



Be Kind, Please Rewind

Adventures in creating a macOS record/replay debugger

LEGAL



REWIND <<

Why?

```
NSURLSessionDataTask *task = [_session
    dataTaskWithRequest:request
    completionHandler:^(NSData *_Nullable data, NSURLResponse *_Nullable response,
        NSError *_Nullable err) {
        ...
        *stop = YES;
```



```
Thread 2 Crashed:: Dispatch queue:
com.apple.NSXPCCConnection.user.com.google.santa.metricservice.63335
0  libobjc.A.dylib      objc_msgSend + 29
1  Foundation            -[NSError copyWithZone:] + 107
2  santametricservice    -[SNTMetricHTTPWriter write:toURL:error:] +
1372
3  santametricservice    -[SNTMetricService exportForMonitoring:] + 475
```


Record/Replay: Prior Art

- PANDA (2020)
 - Whole system!
- WinDbg (2017)
- RR (2014)
- Scribe (2010)
- Jockey (2005)
- Flashback (2004)
- ReTrace
- QuickRec
- Revirt (1999)
 - Whole system!

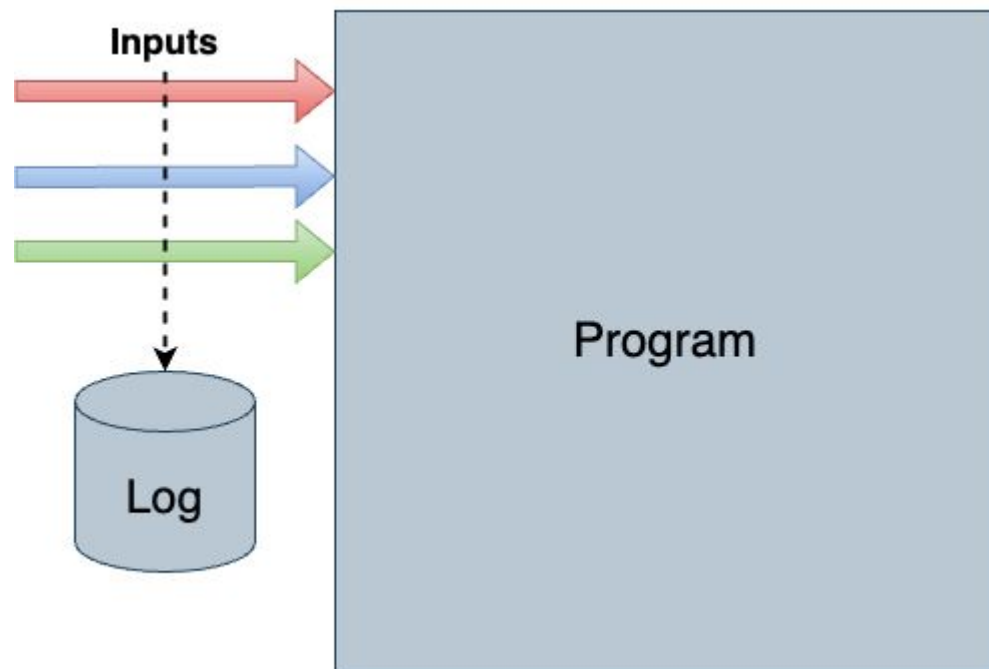
Record/Replay: Prior Art

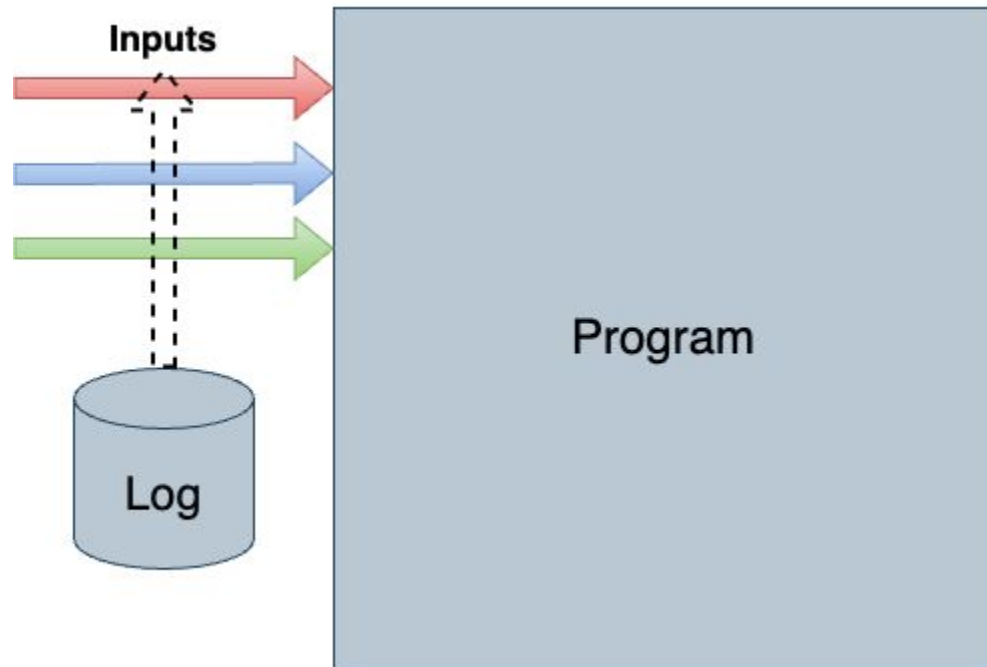
- PANDA (2020)
 - Whole system!
- WinDbg (2017)
- RR (2014)
- Scribe (2010)
- Jockey (2005)
- Flashback (2004)
- ReTrace
- QuickRec
- Revirt (1999)
 - Whole system!



