExitSuccess;
}
```

chromium/sandbox/linux/services/credentials.cc

```
int ChrootToSelfFdinfo(void*) {
    assert(chroot("/proc/self/fdinfo/") == 0);
    assert(chdir("/") == 0);
    _exit(kExitSuccess);
}
```

chromium/sandbox/linux/services/credentials.cc



# Capabilities

@PATI\_GALLARDO

S  
©  
@PATI\_GALLARDO

# CAPABILITIES

(ZYGOTE/ RENDERERS)

Limits what a process is  
privileged enough to do

Api : capset





ZYGOTES + RENDERERS

@PATI\_GALLARDO

49

## DROP CAPABILITIES

Uses `capset()` to drop all or some capabilities

“Linux divides the privileges traditionally associated with superuser into distinct units, known as capabilities, which can be independently enabled and disabled.”

*Man page for capabilities*

Credentials::DropAllCapabilities



# Resource Limits

@PATI\_GALLARDO

# RESOURCE LIMITS

## (RENDERERS)

Limits the access this process  
has to platform resources

Api : setrlimit





RENDERER

@PATI\_GALLARDO

52

## RESOURCE LIMITS : SETRLIMIT

Limits using setrlimit:

1. **RLIMIT\_AS** : Maximum size of the process' virtual memory (address space) in bytes
2. **RLIMIT\_DATA** : Maximum size of the process's data segment

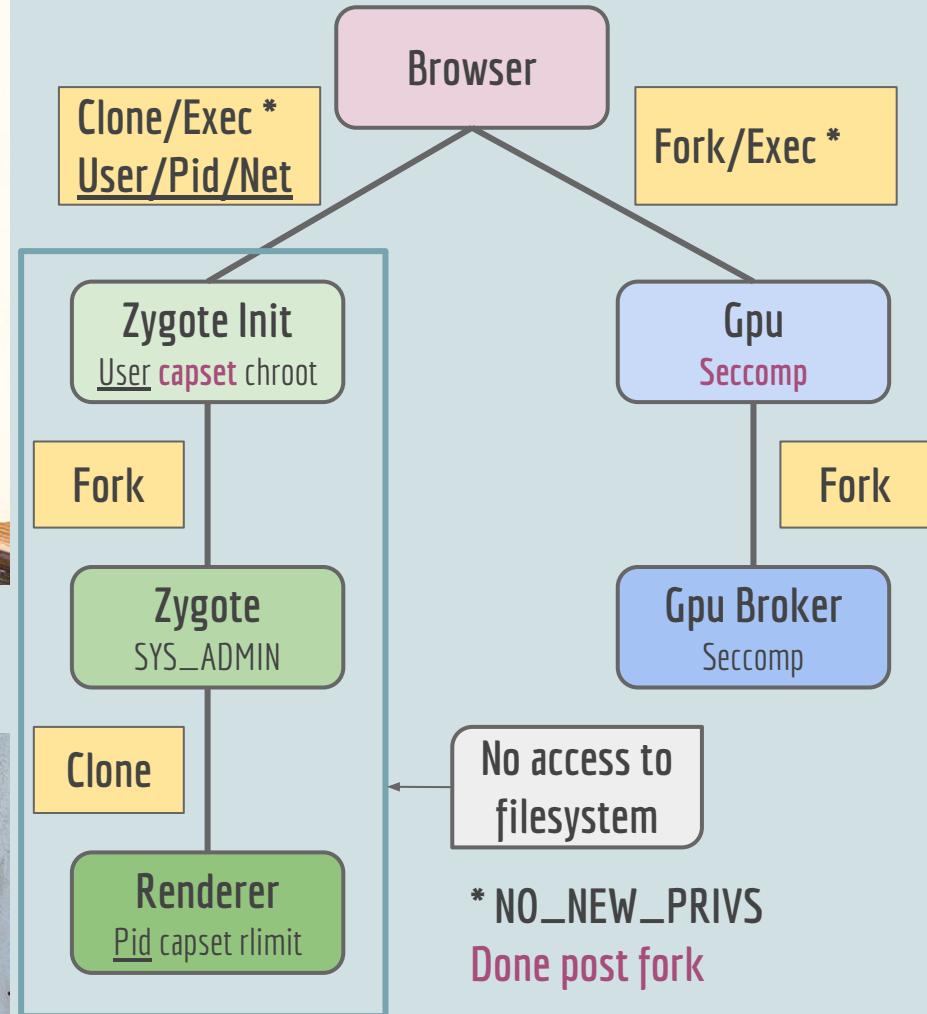
[LinuxSandbox::LimitAddressSpace](#)



TRUST IS RELATIVE

@PATI\_GALLARDO

53



# CHROME: //SANDBOX



@PATI\_GALLARDO

54

## Sandbox Status

SUID Sandbox	No
Namespace Sandbox	Yes
PID namespaces	Yes
Network namespaces	Yes
Seccomp-BPF sandbox	Yes
Seccomp-BPF sandbox supports TSYNC	Yes
Yama LSM Enforcing	Yes

You are adequately sandboxed.

