# How Does Debugger Work A glance on ptrace

By Yuan Liu

# **Topics**

"Focus on lightweight debugger strace and trivial userspace debugger like GDB.
The code samples will refer to Linux Kernel 4.6 x86-64.

# Web version can be found here: <a href="https://lkwq007.github.io/how-does-debugger-work/">https://lkwq007.github.io/how-does-debugger-work/</a>

#### What's covered

- how strace work
  - system calls tracing
- how debuggers work
  - single-step debugging
  - breakpoint
  - access/modify debuggee's memory and registers
- practical uses of ptrace for debuggers

#### What's not covered

- usage of strace
- options of ptrace
- how debugger figures out where to find the functions and variables in the machine code
  - (it's a compiler issue, the debugging information should be provided by the compilers, and then used by debuggers; for more information, you can refer to **DWARF**)

# Introduction

#### **Terms**

- tracer tracee
- debugger debuggee

#### strace

instructional userspace utility for Linux. It is used to monitor and tamper with interactions between processes and the Linux kernel, which include system calls, signal deliveries, and changes of process state.

This slide only talks about how strace traces the system calls.



### debugger

A debugger can start some process and debug it, or attach itself to an existing process.

It can **single-step** through the code, set **breakpoints** and run to them, examine variable values and stack traces.

Many debuggers have advanced features such as executing expressions and calling functions in the debugged process's address space, and even **changing the process's code on-the-fly** and watching the effects.

#### ptrace

ptrace is a system call which a program can use to:

- read and write the tracee's memory and registers
- be notified of system events (recognize exec syscall, clone, exit and all the other syscalls)
- control the tracee's execution
- manipulate signal delivery to the traced process

ptrace is used by strace, GDB and other tools as well.

# **Basic ideas**

#### **Fundamentals**

nearly all about 3 interrupts, will talk them latter.

- int 1h, hardware breakpoint, aka. trap
- int 3h, software breakpoint, aka. 0xcc
- int 80h, system call

# **Basic ideas for ptrace**

ptrace works this way:

- the tracee asks the OS kernel to let its parent trace it
- (or vice-versa, the tracer asks the OS kernel to help it trace the tracee)
- then flags are set, which are used for judgement
- when some conditions are satisfied, the tracee is stopped, and the tracer can access the code and data of the tracee

We can take syscall tracing as an example

# Basic ideas - syscall tracing - arises from attach

• The tracer can use the <a href="PTRACE\_ATTACH">PTRACE\_ATTACH</a> flag when calling <a href="ptrace">ptrace</a> and supply the pid of the tracee:

```
ptrace(PTRACE_ATTACH, target_pid, 0, 0);
```

- After that, the tracee is stoped.
- This is followed by another call to ptrace with the PTRACE\_SYSCALL flag:

```
ptrace(PTRACE_SYSCALL, target_pid, 0, 0);
```

 The tracee will run until it enters a system call, at which point it will be stopped by the Linux kernel.

# **Basic ideas - syscall tracing**

- The tracer will appear as if the tracee has been stopped because it received a SIGTRAP signal (via wait()).
- Then the tracer can inspect the args to the system call, print any relevant information.
- Then, the tracer can call <a href="ptrace">ptrace</a> with <a href="ptrace">PTRACE\_SYSCALL</a> again which will resume the tracee, but cause it to be stopped by the kernel when it exits the system call.

The pattern above continues to trace the entry and exit from system calls allowing the tracer to inspect the tracee and print arguments, return values, timing information, and more.

# **Code insight**

Talk is cheap, show me the code.

#### ptrace

A good place to start looking is the code in the kernel for the ptrace system call.

The code samples below will refer to the Linux Kernel 4.6, the arch-specific part will refer to x86.

In kernel/ptrace.c#L1078, ptrace definition:

SYSCALL\_DEFINE4(ptrace, long, request, long, pid, unsigned long, addr, unsigned long, data)

We can start by understanding what PTRACE\_ATTACH does.