A collection of vintage audio and video equipment is displayed on a wooden surface. In the upper left, a silver cassette deck holds a cassette tape with a yellow label that reads "110+10 60". To its right is a large, silver VCR with its front-loading door open, revealing the internal tape compartment. Below these items is a black remote control with numerous buttons, including a small LCD screen and a "TRANSPORT" button. The equipment is arranged in a way that suggests a focus on audio and video recording and playback technology.

Record / Replay Basics & Goals












Goals for our tool

- Only focusing on user-space programs
- Easy to use and deploy – needs to support a stock MBP with/M1,M2
- No DBI / code instrumentation
- Small investment of effort to maintain
- Fast enough to use on real programs

RR's Requirements for User-Space Replay

Requirement	Does macOS Meet This Out of the Box?
Ability to Record Syscalls	
Ability to Record Syscalls Outside Libc	
Ability to determine if a Syscall is blocking	
Ability to Intercept Signals	

RR's Requirements for User-Space Replay (Part 2)


Requirement	Does macOS Meet This Out of the Box?
Ability to pin a process to a single core (cpuset)	
Ability to trap non-deterministic instructions	
Ability to access reliable and deterministic hardware performance counters	

Recording



Recording: It's not that simple...

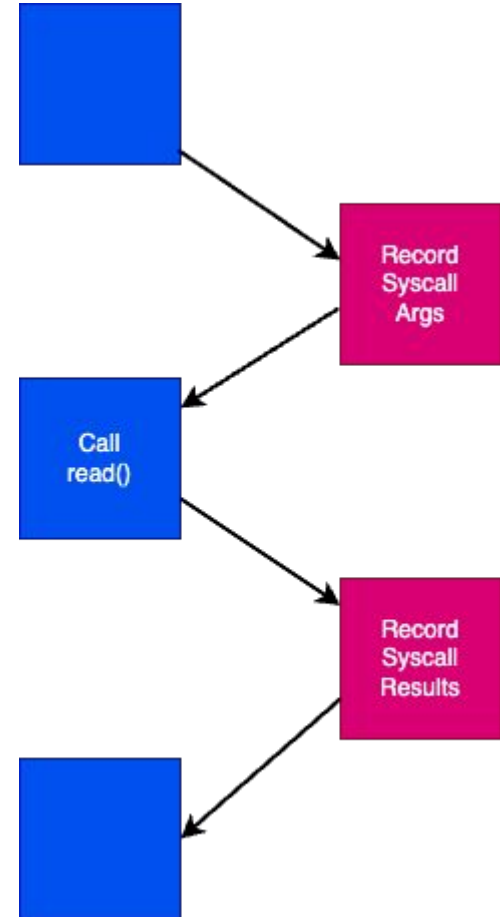
- Mach traps
 - Close enough to syscalls, no big deal
 - ... except for a few traps which don't have normal hook points
- Signals
 - The outside world can “asynchronously” poke the target
- mmap
- Multithreading
 - Well-formed programs shouldn't have issues (data races), but...
 - Aside: thanks Apple for not giving us cpuset - no easy way to pin to one core
- Commpage
 - Similar to vvar (normally accessed via vDSO) on Linux
- Non-deterministic instructions – mrs x0, cntvct_el0

A collage of vintage electronic equipment. At the top left is a silver cassette deck with a tape inside. To its right is a silver VCR with its top cover open, revealing internal components. Below these is a black remote control with numerous buttons. The background is a warm, textured surface.

Recording: Syscalls & Traps

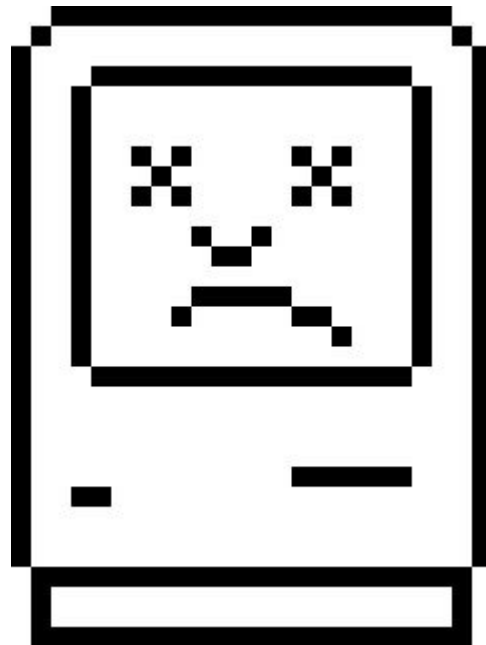
Recording: Syscalls on macOS

- 3 types – BSD, mach traps, machine dependent
- Need pre- and post-hooks for data gathering



Thanks Apple

- Gutted ptrace implementation – no sysemu
- No seccomp-bpf equivalent



One Option: dtrace

- dtrace hooks, storage, etc.
- Not enough to capture arbitrary syscall data though
 - No conditionals for example - not possible to **switch** in “multiplexed” syscalls
- Strictly async
 - How to pause so userland can get what it needs?
 - Luckily there are “destructive” actions
 - `signal(STOP)`
 - `stop()` - `mach_task_suspend`
 - These only take effect *after* the syscall is processed though...

One Option: dtrace

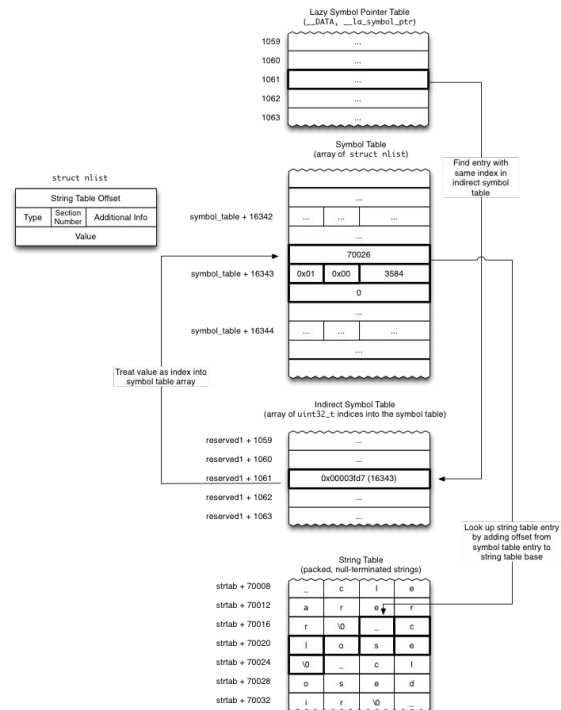



Seatbelt / Sandbox?

