# The Enhancement of Kernel Probing
## - Kprobes Jump Optimization

Masami Hiramatsu

masami.hiramatsu.pt@hitachi.com

Hitachi Systems Development Laboratory

Linux Technology Center
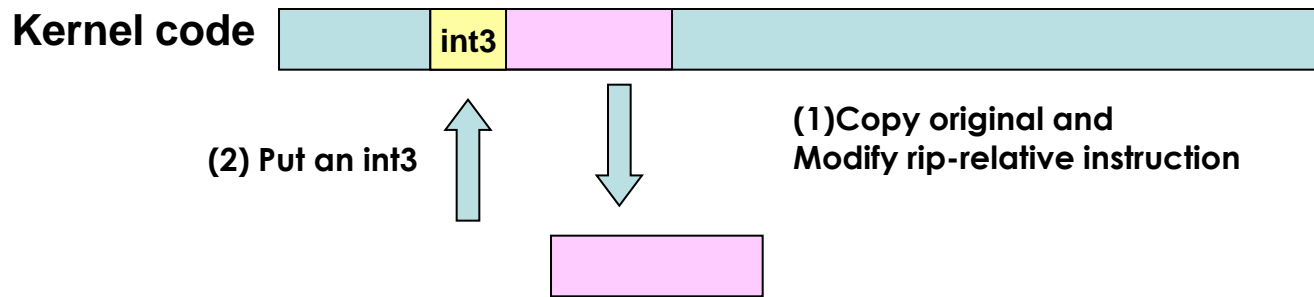
**HITACHI**
Inspire the Next

- Kprobes – Why it is useful
- Kprobes – How it works
- Performance Enhancing Ideas
  - Booster
  - Jump Optimization
- Technical Issues
  - Interrupts
  - Instruction Boundary
    - X86 Instruction Decoder
  - Jumps
  - Cross Code Modifying
- Implementation
  - Transparency of API/ABI
  - Greedy Optimization
  - Reserve Text
- Results
  - Kprobes
  - Kretprobes
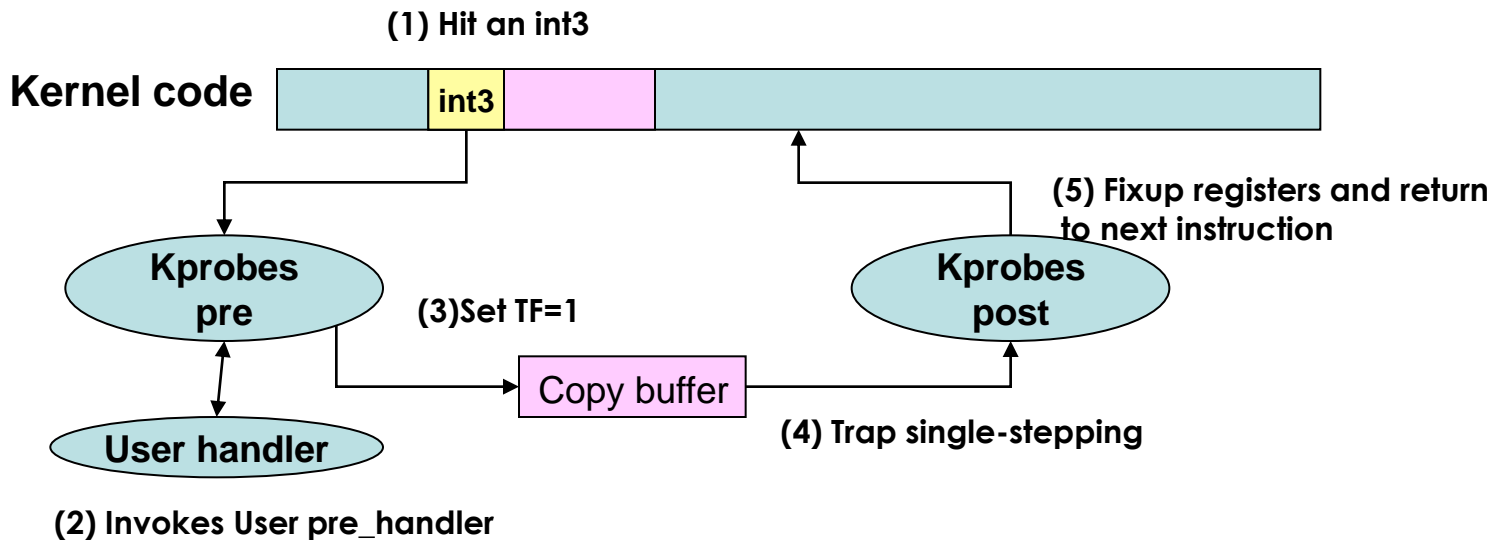  - Results on KVM
- Conclusion

- Kprobes is a dynamic software breakpoint function in the kernel
  - This allows you to add breakpoints inside kernel
    - User can check the kernel internal state almost anywhere
    - This allows user to tweak kernel internal state too (e.g. fault injection, and dynamic patching)
  - Dynamically add and remove the breakpoints.
  - Manage the breakpoint handlers
    - Handling breakpoint exception and call handlers
    - Aggregate probes on the same address
    - Disable probes when a target module is gone
    - Etc.