#### PTRACE\_ATTACH

kernel/ptrace.c#L1097, statements checking the request parameter for PTRACE\_ATTACH.

ptrace\_attach() is called.

```
if (request == PTRACE_ATTACH || request == PTRACE_SEIZE) {
    ret = ptrace_attach(child, request, addr, data);
    /*
    * Some architectures need to do book-keeping after
    * a ptrace attach.
    */
    if (!ret)
        arch_ptrace_attach(child);
    goto out_put_task_struct;
}
```

ptrace\_attach() is at L295 in the same file, it does a few things early on:

- sets flags which will be stored on a structure in the kernel representing the process that is being attached to, L303~L311
- ensures that the task is not a kernel thread, L316
- ensures that the task is not a thread of the current process, L318
- calls \_\_ptrace\_may\_access() to do some security checks,
   L331

```
Then, flags are set, L349: task->ptrace = flags;, here flags = PT_PTRACED.

And, the process is stopped, L355: send_sig_info(SIGSTOP, SEND_SIG_FORCED, task);.

Then ptrace_attach() returns, and ptrace() finishes(L1105, goto out_put_task_struct;)
```

arch\_ptrace\_attach() is Macro that does nothing.

#### PTRACE\_SYSCALL

In this case, ptrace checks if the process is ready for ptrace operations by calling ptrace\_check\_attach(), L1108:

```
ret = ptrace_check_attach(child, request == PTRACE_KILL || request == PTRACE_INTERRUPT);
```

Then ptrace calls arch\_ptrace() which is a function supplied by the CPU architecture specific code, L1113:

```
ret = arch_ptrace(child, request, addr, data);
```

For us, <a href="mailto:arch\_ptrace">arch\_ptrace</a>() is the x86 <a href="ptrace">ptrace</a> code, which can be found in arch/x86/kernel/ptrace.c#L798.

#### arch\_ptrace

x86's arch\_ptrace() contains a long switch-case list, but it has nothing to do for PTRACE\_SYSCALL.

Thereby the default case is hit, arch/x86/kernel/ptrace.c#L908:

```
ret = ptrace_request(child, request, addr, data);
```

the function hands execution back up to <a href="ptrace\_request()">ptrace\_request()</a>, kernel/ptrace.c#L840.

#### ptrace\_request

There are also a long switch-case list, and at L1019, we can find case PTRACE\_SYSCALL: and then ptrace\_resume() is called at L1021:

#### ptrace\_resume()

It starts by setting the TIF\_SYSCALL\_TRACE flag on the thread info structure for the tracee, L745:

```
if (request == PTRACE_SYSCALL)
    set_tsk_thread_flag(child, TIF_SYSCALL_TRACE);
```

A few possible states are checked (L757~L767, as other functions might call <a href="ptrace\_resume">ptrace\_resume</a>() ) and finally the tracee is woken up and execution resumes until the tracee enters a system call, L786:

```
wake_up_state(child, __TASK_TRACED);
```

### **Entering system calls**

we know that the flag TIF\_SYSCALL\_TRACE is set in tracee, then when is the flag TIF\_SYSCALL\_TRACE checked and acted upon?

whenever a system call is made by a program, there is CPU architecture specific code that is executed on the kernel side prior to the execution of the system call itself.

The code that is executed on x86 CPUs when a system call is made is written in assmebly and can be found in arch/x86/entry/entry\_64.S#L141.

#### \_TIF\_WORK\_SYSCALL\_ENTRY

In the assembly function <a href="ENTRY(entry\_SYSCALL\_64">ENTRY(entry\_SYSCALL\_64)</a>, we can see the flag is checked here, L182:

```
testl $TIF_WORK_SYSCALL_ENTRY|TIF_ALLWORK_MASK,
ASM_THREAD_INFO(TI_flags, %rsp, SIZEOF_PTREGS)
jnz entry_SYSCALL64_slow_path
```

The definition of <a href="https://www.syscall\_entry">LTIF\_WORK\_SYSCALL\_ENTRY</a> is shown here, arch/x86/include/asm/thread\_info.h#L141:

If this flag is set, execution moves to entry\_SYSCALL64\_slow\_path . In entry\_SYSCALL64\_slow\_path , do\_syscall\_64 (arch/x86/entry/common.c, L333) is called.

