# runtime: high GC latency

⊙ Open

**aclements** commented

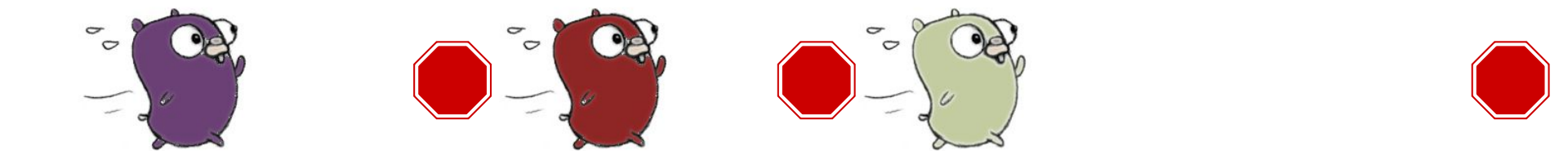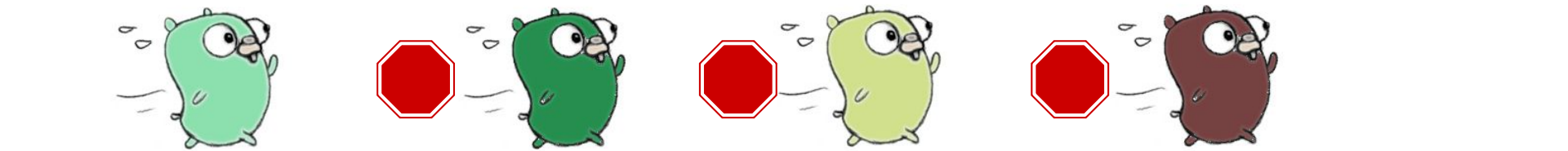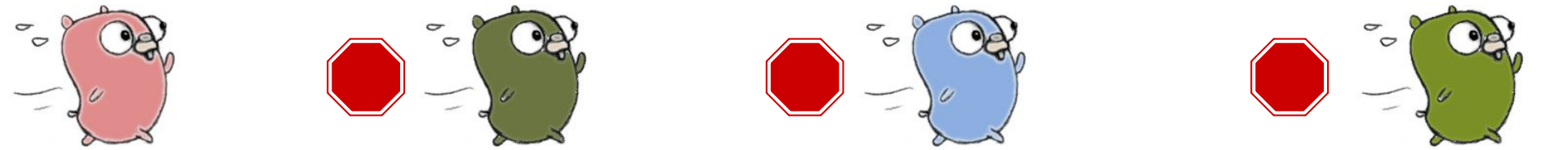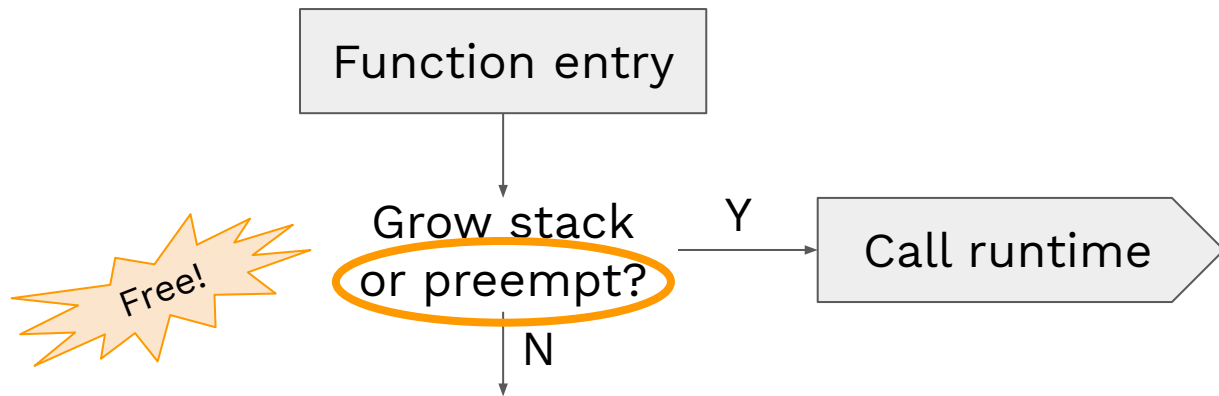Maybe you have a tight loop? ¯\\_(ツ)_/¯