Linux
Technology
Center

- Kprobes uses a breakpoint and a single-step

## Preparing

**Kernel code**

| | | int3 | | |
|---|---|---|---|---|

**(2) Put an int3**

**(1)Copy original and
Modify rip-relative instruction**

## Running

**(1) Hit an int3**

**Kernel code**

| | | int3 | | |
|---|---|---|---|---|

**(5) Fixup registers and return
to next instruction**

**Kprobes
pre**

**Kprobes
post**

**(3)Set TF=1**

**User handler**

Copy buffer

**(4) Trap single-stepping**

**(2) Invokes User pre_handler**

*4*

# • Kprobes uses 2 exceptions

- – Software Breakpoint exception
- – Single-step trap

**Kernel code** | int3

**Breakpoint exception**

**Kprobes pre**

**User handler**

Copy buffer

**Kprobes post**

**Single-step Trap**

## Normal kprobe consumes >1500 cycles/probe
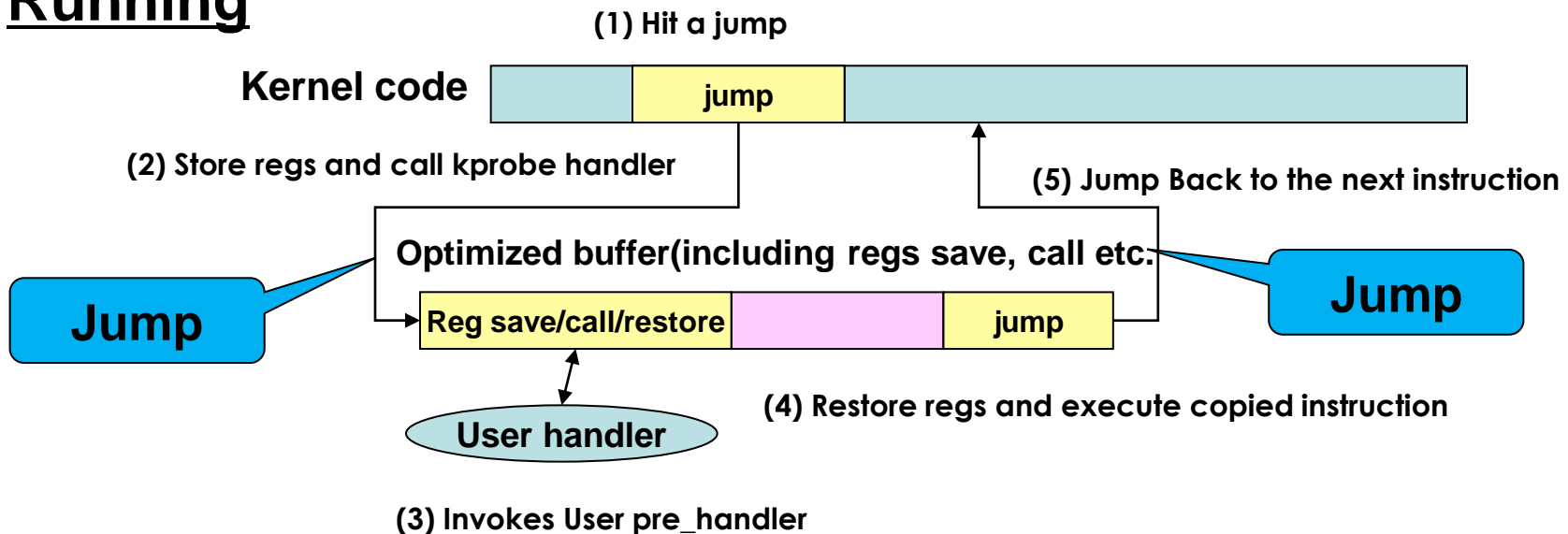
*5*