#### \_TIF\_WORK\_SYSCALL\_ENTRY

In short, on every system call made, the thread info structure for a process has its flags checked for \_TIF\_SYSCALL\_TRACE .

If a flag is set, execution moves to do\_syscall\_64() .

#### do\_syscall\_64()

In do\_syscall\_64() , syscall\_trace\_enter() is called, L341:

```
if (READ_ONCE(ti->flags) & _TIF_WORK_SYSCALL_ENTRY)
    nr = syscall_trace_enter(regs);
```

syscall\_trace\_enter() is defined in the CPU specific ptrace code found in arch/x86/entry/common.c#L192. It calls syscall\_trace\_enter\_phase1() and syscall\_trace\_enter\_phase2() to check if \_TIF\_SYSCALL\_TRACE flag is set. If so, tracehook\_report\_syscall\_entry() is called, L181:

```
if ((ret || test_thread_flag(TIF_SYSCALL_TRACE)) &&
    tracehook_report_syscall_entry(regs))
    ret = -1L;
```

Then tracehook\_report\_syscall\_entry() calls ptrace\_report\_syscall(), include/linux/tracehook.h#L103: return ptrace\_report\_syscall(regs); .

#### ptrace\_report\_syscall()

ptrace\_report\_syscall() would report a syscall: a SIGTRAP is generated when a traced process enters a system call, include/linux/tracehook.h#L66:

```
ptrace_notify(SIGTRAP | ((ptrace & PT_TRACESYSGOOD) ? 0x80 :
0));
```

ptrace\_notify() is implemented in kernel/signal.c#L1913.

ptrace\_notify() calls ptrace\_do\_notify() (L1899,
ptrace\_do\_notify(SIGTRAP, exit\_code, CLD\_TRAPPED); ) which
prepares a signal info structure(L1091~L1097) for
delivery to the process by ptrace\_stop(), L1777:

```
ptrace_stop(exit_code, why, 1, &info);
```

#### **SIGTRAP**

Once the tracee receives a **SIGTRAP**, the tracee is stopped and the tracer is notified that a signal is pending for the process.

The tracer can then examine the state of the tracee and print register values, timestamps, or other information. details will be explains in the debugger part.

This is how strace prints its information to the terminal when you trace a process.

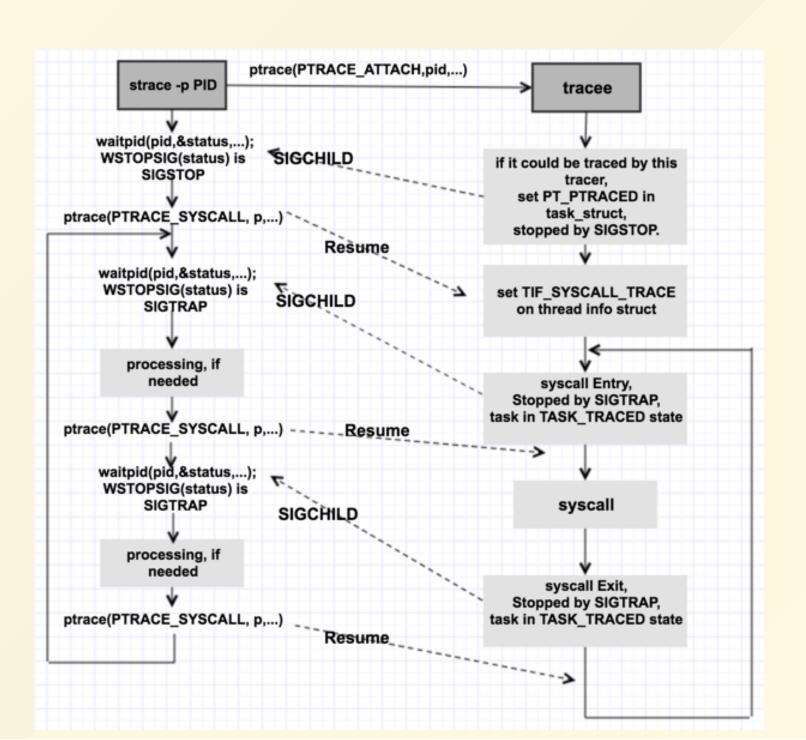
# syscall\_trace\_leave()