@PATI\_GALLARDO

# SOURCES

Chromium/Kernel source + Linux Man Pages + lwn.net

MICHAEL KERRISK

Book: *The Linux Programming Interface*

Course: *Linux Security and Isolation APIs*

ALL ERRORS ARE MY OWN

@PATI\_GALLARDO

*Patricia Aas, Consultant  
TurtleSec*

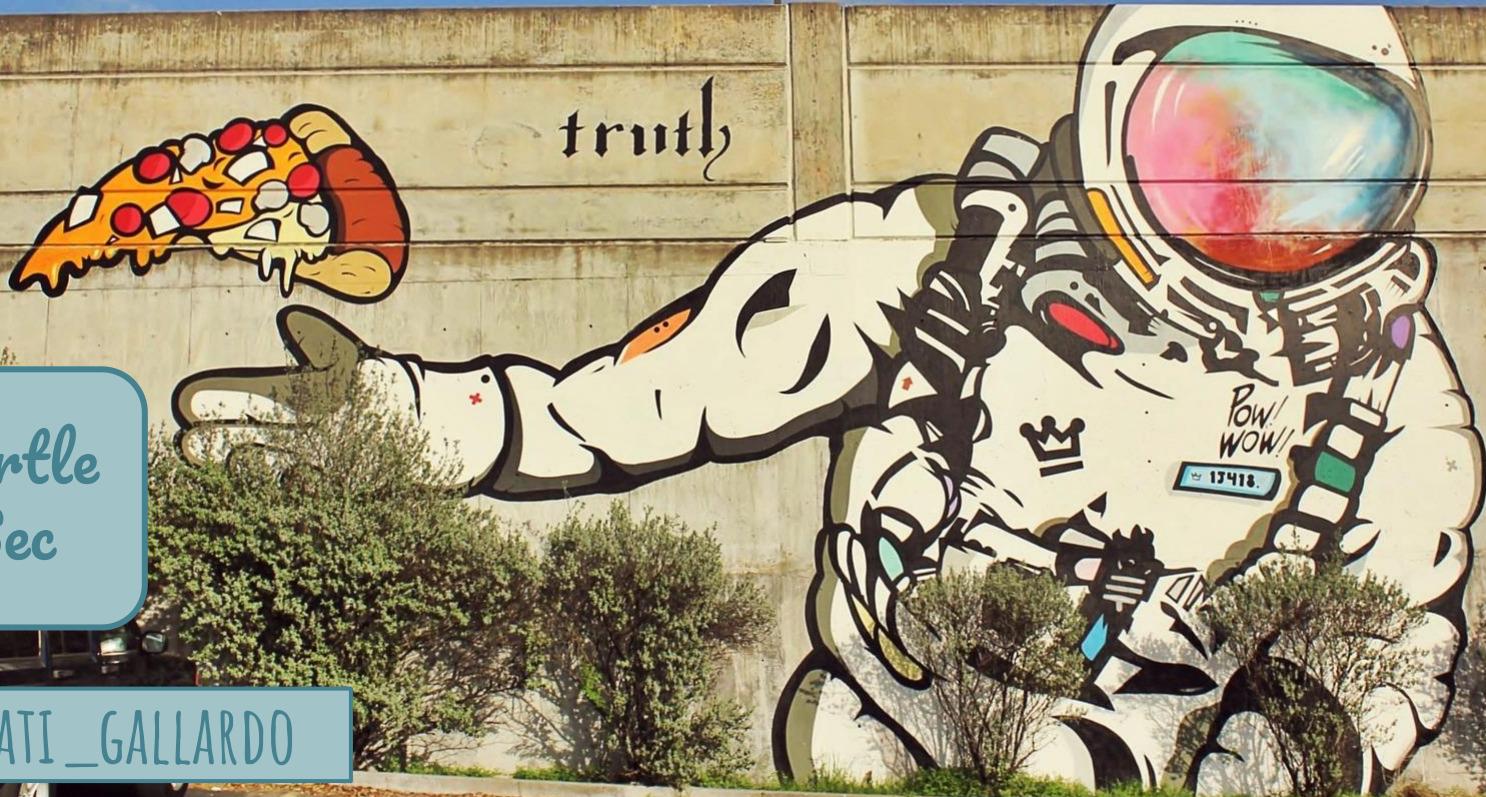
C++ and Application Security

Turtle  
Sec

@PATI\_GALLARDO

Turtle  
Sec

@PATI\_GALLARDO



@PATI\_GALLARDO



## APPENDIX / SOME NOTES

Other APIs  
in use

Not Used Much in Chromium, but...