GO

## Easy: goroutine voluntarily gives up control

```go
// Wait for result.
x := <-ch
```

## Not so easy: need to *preempt* a running goroutine

```go
// Capitalize 100GB of very important data.
strings.ToUpper(bigData)
```

Go 1.0 ◯ Voluntary preemption only

Go 1.2 ◯ **Function call preemption**

Function entry

Grow stack
or preempt? —Y→ Call runtime

*Free!*

N

## Call-free loops delay preemption

```go
// Darken image.
for y := range image {
    for x := range image[y] {
        image[y][x] *= 0.5
    }
}
```

47x throughput drop!

Large base64 decode

48 cores

~5ms

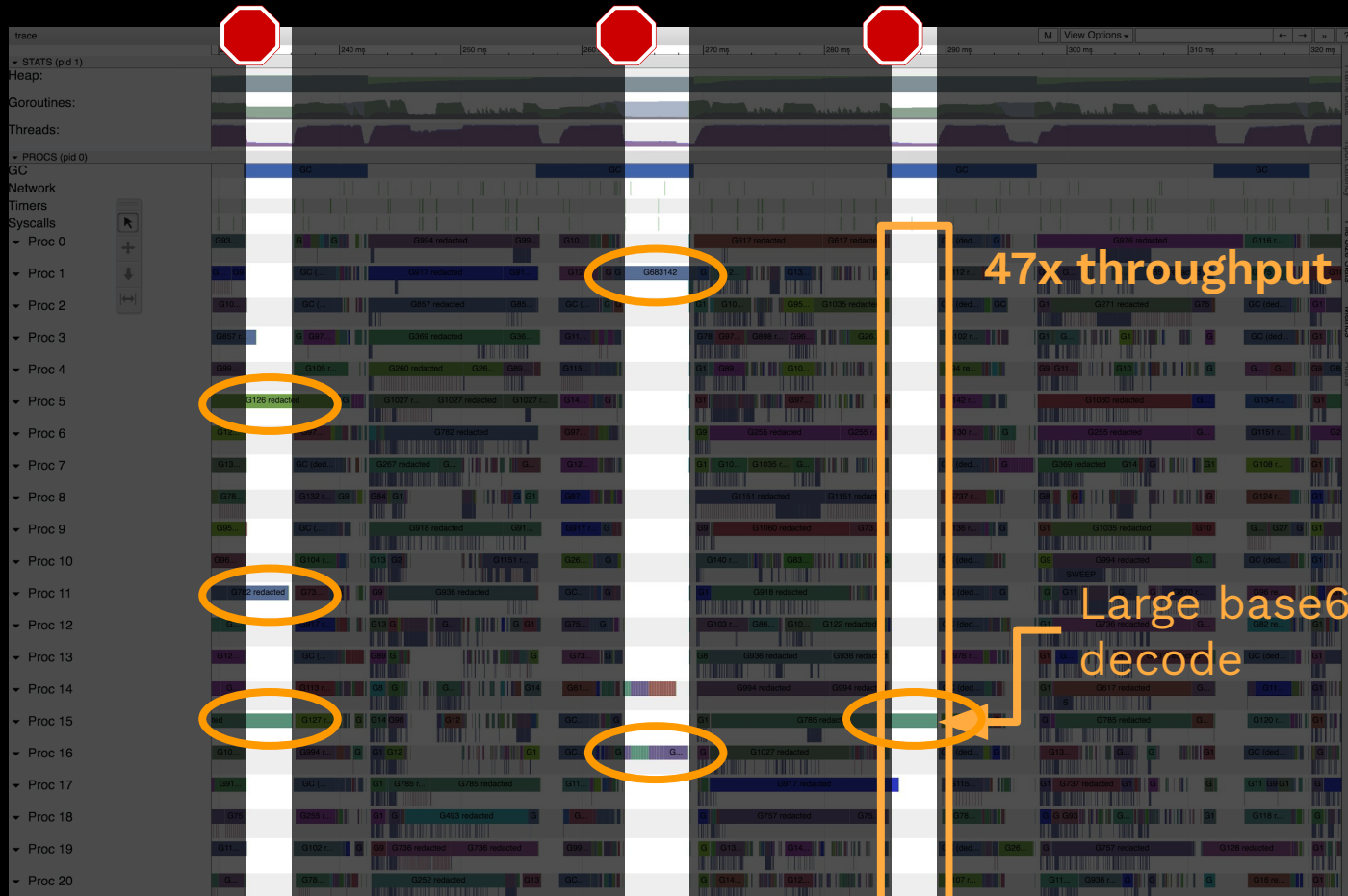Source: https://golang.org/issue/17831

## Call-free loops delay preemption

```go
// Darken image.
for y := range image {
    for x := range image[y] {
        image[y][x] *= 0.5
    }
}
```