A similar code path is executed for the exit path of the system call:

- syscall\_return\_slowpath() is called by do\_syscall\_64()
- syscall\_slow\_exit\_work() is called by syscall\_return\_slowpath()
- this function calls tracehook\_report\_syscall\_exit()
- which also calls ptrace\_report\_syscall(), just like the entry path

And this is how tracing processes are notified when a system call completes.

So they can harvest the **return value**, **timing information**, or **anything else** needed to print useful output for the user.



# More on PT\_PTRACED

#### PTRACE\_TRACEME

another way to set the ptrace flag of is PTRACE\_TRACEME. (PTRACE\_SEIZE is just like PTRACE\_ATTACH, but will not stop the tracee).

kernel/ptrace.c#L1084, statements checking the request parameter for PTRACE\_ATTACH, ptrace\_traceme() is called:

```
if (request == PTRACE_TRACEME) {
    ret = ptrace_traceme();
    if (!ret)
        arch_ptrace_attach(current);
    goto out;
}
```

#### PTRACE\_TRACEME

ptrace\_traceme() is at L409, it does a few things:

- check whether this current process is already being traced, L415
- determine whether another process may trace the current, L416
- Check PF\_EXITING to ensure real\_parent has not passed exit\_ptrace(), L422

Then PT\_PTRACED is set, current->ptrace = PT\_PTRACED;, at L423.

And the function returns, ptrace() finishes (L1088, goto out; )

#### PTRACE\_TRACEME

you can see the current process is not stopped, but:

- all subsequent calls to exec() by this process will cause a SIGTRAP to be sent to it
- gives the parent a chance to gain control before the new program begins execution.

#### Note:

- A process probably shouldn't make this request if its parent isn't expecting to trace it.
- child should call ptrace\_traceme(), then exec() the specified program.

#### execve

we can checked this in sys\_execve(), fs/exec.c#L1806:

sys\_execve() calls do\_execve(), then do\_execve() calls do\_execveat\_common(), L1724.

In the do\_execveat\_commom(), exec\_binprm() is called at L1686.

Finally in exec\_binprm(), can find ptrace\_event(PTRACE\_EVENT\_EXEC, old\_vpid); is executed at L1570.

#### ptrace\_event

located at include/linux/ptrace.h#L143. with the event PTRACE\_EVENT\_EXEC and PT\_PTRACED flag, the tracee would receive a SIGTRAP signal, L148~151:

#### PT\_PTRACED

Indicates that this process is to be traced by its parent. Any signal (except SIGKILL) delivered to this process will cause it to stop and its parent to be notified via wait().

# About debuggers

ptrace allow a (debugger) process to access low-level information about another process (the debuggee). In particular, the debugger can:

- read and write the debuggee's memory:
   PTRACE\_PEEKTEXT, PTRACE\_PEEKDATA, PTRACE\_PEEKUSER,
   PTRACE\_POKETEXT, PTRACE\_POKEDATA, PTRACE\_POKEUSER
- read and write the debuggee's CPU registers:
   PTRACE\_GETREGSET , PTRACE\_SETREGS
- be notified of system events: PTRACE\_O\_TRACEEXEC,
  PTRACE\_O\_TRACECLONE, PTRACE\_O\_EXITKILL, PTRACE\_SYSCALL
- control its execution: PTRACE\_SINGLESTEP , PTRACE\_KILL , PTRACE\_CONT
- alter its signal handling: PTRACE\_GETSIGINFO, PTRACE\_SETSIGINFO

things about PTRACE\_SYSCALL is clearly explained above. the rest will talk about other flags briefly and how debuggers work.