# Linux Technology Center

- ## Kprobes Booster skips Trap exception
  - – Add a jump which jumps back to next instruction
  - – Execute copied instruction and the jump
  - – Some instructions can't be boosted
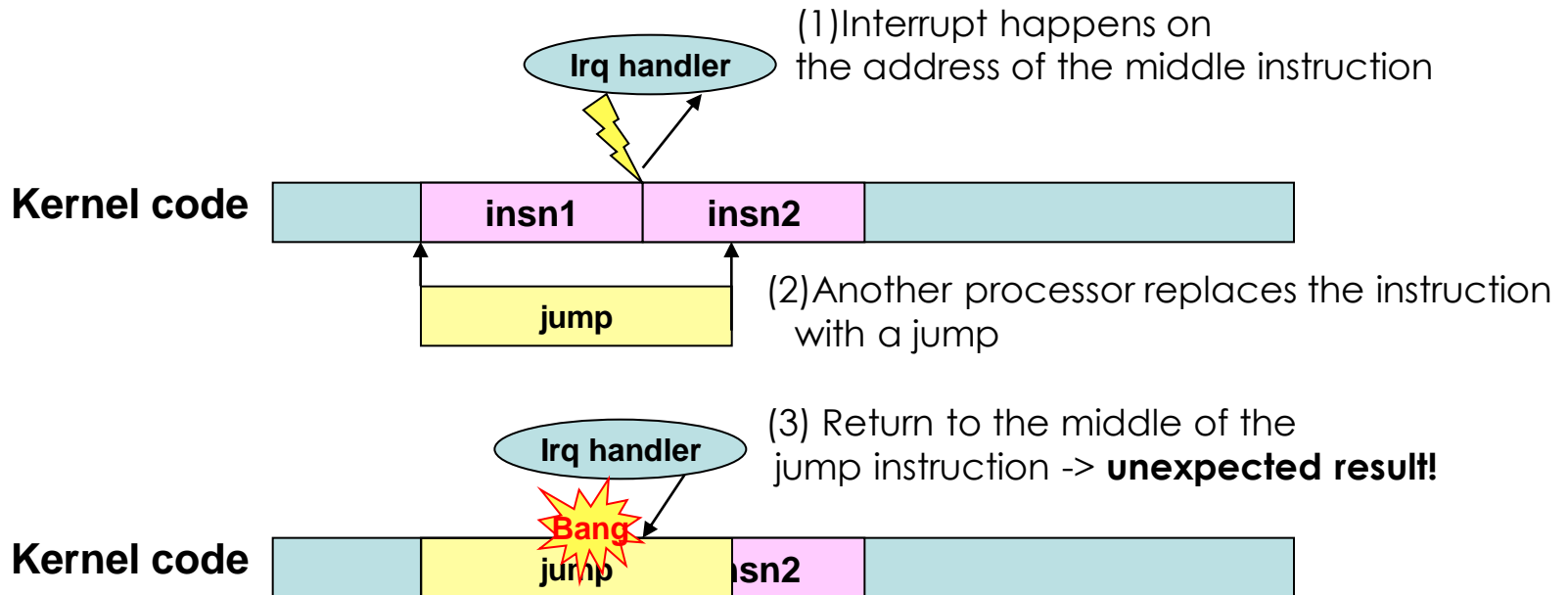    - • Call, near jump, etc

## Running

**(1) Hit an int3**

**Kernel code**    `int3`

**Breakpoint exception**

**(4) Jump Back to the next instruction**

**Kprobes pre**

**(3) Return to the copied instruction**

`jump`

**Jump**

**User handler**

**(2) Invokes User pre_handler**

# • Kprobes Jump Optimization

- – Skips software breakpoint too
  - • No exception: Reduce the overhead drastically
- – It's not easy – of course.
  - • This will replace **several instructions** with one jump
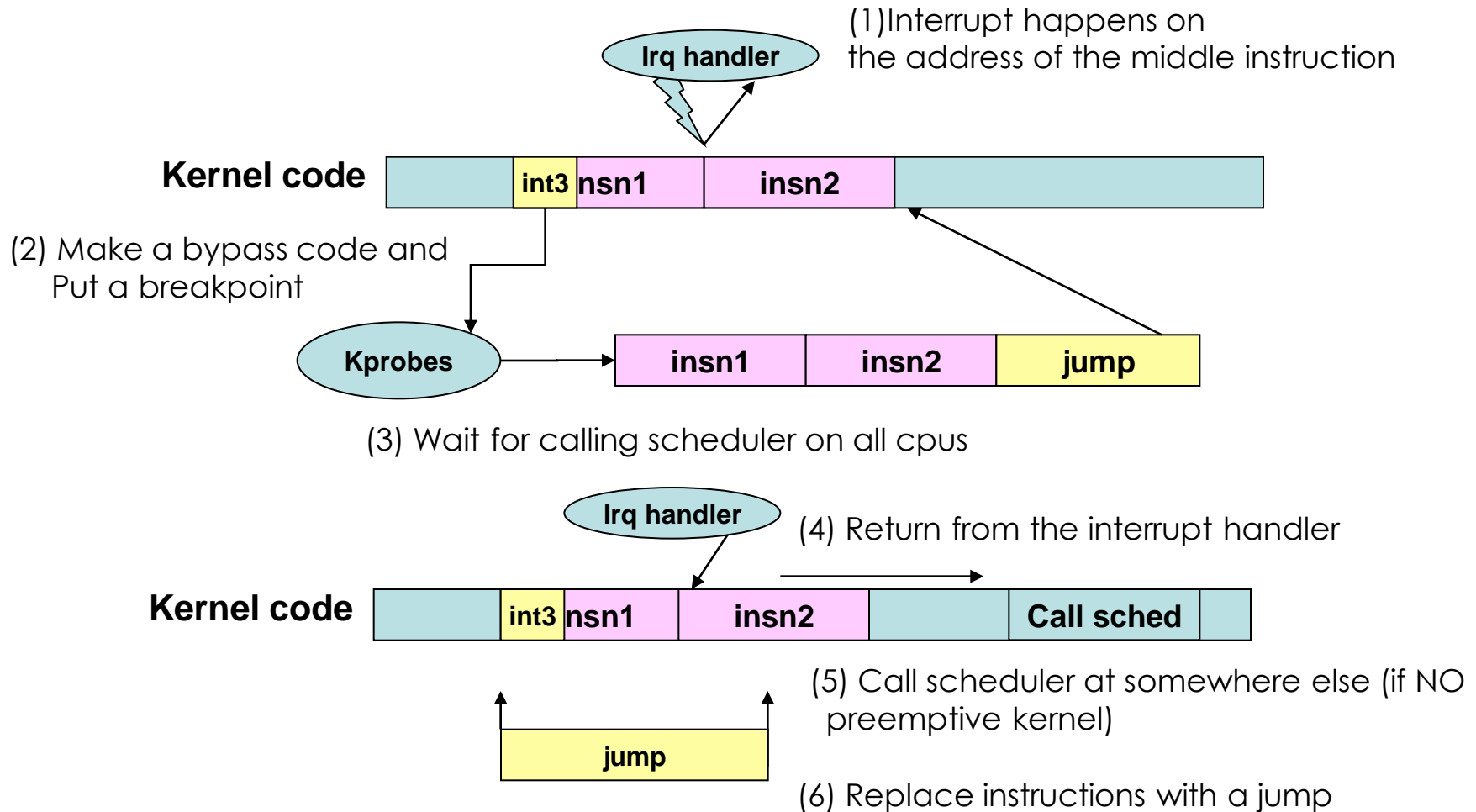    - – Kprobes just replace one instruction.