## Unbounded call-free loops could deadlock the scheduler

```go
// Wait for completion.
for atomic.LoadUint32(&status) == 0 { }
```

```go
func (c *queue) wait() {
    for {
        if atomic.AddInt32(&c.RWMutex.readerCount, 1) < 0 {
            c.RWMutex.rLockSlow()
        }
        flushed := c.flushed
        if r := atomic.AddInt32(&c.RWMutex.readerCount, -1); r < 0 {
            c.RWMutex.rUnlockSlow(r)
        }
        if flushed {
            return
        }
    }
}
```

Uncontended lock + unlucky goroutine schedule ⇒ deadlock

**Sync package change introduced a deadlock into this application code**
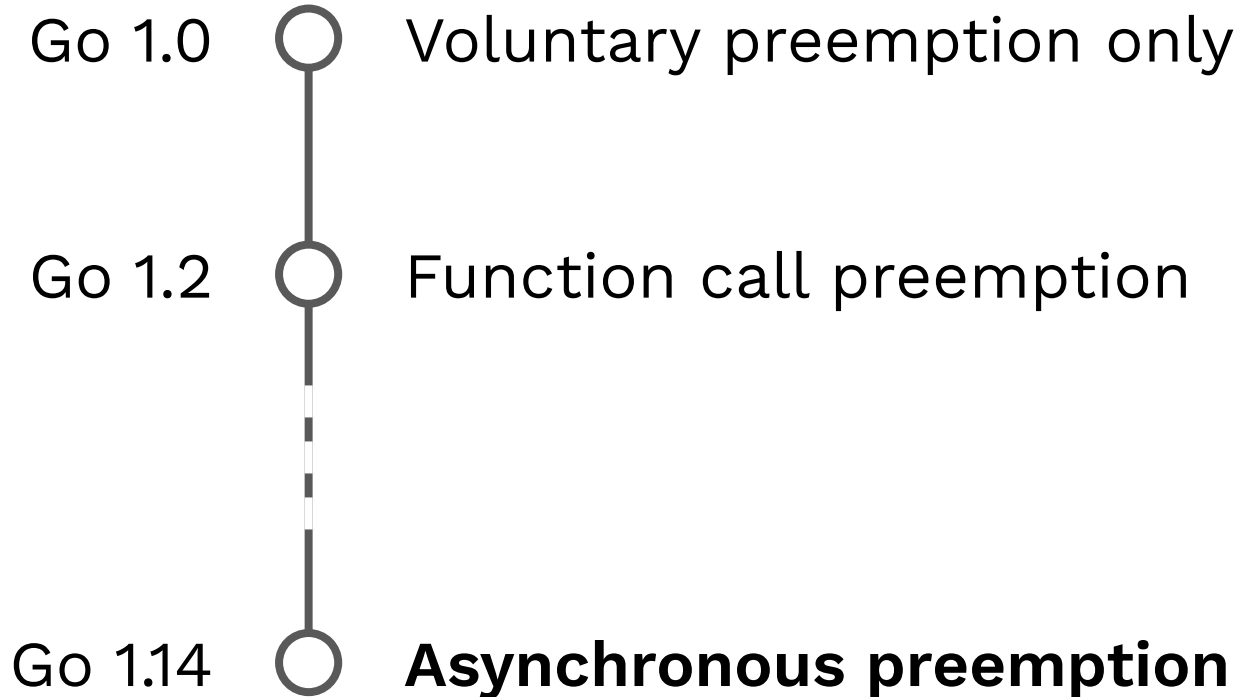
## Conventional wisdom

```go
// Darken image.
for y := range image {
    for x := range image[y] {
        image[y][x] *= 0.5
    }
    runtime.Gosched()    // Fixed! YOLO
}
```