# More on ptrace

- how to access CPU registers
  - easy with copy\_regset\_to/from\_user . both
     PTRACE\_GETREGS and PTRACE\_SETREGS are handled in
     arch\_ptrace() at arch/x86/kernel/ptrace.c#L840
     and #L847 with copy\_regset\_to\_user() and
     copy\_regset\_from\_user() respectively.
- how to access debuggee's memory
  - copy\_to/from\_user may work well. Actually with respect to PTRACE\_PEEKDATA, PTRACE\_PEEKTEXT and PTRACE\_POKEDATA, PTRACE\_POKETEXT (Linux does not have separate text and data address spaces, so these two requests are currently equivalent)
  - o ptrace\_requests just calls generic\_ptrace\_peekdata() and generic\_ptrace\_pokedata() respectively, which using access\_process\_vm() to read/write the destination address. access\_process\_vm() calls \_\_access\_remote\_vm(), and in \_\_access\_remote\_vm(), copy\_to\_user\_page or copy\_from\_user\_page is called to write or read the desired data, mm/memory.c#L3768~L3774.

- how to change the signal handling
  - PTRACE\_GETSIGINFO and PTRACE\_SETSIGINFO are handled in ptrace\_request() by ptrace\_getsiginfo() and ptrace\_setsiginfo() respectively. ptrace\_getsiginfo() just copies the signal from child->last\_siginfo and ptrace\_setsiginfo() just updates field child->last\_siginfo , kernel/ptrace.c#L617 and L633
- how to know system events
  - ptrace is hooked (function ptrace\_event(), include/linux/ptrace.h#L143) in many scheduling operations, so that it can send a SIGTRAP signal to the debugger if requested (PTRACE\_O\_TRACEEXEC option and its family). just like the PTRACE\_EVENT\_EXEC mentioned above.

### single step debugging

Use <a href="https://process.com/pt/pt/">ptrace(PTRACE\_SINGLESTEP, pid, 0, 0)</a> for single step debugging. What this does is tell the OS - restart the child process, but stop it after it executes the next instruction.

actually this feature is implemented with the Trap Flag and int 1h in x86.

like PTRACE\_SYSCALL, PTRACE\_SINGLESTEP is also handled in ptrace\_resume(), kernel/ptrace.c#L840:

```
return ptrace_resume(child, request, data);
```

then ptrace\_resume() calls user\_enable\_single\_step(child) at L764.

## single step debugging

```
Inside user_enable_single_step(), enable_step(child, 0); is
called at kernel/step.c#L210. Then enable_step() calls
enable_single_step(child) at L204, in which some
important flags are set: L131
set_tsk_thread_flag(child, TIF_SINGLESTEP); , L136:
regs->flags |= X86_EFLAGS_TF; , L159:
set_tsk_thread_flag(child, TIF_FORCED_TF); .
```

After user\_enable\_single\_step() is called, ptrace\_resume() executes wake\_up\_state(child, \_\_TASK\_TRACED); at kernel/ptrace.c#L786, the tracee will be woken up by the scheduler.

Once the tracee is scheduled to run, the regs is restored, causing TF set, the CPU enters single-step mode. What means the system will execute the next instruction and do a int 1h interrupt.

### int 1h

in Linux, int 1h handler is defined at arch/x86/kernel/traps.c#L607:

```
dotraplinkage void do_debug(struct pt_regs *regs, long
error_code)
```

at L708, SIGTRAP is sent to the current task:

```
send_sigtrap(tsk, regs, error_code, si_code);
```

at this moment, tracee is stopped, and tracer knows tracee is stopped.