## Running

**(1) Hit a jump**

**Kernel code** | jump |

**(2) Store regs and call kprobe handler**

**(5) Jump Back to the next instruction**

**Optimized buffer(including regs save, call etc.**

**Jump**

| Reg save/call/restore | | jump |

**Jump**

**User handler**

**(4) Restore regs and execute copied instruction**

**(3) Invokes User pre_handler**

- Interrupts can happen on other processors

(1)Interrupt happens on the address of the middle instruction

**Irq handler**

**Kernel code**

| insn1 | insn2 |
|-------|-------|

**jump**

(2)Another processor replaces the instruction with a jump

(3) Return to the middle of the jump instruction -> **unexpected result!**

**Irq handler**

**Bang**

**Kernel code**

| jump | sn2 |
|------|-----|

Make sure no process is interrupted on
the address where will be replaced by the jump

# • Make a bypass and wait for scheduler

(1) Interrupt happens on
the address of the middle instruction

**Irq handler**

**Kernel code**  | int3 | nsn1 | insn2 |

(2) Make a bypass code and
Put a breakpoint

**Kprobes** → | insn1 | insn2 | jump |

(3) Wait for calling scheduler on all cpus

**Irq handler**

(4) Return from the interrupt handler

**Kernel code**  | int3 | nsn1 | insn2 | Call sched |

(5) Call scheduler at somewhere else (if NO
preemptive kernel)

| jump |

(6) Replace instructions with a jump

- # x86 is a CISC processor
  - Instructions vary in length
  - How many bytes do we need to copy?

**Kernel code** | | insn1 | insn2 | |
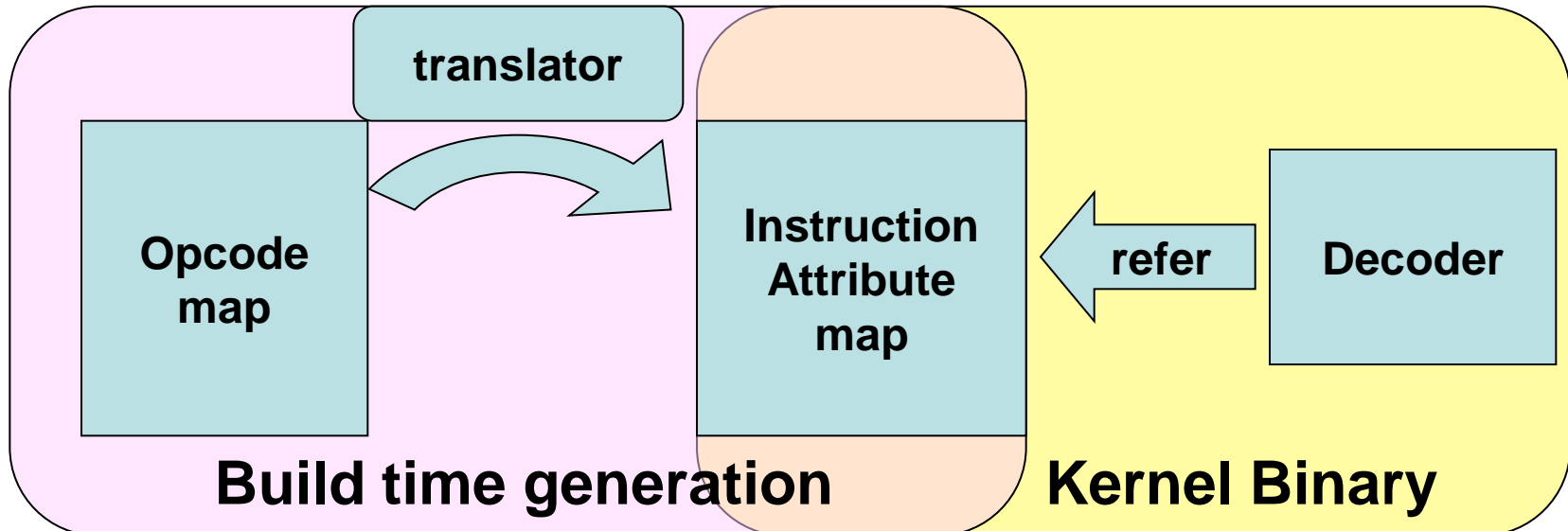
Copy

←————————→ How many bytes?