- Seatbelt is wired up into every syscall maybe?
- Trace mode for recording
 - Not a good API, minimal log entries
- No way to not kill on replay

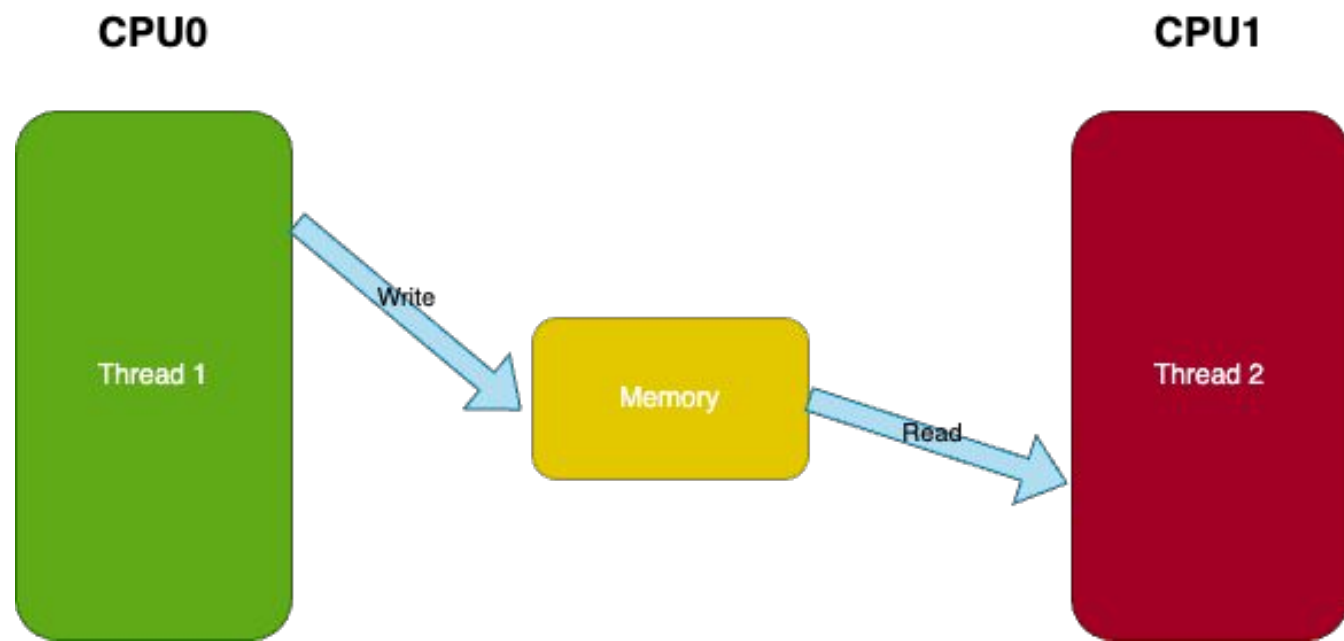
Interposing / Dynamic Interposing / Symbol Rebinding

- macOS is a BSD!
 - ABI compatibility is at the libc level not the kernel
 - Can we just hook `libsystem_kernel`?
- We can interpose on the symbol
 - Could use [fishhook](#)
- Doesn't catch direct syscalls...



A collection of vintage electronic equipment is displayed on a wooden surface. In the upper left, a silver cassette deck is shown with a cassette tape inserted; the tape's label reads "SONY 10+10 60". To its right is a silver VCR with its front-loading door open, revealing the internal tape compartment. Below these items is a black remote control with numerous buttons, including "VTR", "CLOCK", "MEMORY", "AUDIO RECORDER", "REVIEW", "PLAY", "PAUSE/STILL", "STOP", "REC", and "SEARCH". The background is a warm-toned wooden surface.

Recording: Dealing with Data Races



How Does RR Handle This?

- Only runs one thread to run at a time (non-parallel)
- Limits threads to the same core using processor affinity
- Schedules threads and records the choice in the log (**can mixup order on replay to find bugs**)