Function call preemption problems:

- Loops cause scheduling latency

- Latency reduces application throughput

- Loops can cause scheduler deadlocks

Difficult to diagnose

Difficult to fix

## We need **loop preemption**.

Go 1.0 ⭘ Voluntary preemption only

Go 1.2 ⭘ Function call preemption

Go 1.14 ⭘ **Asynchronous preemption**

# Asynchronous preemption

| ~25μs | 0 | +0.3% |
|:---:|:---:|:---:|
| typical preemption bound | performance overhead | binary size |

~~Gosched~~

# False starts

GO

8%

slowdown

**Adding instructions to loops isn't viable.**

Break in to the goroutine

Single step to back-edge

Redirect to preemption path

**Keep it simple. Consider user experience.**

CPU state
PC 0x000000000046df61
R1 0x000000c00004ef00
R2 0x00000000000000ca
…

Which are pointers?

Zero overhead, but doesn't work with GC.

# Interlude: How does the GC normally find pointers?

```
...
MOVQ      AX, 8(SP)
LEAQ      -1(CX), AX
MOVQ      AX, 16(SP)
MOVQ      BX, 24(SP)
CALL      flag.(*FlagSet).Parse(SB)
XCHGL     AX, AX
MOVQ      log.std(SB), AX
MOVQ      AX, (SP)
MOVQ      $0, 8(SP)
CALL      log.(*Logger).SetFlags(SB)
MOVQ      flag.CommandLine(SB), AX
...
```

GC stack map

```
SP+ 0 000000c0000ac180
SP+ 8 000000c0000b81b0
SP+16 0000000000000000
SP+24 0000000000000000
SP+32 0000000000203000
SP+40 000000c000014c40 ...
SP+48 0000000000a3b4a8
SP+56 000000c00006ce58
SP+64 0000000000000120
SP+72 000000000097c140
SP+80 000000c000004f2a ...
SP+88 0000000000ea5138
...
```

0   8  16  24  32  40  48  56 ...
Stack frame offset

GC stack map

```
...
MOVQ    AX, 8(SP)
LEAQ    -1(CX), AX
MOVQ    AX, 16(SP)
MOVQ    BX, 24(SP)
CALL    flag.(*FlagSet).Parse(SB)
XCHGL   AX, AX
MOVQ    log.std(SB), AX
MOVQ    AX, (SP)
MOVQ    $0, 8(SP)
CALL    log.(*Logger).SetFlags(SB)
MOVQ    flag.CommandLine(SB), AX
...
```

| 0 | 8 | 16 | 24 | 32 | 40 | 48 | 56 | ... |
|---|---|----|----|----|----|----|----|-----|
| - | - | -  | P  | -  | P  | -  | P  | ... |
| - | - | -  | P  | -  | P  | -  | P  | ... |
| - | - | -  | P  | -  | P  | -  | P  | ... |
| - | - | P  | P  | -  | P  | -  | P  | ... |
| - | - | P  | P  | -  | P  | -  | P  | ... |
| - | - | -  | -  | -  | P  | -  | P  | ... |
| - | - | -  | -  | -  | P  | -  | P  | ... |
| P | - | -  | -  | -  | P  | -  | P  | ... |
| P | - | -  | -  | -  | P  | -  | P  | ... |
| P | - | -  | -  | -  | P  | -  | P  | ... |
| - | - | -  | -  | -  | P  | -  | P  | ... |