## handle\_signal

Also in the handle\_signal() at arch/x86/kernel/signal.c#L736, the function tests the TIF\_SINGLESTEP and disables the single step, L736~738

```
stepping = test_thread_flag(TIF_SINGLESTEP);
if (stepping)
    user_disable_single_step(current);
```

user\_disable\_single\_step() just clear task flag clear\_tsk\_thread\_flag(child, TIF\_SINGLESTEP); and clear TRAP FLAG ( TIF\_FORCED\_TF is used here ), kernel/step.c#L232~233:

```
if (test_and_clear_tsk_thread_flag(child, TIF_FORCED_TF))
  task_pt_regs(child)->flags &= ~X86_EFLAGS_TF;
```

All done!

### breakpoints

breakpoints are implemented on the CPU by a special trap called int 3h.

The int 3h instruction generates a special one byte opcode (0xcc) that is intended for calling the debug exception handler.

So how to use it?

- 1. debugger reads (ptrace peek) the instruction stored at this address, and backups it.
- 2. write <code>0xcc</code> at this location.
- 3. when debuggee reaches this location, int 3h is triggered.
- 4. this interruption causes a **SIGTRAP** sent to debuggee, then the debugger gets noticed.
- 5. the debugger gets the signal, and checks the debuggee's instruction pointer ( check whether the IP address is in its breakpoint list ).
- 6. now that the debuggee is stopped at the breakpoint, the debugger can let its user do what ever s/he wants.
- 7. to continue, the debugger needs to:
  - 1. write the correct instruction back.
  - 2. unwind the IP back by 1.
  - 3. single-step it.
  - 4. write the <code>0xcc</code> back.
  - 5. let the execution flow normally.

#### as for ptrace

- use PTRACE\_PEEKTEXT and PTRACE\_POKETEXT to alter the instruction
- use PTRACE\_GETREGS and PTRACE\_SETREGS to access and update IP register
- use PTRACE\_CONT or PTRACE\_SINGLESTEP to resume the tracee.

PTRACE\_CONT is handle by ptrace\_resume() like PTRACE\_SINGLESTEP, in this case:

- TIF\_SYSCALL\_TRACE is cleared, L748: clear\_tsk\_thread\_flag(child, TIF\_SYSCALL\_TRACE);
- single\_step is disabled, L766: user\_disable\_single\_step(child);
- then the tracee will be woken up by wake\_up\_state(child, \_\_TASK\_TRACED);, L786.

#### int 3h

the int 3h handler is defined at arch/x86/kernel/traps.c#L468:

```
dotraplinkage void notrace do_int3(struct pt_regs *regs, long
error_code)
```

we can find the SIGTRAP is sent at L506:

```
do_trap(X86_TRAP_BP, SIGTRAP, "int3", regs, error_code,
NULL);
```

### conditional breakpoint

conditional breakpoints are **normal breakpoints**, except that:

- the debugger checks the conditions before giving the control to the user.
- if the condition is not matched, the execution is silently continued.

### watchpoints

watchpoints are implemented (if available) with the help of the processor:

- write in registers which addresses should be monitored
- it will raise an exception when the memory is read or written

If this support is not available, or if you request more watchpoints than the processor supports, then the debugger falls back to "hand-made" watchpoints:

- single-step the program
- check if the current operation touches a watchpointed address

# **Demos**

strace.c singlestep.c breakpoint.c

### References

- http://man7.org/linux/man-pages/man1/strace.1.html
- http://man7.org/linux/man-pages/man2/ptrace.2.html
- https://github.com/torvalds/linux/tree/v4.6
- https://blog.nelhage.com/2010/08/write-yourself-anstrace-in-70-lines-of-code/
- https://blog.0x972.info/?d=2014/11/13/10/40/50how-does-a-debugger-work

Thanks!