CPU0

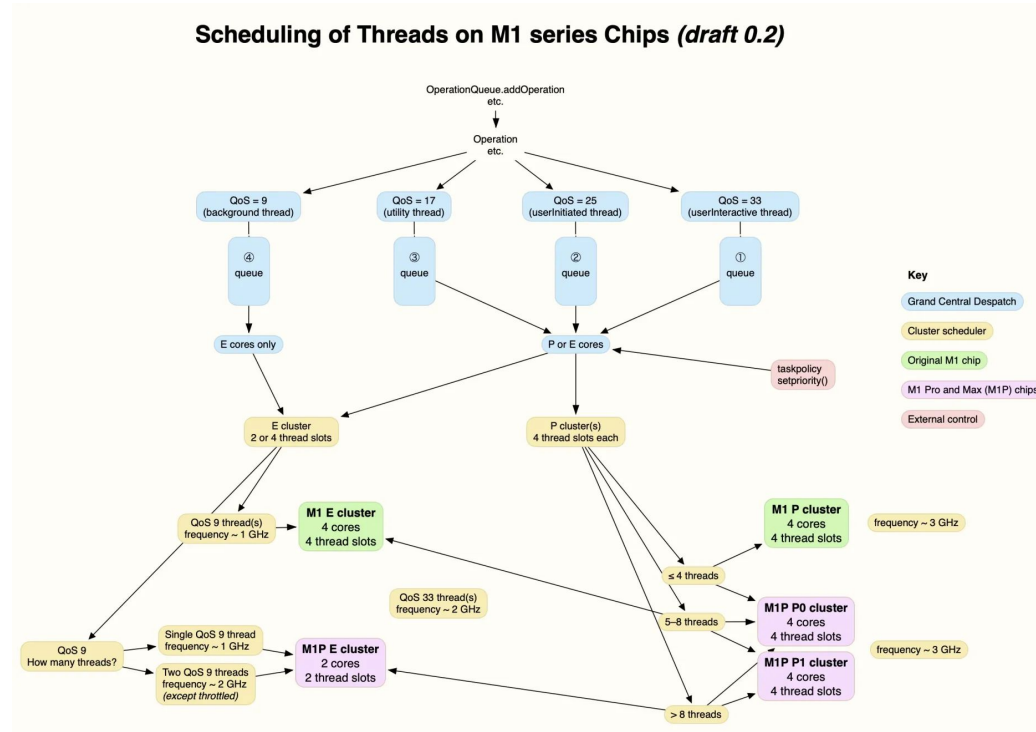


Thread scheduling on macOS not guaranteed

- No `cpu_set(3)`
- Can we use `THREAD_AFFINITY_POLICY`?

```
/*  
 *      thread_bind:  
 *  
 *      Force the current thread to execute on the specified processor.  
 *      Takes effect after the next thread_block().  
 *  
 *      Returns the previous binding.  PROCESSOR_NULL means  
 *      not bound.  
 *  
 *      XXX - DO NOT export this to users - XXX  
 */  
processor_t  
thread_bind(  
    processor_t      processor)  
{  
    thread_t      self = current_thread();  
    processor_t    prev;
```

P-cores and E-cores



From: <https://eclecticlight.co/2022/01/13/scheduling-of-processes-on-m1-series-chips-first-draft/>

Can we shutdown cores?

- In the old OS X internals books there was an example showing how to shutdown cores using `processor_exit`
- Can we just limit ourselves to a core?

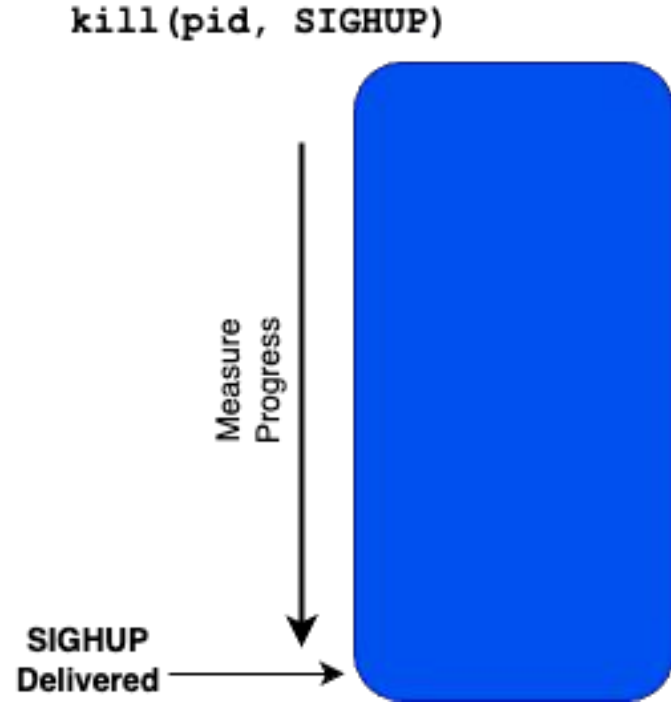
```
[ user@waterville ~ ]
$ sudo ./print_processors
Password:
Number of processors: 12
CPU: slot 0 (master)
CPU: slot 1
//snipped.
CPU: slot 11
[ user@waterville ~ ]
$ sudo ./processor_xable
processor_exit: (os/kern) service not supported
```

A collage of vintage electronic equipment. In the top left is a silver cassette deck with a tape inside. To its right is a silver VCR with its top cover open, revealing internal components. Below these is a black remote control with numerous buttons. The background is a warm, textured surface. A semi-transparent black rectangle with white text is centered over the middle of the image.

Recording: Asynchronous Events

Signals & Scheduling

- Need to be able to intercept signals and record register state of where the signal was delivered or program interrupted for scheduling.
- Need to know where you are in the programs execution so you can inject your signals in the right place during replay
- Replay: when using something interrupt driven must account for late firing interrupts



Using PMUs from macOS

- RR works on Asahi Linux and uses the PMU can we?
 - Uses the count of retired conditional branches as progress indicator (0x8c)
 - Can reset for an interrupt when replaying
- macOS does not have an interface for setting PMUs from EL0

```
[ user@waterville ~/src/pmu_counters ]  
$ sudo ./counter_test  
loaded db: a15 (Apple A15)  
number of fixed counters: 2  
number of configurable counters: 8  
counters value:  
    cycles: 41865278  
    instructions: 91998218  
    branches: 21071096  
    branch-misses: 53779  
[ user@waterville ~/src/pmu_counters ]  
$ sudo ./counter_test  
loaded db: a15 (Apple A15)  
number of fixed counters: 2  
number of configurable counters: 8  
counters value:  
    cycles: 41946121  
    instructions: 92093331  
    branches: 0  
    branch-misses: 0
```

panic(cpu 5 caller 0xfffffe0017c66cd8): kperf: timer fired
at 2793246644070, but sampling is disabled
@kptimer.c:328
Debugger message: panic

Supporting Nondeterministic Execution in Fault-Tolerant Systems*

J. Hamilton Slye

Dept. of Electrical and Computer Engineering
Carnegie Mellon University
ham+@cmu.edu

E.N. Elnozahy

Department of Computer Science
Carnegie Mellon University
mootaz@cs.cmu.edu

Abstract

We present a technique to track nondeterminism resulting from asynchronous events and multithreading in log-based rollback-recovery protocols. This technique relies on using a software counter to compute

with the end users [11]. Efficient tracking of nondeterminism is thus crucial to supporting interactive applications [14].