- # Check non-relocatable instructions
  - Some IP-related instructions can't execute directly on copy buffer (Call, relative-jump, etc)
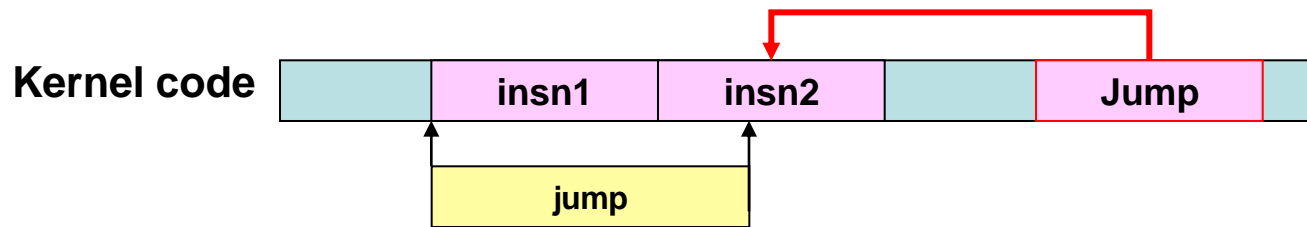  - How can we find those instructions if it is in the middle?

**Kernel code** | | insn1 | IP-related | |

From where does it start?

## We need something to decode instructions!

# Linux Technology Center

- ## Introduce in-kernel x86 instruction decoder
  - ### Simple instruction decoder
    - Just ~350 logical lines including AVX(Intel® Advanced Vector Extensions) decoding support
  - ### Generic & easy maintain
    - Based on x86 opcode map (in Intel's software developers manual)
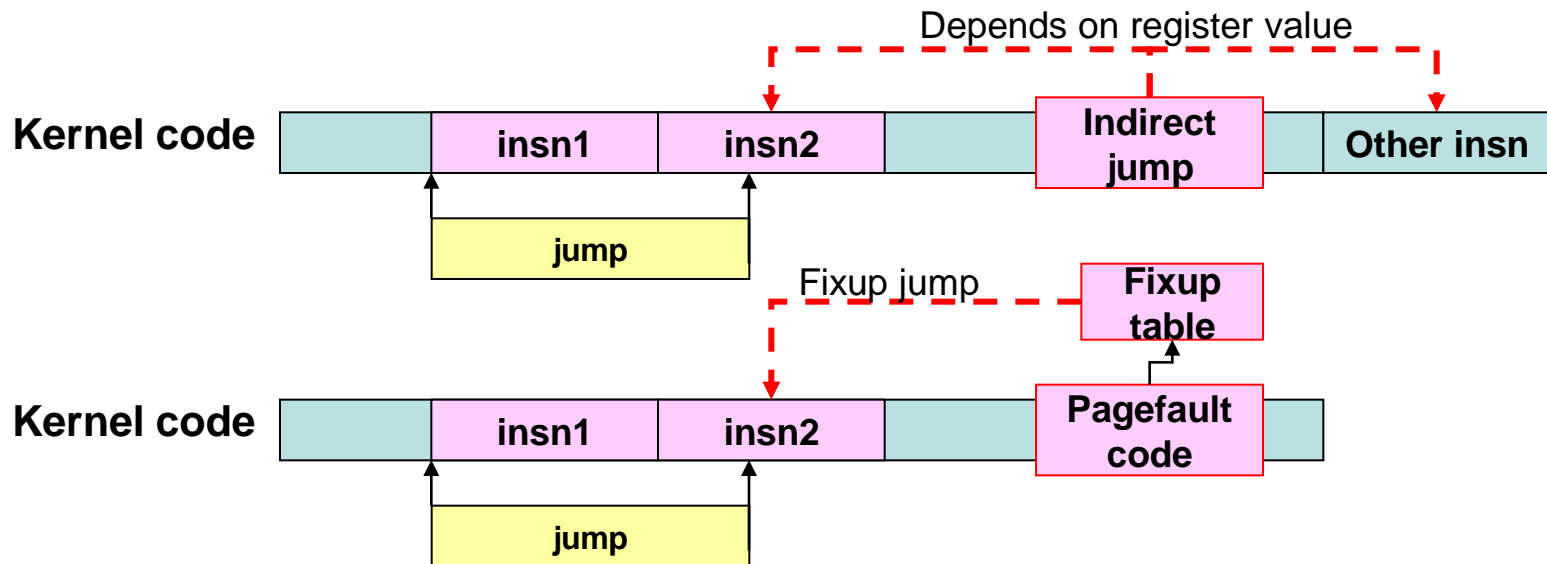    - Generate instruction attribute map from the opcode map when compiling kernel