# YAMA LSM

Limits the access that other process have to this process - especially ptracing

Status is checked by reading :  
`/proc/sys/kernel/yama/ptrace_scope`

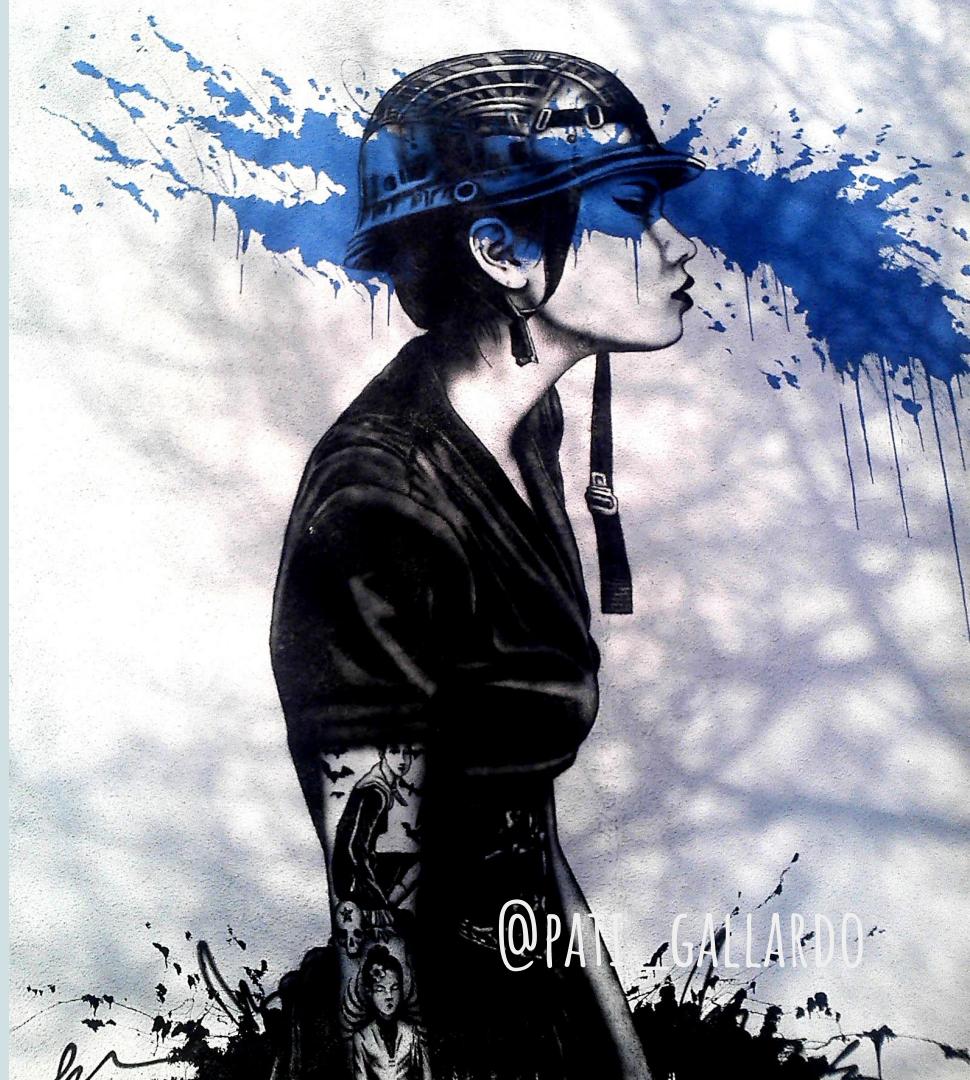


# SETUID / SETGID

## (LEGACY)

Increases what a process is privileged enough to do, by using the privilege of the executables owner

Api : set\*uid / set\*gid



# CGROUPS (CHROMEOS)

Limits the resources available to a process. Used in ChromeOS - not covered here.