Different flavors of logging have been suggested with different performance and resilience character-

Decrement Register

Branch to handler if register = 0

Options without PMU or DBI

- We can count the number of syscalls and then single step forward then inject the signal (set a breakpoint and invoke the signal handler)
- Do what scribe(10) does and simply deliver the signal at the next syscall and replay interrupted syscall (special case for signals like SIGSEGV that originate in user space.)
- If we need to go further than say 10,000 instructions we can use an high res clock (e.g. pacman) to trap back to us

```
...// set up the registers for our counter
...asm!(
    ...    "mov x3, {running}",
    ...    "mov x4, x0", // store backed up copy of
    ...    "ldr x0, [x0]", // load the value of coun
    ...    "1:",
    ...    "isb",
    ...    "ldr x2, [x3]", // load the value of run
    ...    "sub x0, x0, 1", // decrement the counter
    ...    "cbz x2, 456f", // break the loop if we'
    ...    "cbnz x0, 1b", // if we've not underflo
    ...    "123:",
    ...    "mov {interrupted}, 1", // set interrupte
    ...    "456:",
    ...    "str xzr, [x3]", // set running to false
    ...    "str x0, [x4]", // store the value of co
    ...    in("x0") count_addr,
    ...    interrupted = out(reg) interrupted,
    ...    running = in(reg) running_addr,
    ...);
```



A collage of vintage audio and video equipment. At the top left is a cassette tape with a yellow label that reads "10+10 15". To its right is a VCR with its top cover open, revealing internal components. Below these is a black remote control with various buttons labeled "VTR", "CLOCK", "MEMORY", "AUDIO RECORDER", "REVIEW", "PLAY", "EJECT", "PAUSE/STILL", "STOP", "REC", "FIND", "STILL", "SEARCH", and "TRANSPORT". To the right of the remote is a VHS tape with a label that reads "VHS" and "Hi-Fi REC LEVEL". The word "Darling" is written in white text across the center of the image.

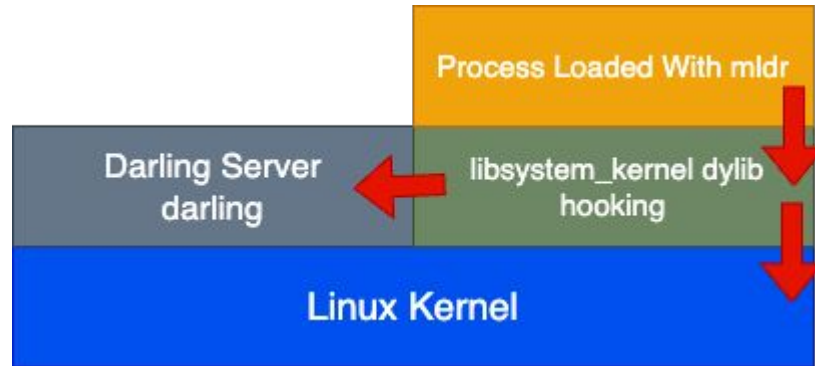
Darling

Darling

- “A Translation Layer that lets you run macOS software on Linux”
- Uses a custom loader, interposing of `libsystem_kernel`, a lot of duct tape code and userland a server to translate macOS syscalls to Linux syscalls
- Can run software like xcode on Linux



High Level: How Darling Works

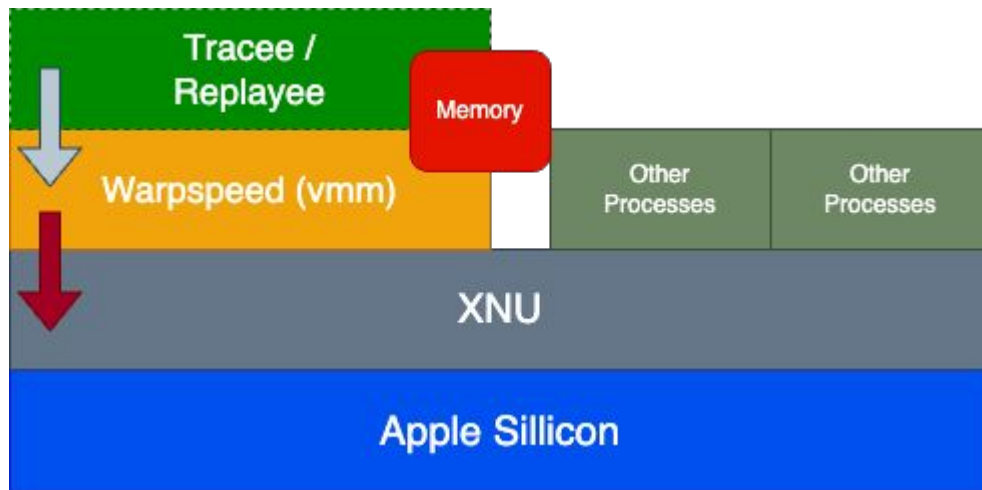


A collage of vintage audio and video equipment. At the top left is a silver cassette tape with a yellow label that reads "10+10 60". To its right is a large, dark-colored VCR with its top cover open, revealing internal components like the tape compartment and reels. Below the VCR is a black VHS tape with a label that says "VHS" and "Hi-Fi REC LEVEL". In the foreground, a black remote control with numerous buttons is visible. The word "Warp speed" is written in white, stylized font across the center of the image, overlaid on a semi-transparent black rectangular background.