**translator**

**Opcode map**

**Instruction Attribute map**

**refer**

**Decoder**

**Build time generation**

**Kernel Binary**

*11*

- ## x86 instruction decoder has two parts
  - insn
    - Data structure represents an instruction
    - insn_init() and insn_get_XXX()
    - users usually use this part
  - inat
    - Instruction attribute maps for decoding
      - Each opcode has attributes

```
struct insn;
int x86_64 = 0; /* depends on the arch */
insn_init(&insn, target_address, x86_64);
insn_get_length(&insn); /* insn_get_length() decodes the entire instruction */
printk("opcode size:%d, instruction length:%d¥n", insn.opcode.size ,insn.length);
```
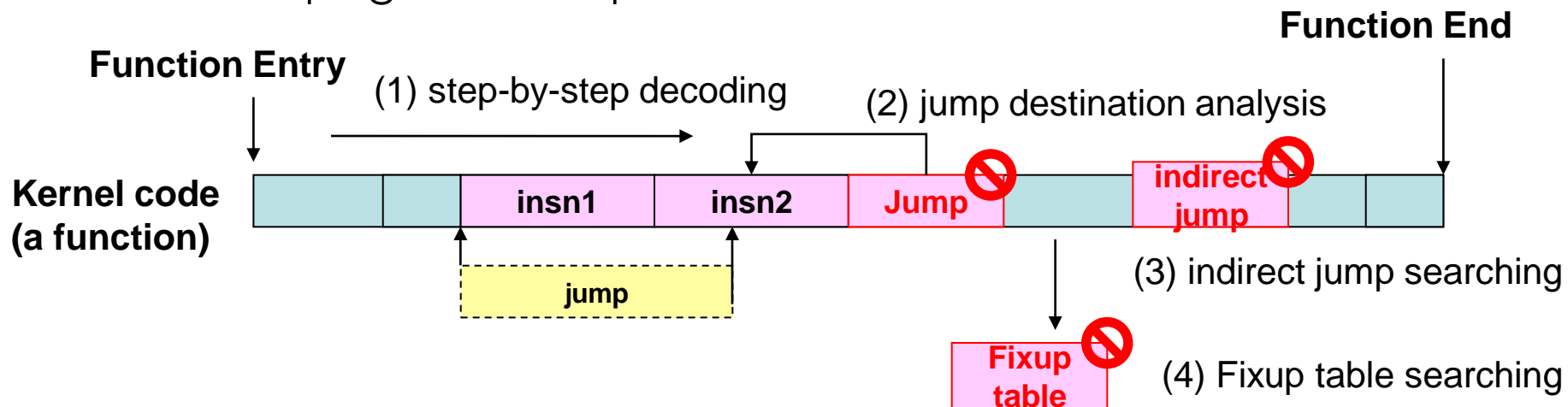
# Linux Technology Center

- ## There are some jump-in issues
  - ### Kernel jumps into the middle of target instructions

**Kernel code** | insn1 | insn2 | | Jump |

**jump**

  - ### Kernel MAY jump into the middle of target instructions

Depends on register value

**Kernel code** | insn1 | insn2 | | Indirect jump | | Other insn |

**jump**

Fixup jump — **Fixup table**

**Kernel code** | insn1 | insn2 | | Pagefault code |

**jump**

- Check a target function to find those jumps
  - Decode an **entire function**
  - Check pagefault fixup table

**Function End**

**Function Entry**

(1) step-by-step decoding

(2) jump destination analysis

**Kernel code (a function)**

| | | insn1 | insn2 | Jump 🚫 | | indirect jump 🚫 | |

jump

(3) indirect jump searching

**Fixup table** 🚫

(4) Fixup table searching

- Reject optimization and just use normal kprobe
  - If a jump destination is the middle of target
  - If the function including indirect jump
  - If the function including an address in fixup-table

*14*

- # Cross modifying code needs a special operation
  - ## Documented method
    - Intel® 64 and IA-32 Architectures Software Developer's Manual Vol. 3   8.1.3
  - ## Stop-machine and modify code
    - This can't use in NMI handler, but kprobes itself doesn't allow to probe NMI handler too.
  - ## Stop-machine is slow, so modifying should be batched.
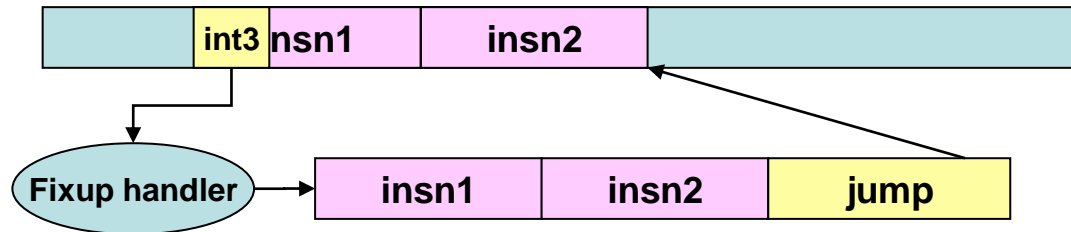
(1) Stop other processors