Stack frame offset

GC stack map

+10%

binary size

+ 🫨

testing

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| − | − | − | P | − | P | − | P | · · · |
| − | − | − | P | − | P | − | P | · · · |
| − | − | − | P | − | P | − | P | · · · |
| − | − | P | P | − | P | − | P | · · · |
| − | − | P | P | − | P | − | P | · · · |
| − | − | − | − | − | P | − | P | · · · |
| − | − | − | − | − | P | − | P | · · · |
| P | − | − | − | − | P | − | P | · · · |
| P | − | − | − | − | P | − | P | · · · |
| P | − | − | − | − | P | − | P | · · · |
| − | − | − | − | − | P | − | P | · · · |

# Signals are a good idea, but metadata isn't.

Keep it simple and consider end-to-end user experience.
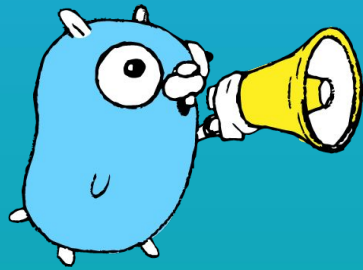
Explicit loop preemption checks are too costly for Go.

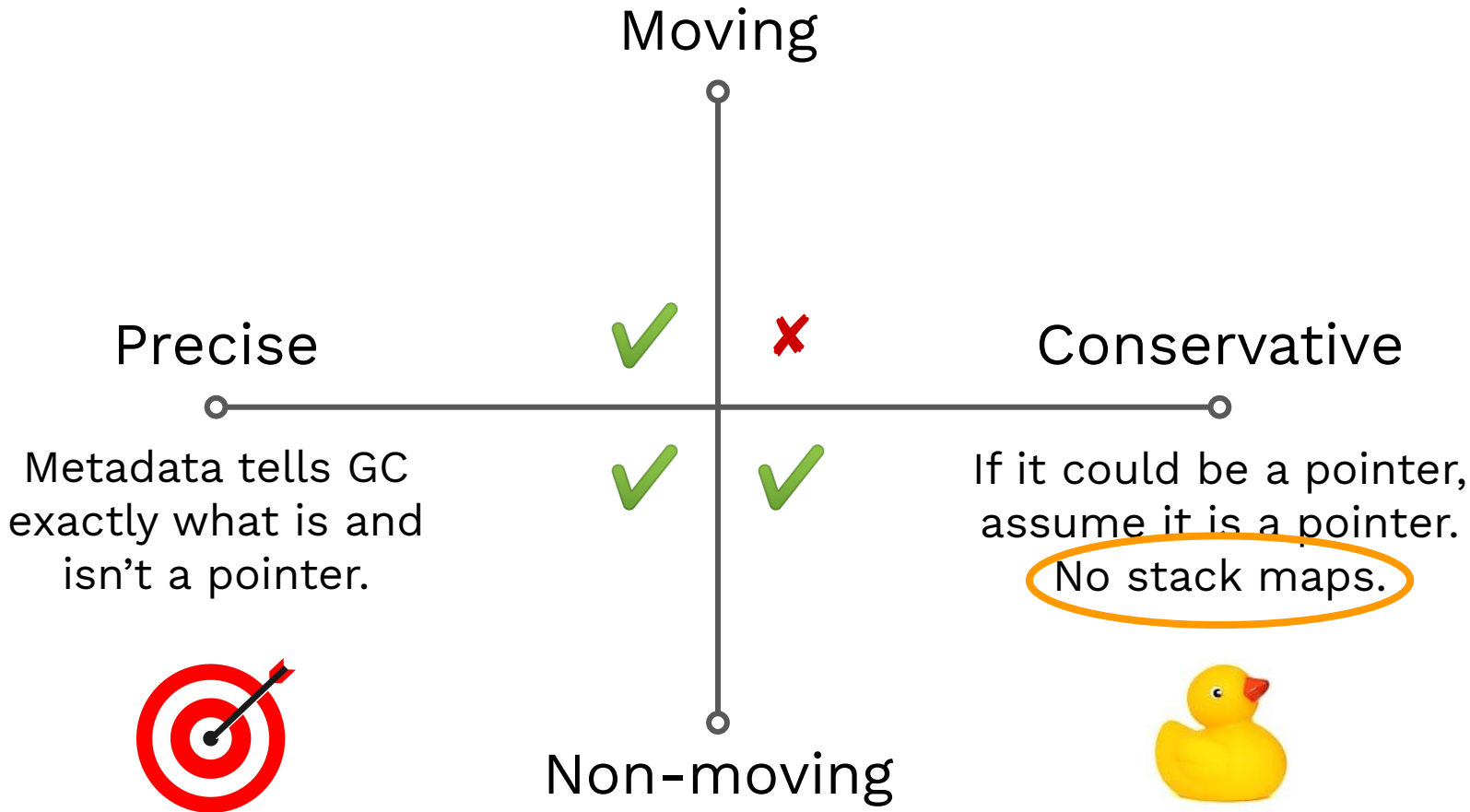Signals have zero ongoing cost, but are only half the answer.