Warp speed

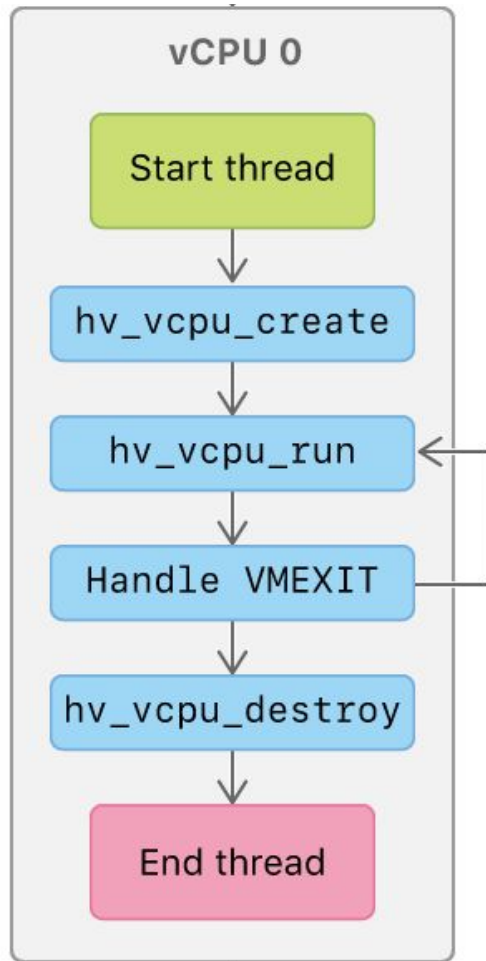
Warpspeed

- Isolate target inside a VM with 1 core
- Proxy syscalls
- Both signal slide + SoftPMU to approximate program progression
- Manual thread scheduling



Hypervisor.framework

- Super light-weight framework
 - Little as possible in the kernel
- Usage:
 - Create a VM
 - Map memory (from hypervisor address space)
 - Create vCPU
 - Set regs
 - Run
 - Trap out to userland on VM exit
 - GOTO 5
 - *That's it*

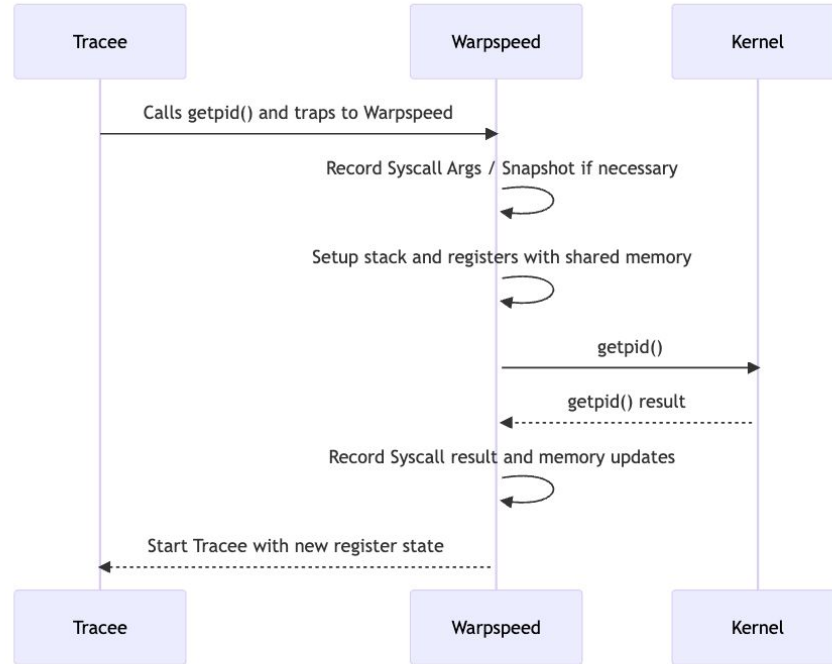


Warpspeed: VM/Hypervisor

- Use modified darling's loader (mldr) to map in target program and dyld
 - Load in shared cache
 - “Share” an address space with the guest
 - 1:1 map the regions of the loaded target into VM at the same virtual address
 - Trap out and forward syscalls
-
- All based on Hyperpom (Rust!)
-
- Lets us control the execution of the program perfectly
 - Only have one virtual core
 - Manually schedule threads