stop

| CPU0 | | CPU1 | | ... | | CPUN |

(2) Write a jump

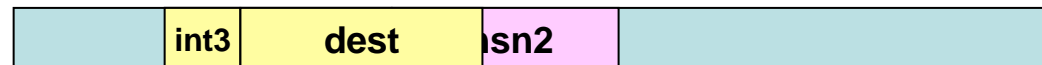| | insn1 | insn2 | |

jump

(3) Serializing and continue to run on other processors

# Int3 bypass method

- Make a bypass by using int3 while XMC
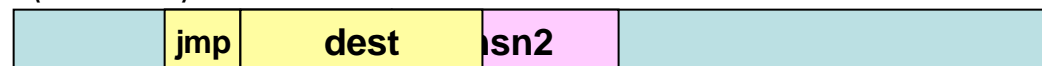- No stop machine required
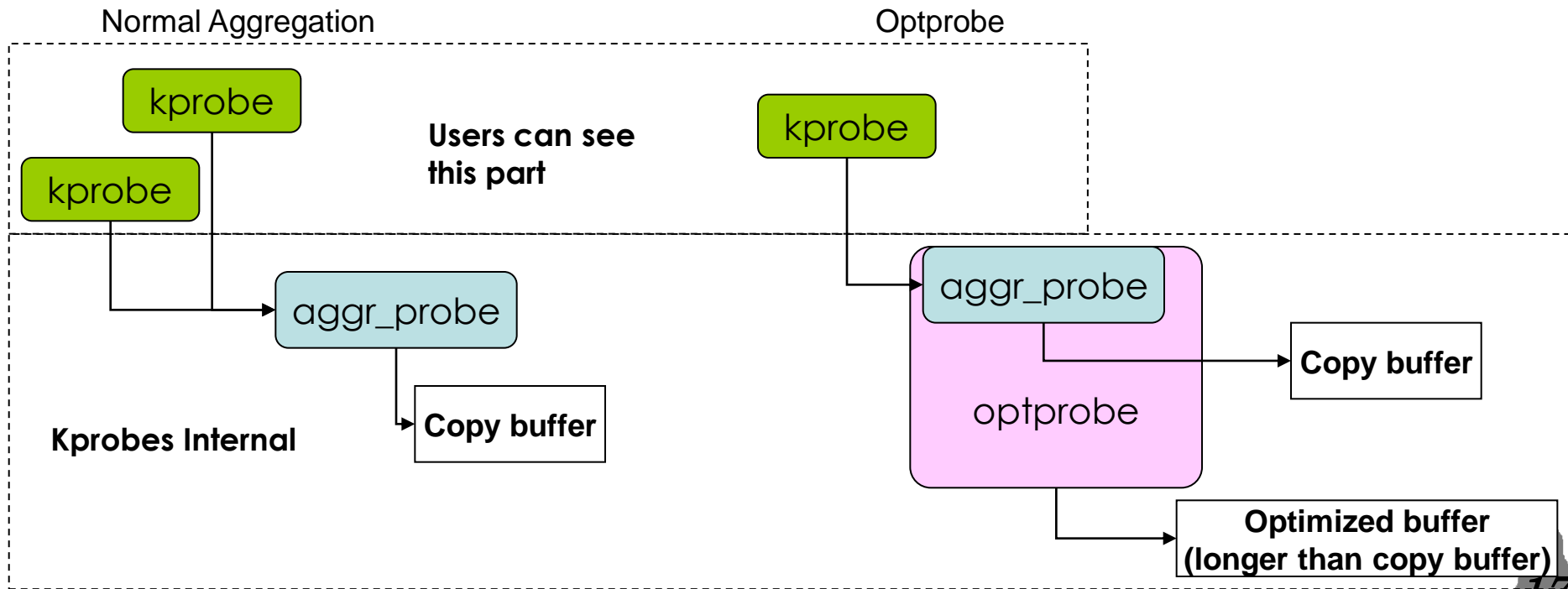- Still be under discussion

(1) Make a bypass

| | int3 | nsn1 | insn2 | |

Fixup handler → | insn1 | insn2 | jump |

(2) Write a jump destination
and sync all processors(send IPI)

| | int3 | dest | nsn2 | |

(3) Write a jump opcode
and sync all processors(send IPI)
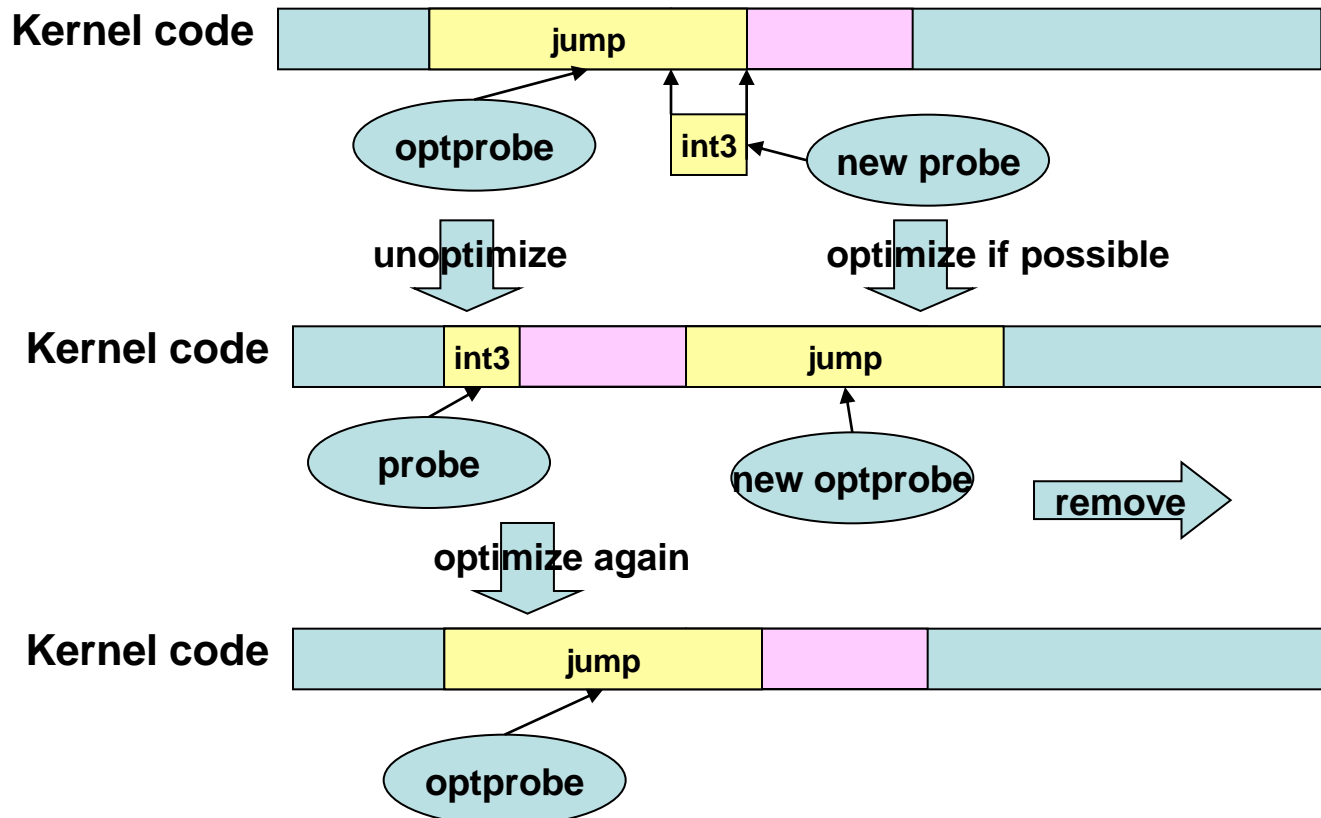
| | jmp | dest | nsn2 | |