/proc/<PID>/cgroup

/proc/cgroups



# PROCESS GROUPS (CHROMEDRIVER)

Can be used to treat a group of processes as one.

Used with the 'detach' property of Chromedriver

Api : setpgid



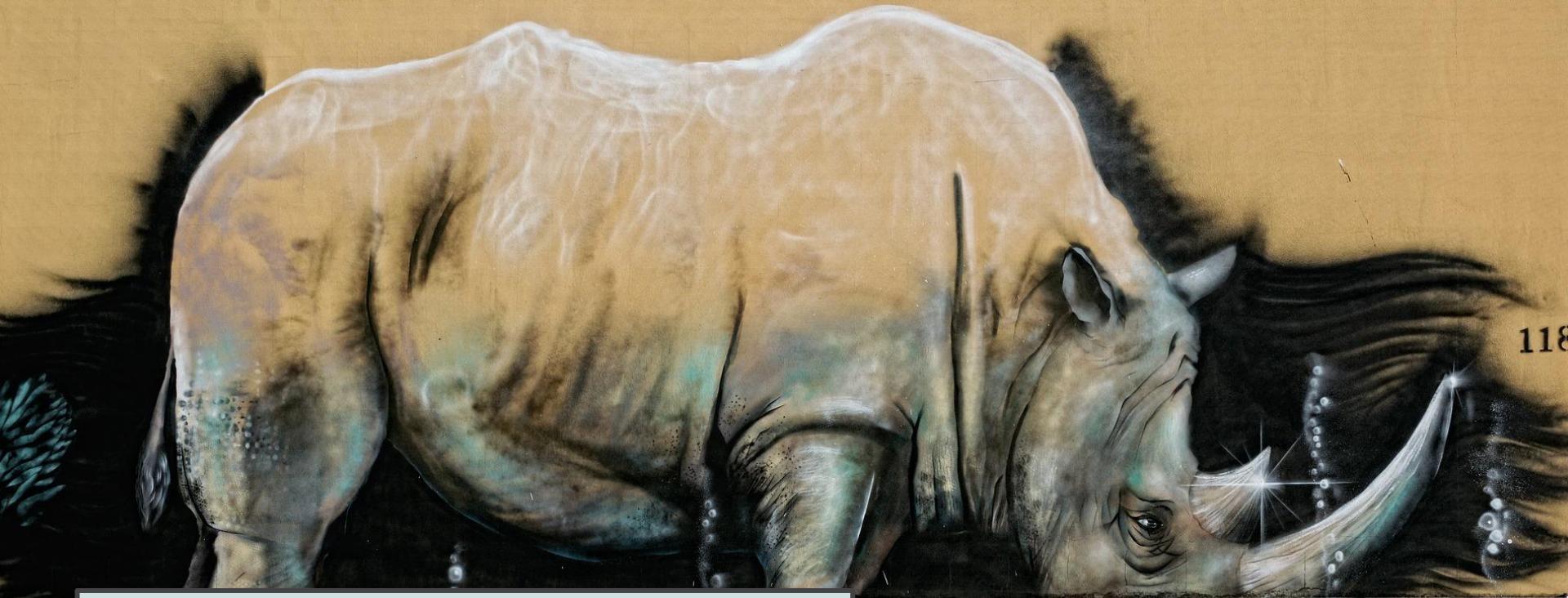
# Challenges



JARVS.