Warpspeed: VM/Hypervisor



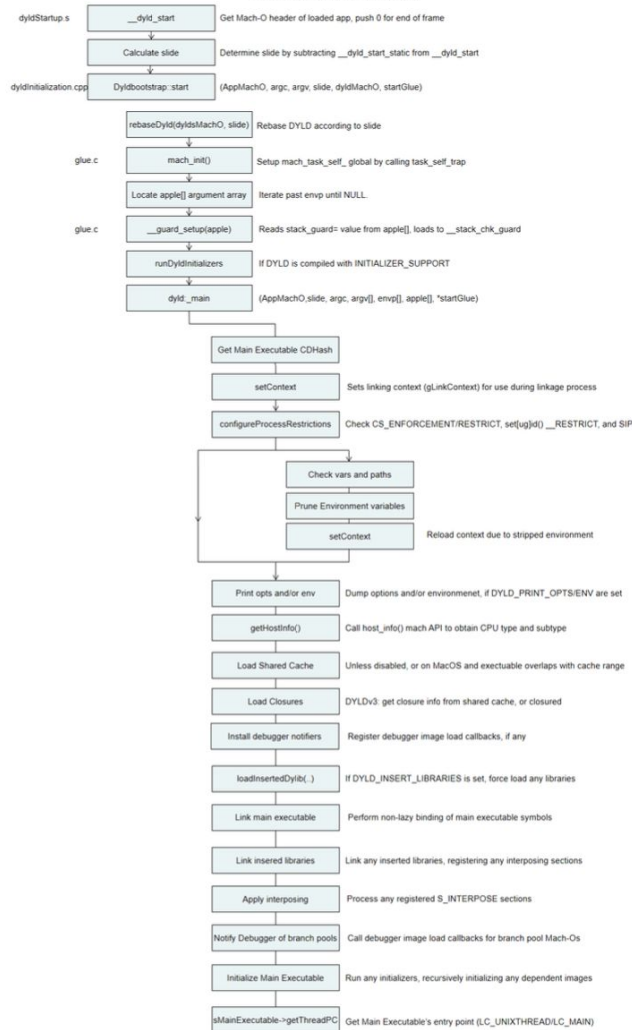

```

[2023-06-09T03:19:56Z] [INFO] [warp-speed] 287: Forwarding syscall 90x[0x1, 13400000, 9a7, 0, 0, 0, 0, 9a8, 0, 101, 00000000, 1aef, 13400000, 20b0fffa0, 20b0fffa0]
[2023-06-09T03:19:56Z] [INFO] [warp-speed] Changed page: [ ]
[2023-06-09T03:19:56Z] [INFO] [warp-speed] Returning 0x0 40000300
[2023-06-09T03:19:56Z] [INFO] [warp-speed] ELR.EI: 0x28010000
[2023-06-09T03:19:56Z] [INFO] [warp-speed] Incoming syscall 0x0[0x, 13400000, 13400000, 0, 00000002, 2, 2, 9a9, 103, 103, 0, 27, 13400000, 1, 13400000, 00000002]
[2023-06-09T03:19:56Z] [INFO] [warp-speed] 288: Forwarding syscall 60x[0x, 13400000, 13400000, 0, 00000002, 2, 2, 9a8, 101, 103, 0, 27, 13400000, 1, 13400000, 00000002]
[2023-06-09T03:19:56Z] [INFO] [warp-speed] Changed page: [ ]
[2023-06-09T03:19:56Z] [INFO] [warp-speed] Returning 0x0 40000300
[2023-06-09T03:19:56Z] [INFO] [warp-speed] ELR.EI: 0x28010000
[2023-06-09T03:19:56Z] [INFO] [warp-speed] Incoming syscall 0x0[00000000, 10000000, 0, 1010010e, 0, 10, 0, 9a9, 103000000, 0, 120020020, 1, 1, 13400000, 20b0fffa0, 20b0fffa0]
[2023-06-09T03:19:56Z] [INFO] [warp-speed] openC:/System/Windows/System/Library/FeatureFlags/GlobalWindowsOverrides.plist
[2023-06-09T03:19:56Z] [INFO] [warp-speed] 289: Forwarding syscall 50x[0x00000000, 10001000, 0, 1010010e, 0, 10, 0, 9a8, 1010010e, 0, 120020020, 1, 1, 13400000, 20b0fffa0, 20b0fffa0]
[2023-06-09T03:19:56Z] [INFO] [warp-speed] Changed page: [ ]
[2023-06-09T03:19:56Z] [INFO] [warp-speed] Returning 2x0 40000300
[2023-06-09T03:19:56Z] [INFO] [warp-speed] ELR.EI: 0x28010000
[2023-06-09T03:19:56Z] [INFO] [warp-speed] Incoming syscall 0x0[1010010e, 10000000, 0, 1010010e, 0, 10, 0, 9a9, 103000000, 0, 0000000000000000, 120020020, 1, 1, 13400000, 20b0fffa0, 20b0fffa0]
[2023-06-09T03:19:56Z] [INFO] [warp-speed] openC:/System/Windows/System/Library/FeatureFlags/WindowsOverrides.plist
[2023-06-09T03:19:56Z] [INFO] [warp-speed] 290: Forwarding syscall 50x[0x00000000, 10001000, 0, 1010010e, 0, 10, 0, 9a8, 1010010e, 0000000000000000, 120020020, 1, 1, 13400000, 20b0fffa0, 20b0fffa0]
[2023-06-09T03:19:56Z] [INFO] [warp-speed] Changed page: [ ]
[2023-06-09T03:19:56Z] [INFO] [warp-speed] Returning 2x0 40000300
[2023-06-09T03:19:56Z] [INFO] [warp-speed] ELR.EI: 0x28010000
[2023-06-09T03:19:56Z] [INFO] [warp-speed] Incoming syscall 0x0[00000000, 10000000, 0, 1010010e, 0, 10, 0, 9a9, 103000000, 0, 0000000000000000, 120020020, 1, 1, 13400000, 20b0fffa0, 20b0fffa0]
[2023-06-09T03:19:56Z] [INFO] [warp-speed] openC:/System/Library/FeatureFlags/SettingsOverrides.plist