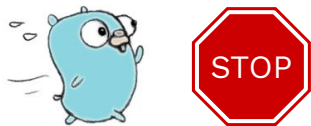Stack maps everywhere bloats binaries and doesn't "keep it simple."

GO

# Asynchronous preemption in Go 1.14

Signal-based preemption with conservative innermost frame scanning

Moving

Precise ✔ ✘ Conservative

Metadata tells GC
exactly what is and
isn't a pointer.

If it could be a pointer,
assume it is a pointer.

No stack maps.

✔ ✔

Non-moving

Send a signal to the goroutine
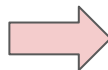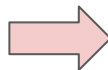


```
MOVQ     BX, 24(SP)
CALL     flag.(*FlagSet).Parse(SB)
XCHGL    AX, AX
MOVQ     log.std(SB), AX
MOVQ     AX, (SP)
MOVQ     $0, 8(SP)
CALL     log.(*Logger).SetFlags(SB)
MOVQ     flag.CommandLine(SB), AX
...
```

Examine where it stopped



Conservatively scan function



```
CALL     main.setup(SB)
```

```
CALL     main.main(SB)
```

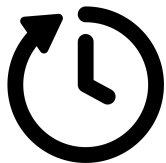Precisely scan rest of stack



```
CALL     runtime.main(SB)
```

GO

# Limitation: Unsafe points

```
// Increment pointer.
y := unsafe.Pointer(uintptr(unsafe.Pointer(x)) + 1)
```

```
var big [1<<30]byte  ⟹  runtime.memclr(...)
```

```
z := big             ⟹  runtime.memmove(...)
```
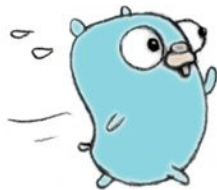
Wait and retry signal

# Implementation

# Send a signal to the goroutine



|  |  |  | Spurious okay? | Coalesce okay? | GDB pass? |
|---|---|---|---|---|---|
| 21 | SIGTTIN | Terminal input in background | | | |
| 22 | SIGTTOU | Terminal output in background | | | |
| 23 | SIGURG | Urgent condition on socket | ✓ | ✓ | ✓ |
| 24 | SIGXCPU | CPU time limit exceeded | | | |
| 25 | SIGXFSZ | File size limit exceeded | | | |

Examine where it stopped



CALL runtime.asyncPreempt

runtime.asyncPreempt

Save CPU state,
call scheduler

runtime.sighandler

If safe to preempt,
simulate call

# Insight: get out of the signal handler ASAP
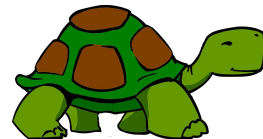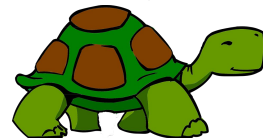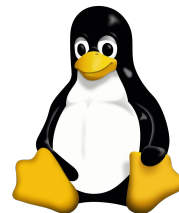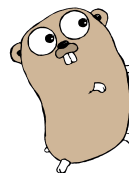
GO

# Debug the Linux kernel?!

New thread
+ Alternate signal stack
+ First signal
+ Lots of CPU activity
+ Linux 5.2 compiled with GCC 9
_____

   SSE register corruption
⇒ Memory corruption

Bad preemption
point in the kernel!

**Moving GC** requires

⇓

Precise stack scanning ⬦ Preempt anywhere

Metadata at preempt points Metadata-less conservative GC

Back-edge preemption / Zero-cost signal preemption
Forward simulation

⇓

Requires **non-moving GC**