# Debugging

@PATI\_GALLARDO



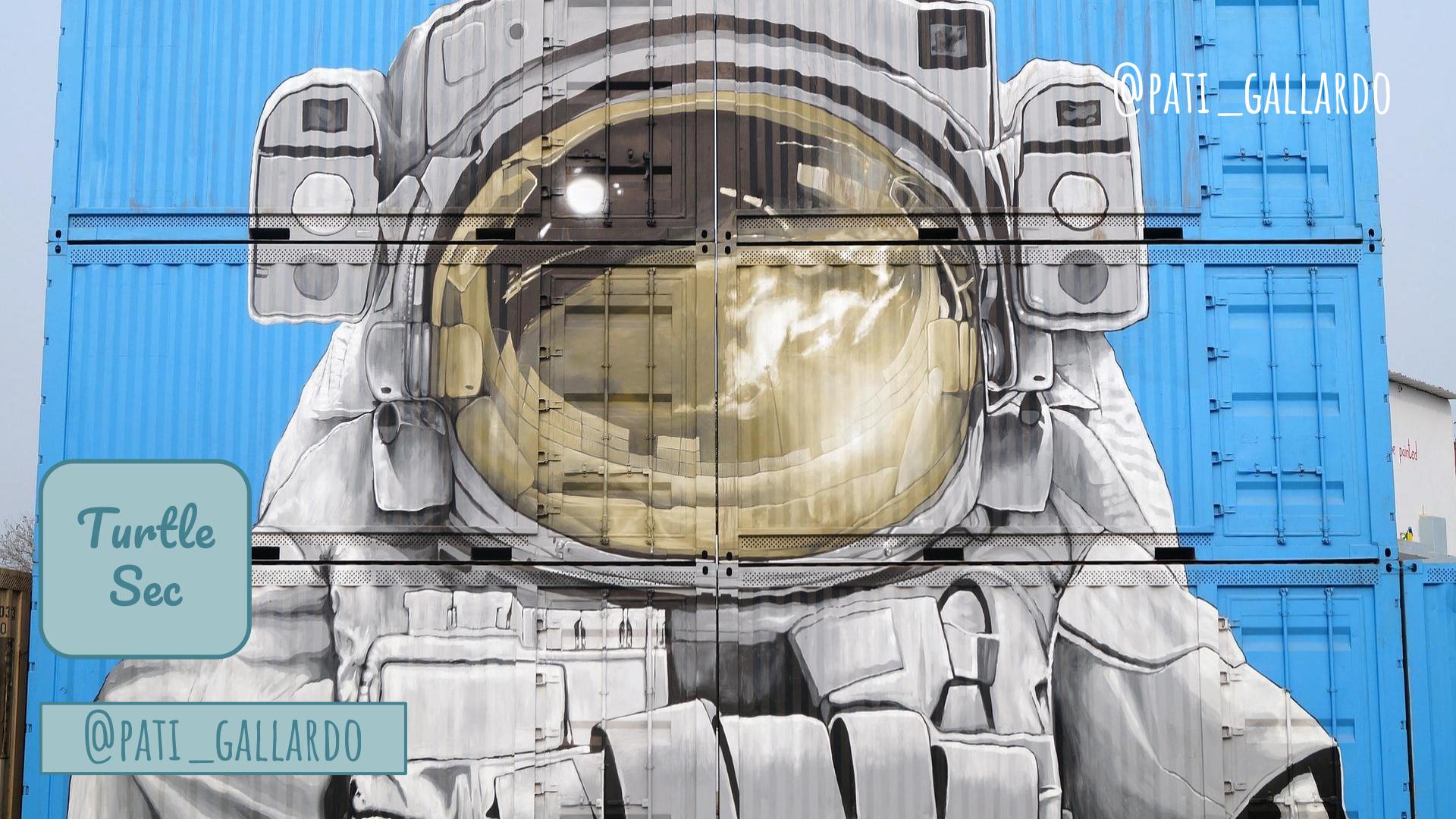
# Crash Reporting

@PATI\_GALLARDO

Photos from pixabay.com

Patricia Aas, *TurtleSec*  
@pati\_gallardo

*Turtle  
Sec*

A large-scale mural of a sea turtle, painted on the side of a blue shipping container. The turtle is depicted in a three-quarter view, facing right. Its shell is a light tan or yellowish color with dark, irregular patterns. The head and front flippers are white, while the back flippers and the lower part of the body are grey. The mural is highly detailed, showing textures like scales and fabric. The background is a solid blue. In the top right corner of the mural, there is a small signature that reads "e painted".

@PATI\_GALLARDO

Turtle  
Sec

@PATI\_GALLARDO