```


dyld

Figure 7-2: The flow of dyld



Warpspeed: Unimplemented Features

- LLDB/GDB interface
- Optimizing/compressing log format
- The hypervisor itself is responsible for performing the syscalls
 - What happens on a blocking call?
 - Could deadlock on mutex wait
- Handling blocking syscalls
 - Manually enumerate and perform some non-blocking alternative
 - or...

```

20.14 g eqv pMMAAe-Cl + CuCl2 -> CuCl2 + 2.0
20.14 g eqv pMMAAe-Cl + CuCl2 -> CuCl2 + 2.0

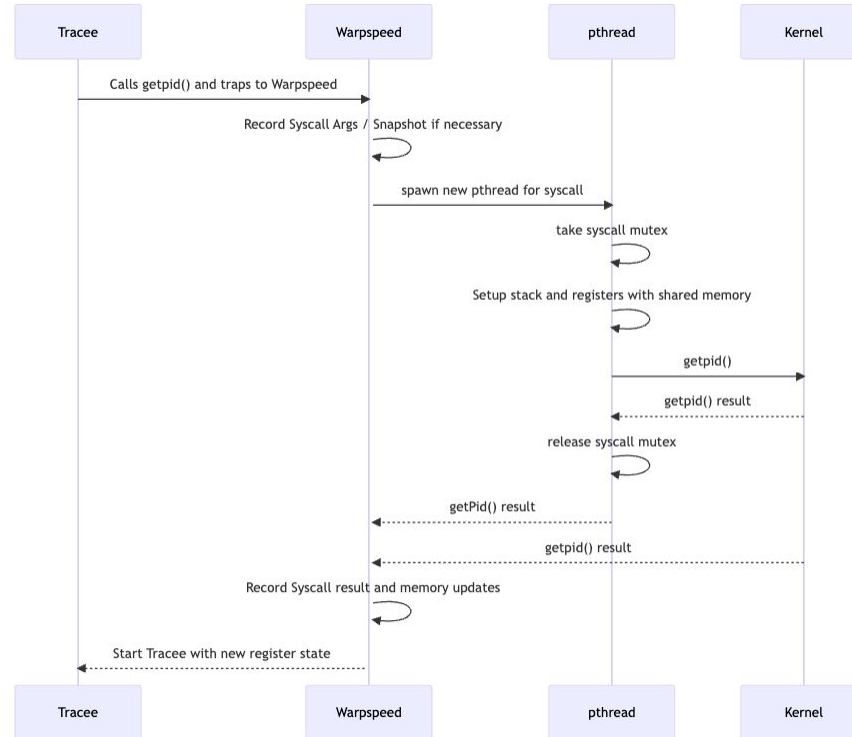
```

```

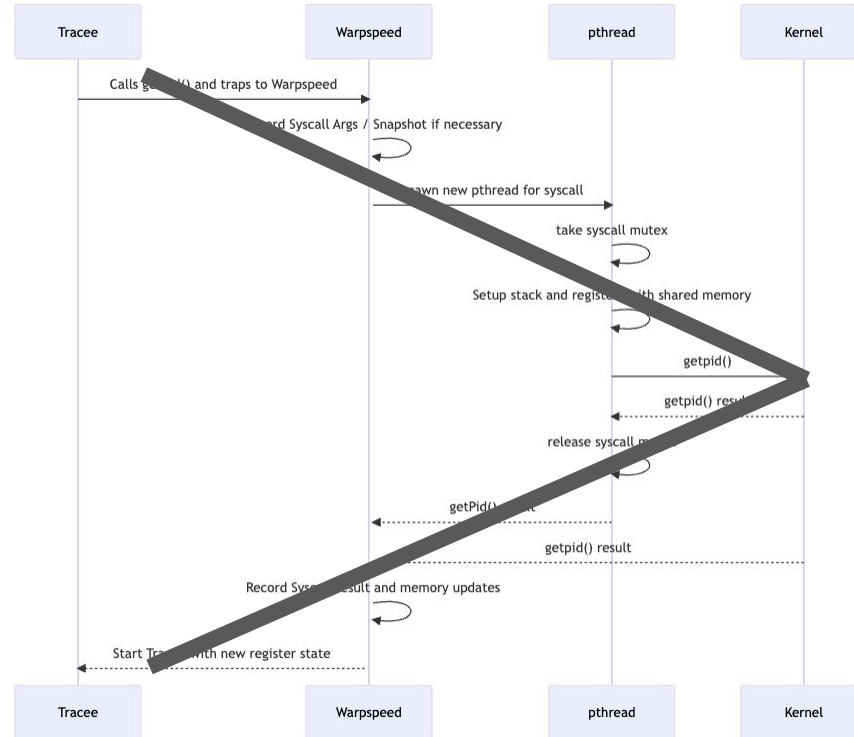
$ cd /usr/src/yaml-0.1.7/build -E make ./libyaml.so.0.1.7
Threads: 188760
[process/debug_blocking.rs MZ: "1.7"] = "1.7"
[process/debug_blocking.rs MZ: event = fd_buf {
  timer_pos: 1459945136,
  arg: 1
    0,
    0,
    0,
    0,
    64%,
  },
  debug_id: 20071546,
  uid: 0,
  moved: 0,
}
[process/debug_blocking.rs MZ: "1.7"] = "1.7"
[process/debug_blocking.rs MZ: event = fd_buf {
  timer_pos: 1459945136,
  arg: 1
    0,
    0,
    0,
    0,
    64%,
  },
  debug_id: 20071546,
  uid: 2,
  moved: 0,
}

```

Warpspeed: VM/Hypervisor




Warpspeed: VM/Hypervisor



WarpSpeed: Outstanding Issues

- MMIO
- Entitlements

A collage of vintage VCR and VHS tape components. At the top left is a VHS tape with a yellow label that reads "10+10 15". To its right is the internal mechanism of a VCR, showing the tape reels and transport mechanism. Below these is a close-up of a VCR's control panel, featuring buttons for "CLOCK", "PRIORITY", "FREQ", "TRACKING", "NEXT", "MEMORY", "AUDIO RECORDING", "SP/LP", "REVIEW", "PLAY", "EJECT", "PAUSE/STILL", "STOP", "REC", and "VTR". To the right of the control panel is a VHS tape with a label that reads "Hi-Fi REC LEVEL 01".

That's Only Half the Battle

Replay

- If you can figure out recording, replay is much simpler
 - Set breakpoints where something happened in recording
 - Mimic side-effects
 - Continue
- SoftPMU needed here in case we end up with an async event in a hot loop

Replay: GUI

- UI is core to macOS
- How can we “pass through” events on replay to the OS (to see the app running) while not introducing nondeterminism?
 - *In theory* it will “just work”
 - No (easy) way to show the UI on replay though

Recap

- Tool is WIP
- But principles work!
- Stay posted for more

A collage of vintage VCR and video cassette tape images. In the top left, a VHS tape is shown with its label. To its right is a close-up of a VCR's internal mechanism. Below these, a VCR control panel with various buttons like 'REWIND', 'PLAY', and 'STOP' is visible. In the bottom right, another VCR control panel is shown. A semi-transparent dark rectangle is overlaid in the center, containing the word 'Questions?'.

Questions?