# Optimization without changing APIs

- – Optimized kprobe is hidden in aggr_probe
  - Aggr_probe is usually used for aggregating multiple probes on the same address
- – User don't know their probe is optimized or not.

Normal Aggregation                                    Optprobe

kprobe

kprobe

**Users can see this part**

kprobe

aggr_probe

**Copy buffer**

aggr_probe

**Copy buffer**

optprobe

**Kprobes Internal**

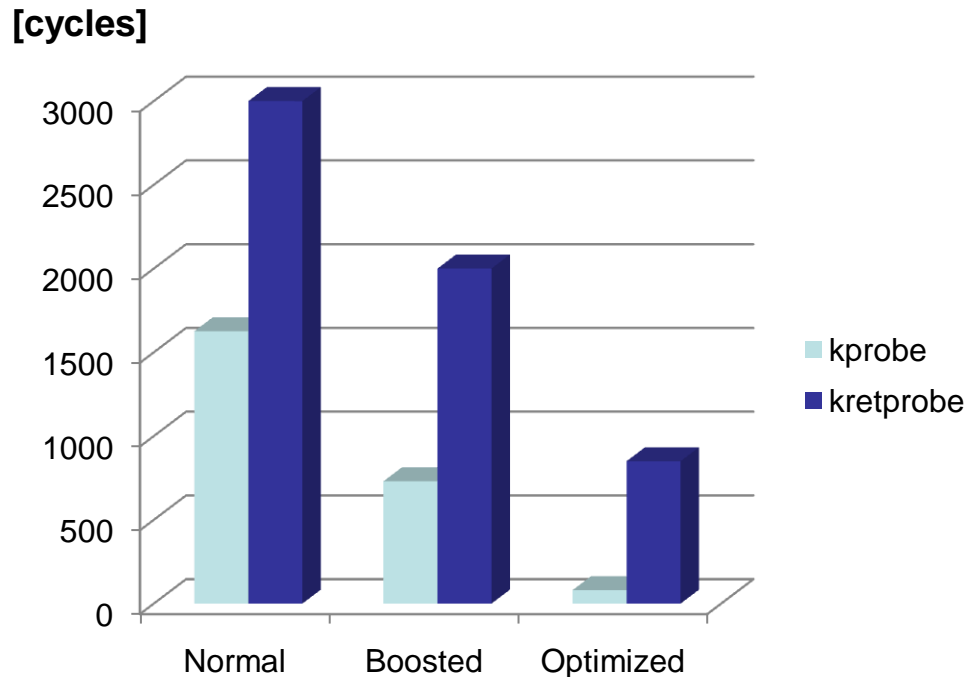**Optimized buffer (longer than copy buffer)**

- Optimization is transparently done (No explicit APIs)
  - Jump code modifying is done in background
  - Some probe state changes requires unoptimizing
    - Unoptimizing is also done in background
  - Only one knob for debugging
    - /proc/sys/debug/kprobes-optimization

# Optimizing/Unoptimizing probes automatically

- Kprobes tries to optimize probes every state change if possible
  - A probe removed from the instruction next to another probe
  - An aggregated probe which has a post_handler is removed

**Kernel code** | jump | | | |

**optprobe** **int3** **new probe**

**unoptimize** **optimize if possible**

**Kernel code** | int3 | | jump | |

**probe** **new optprobe** **remove**

**optimize again**

**Kernel code** | jump | | |

**optprobe**

- ## Some other functions can modify text too
  - Ftrace, alternatives, jump labels
  - Only kprobes is modifying code anywhere
  - Introduce text_reserve interface
    - Checking specified area can be modified by other functions
    - If so, kprobes gives up putting a probe on it.

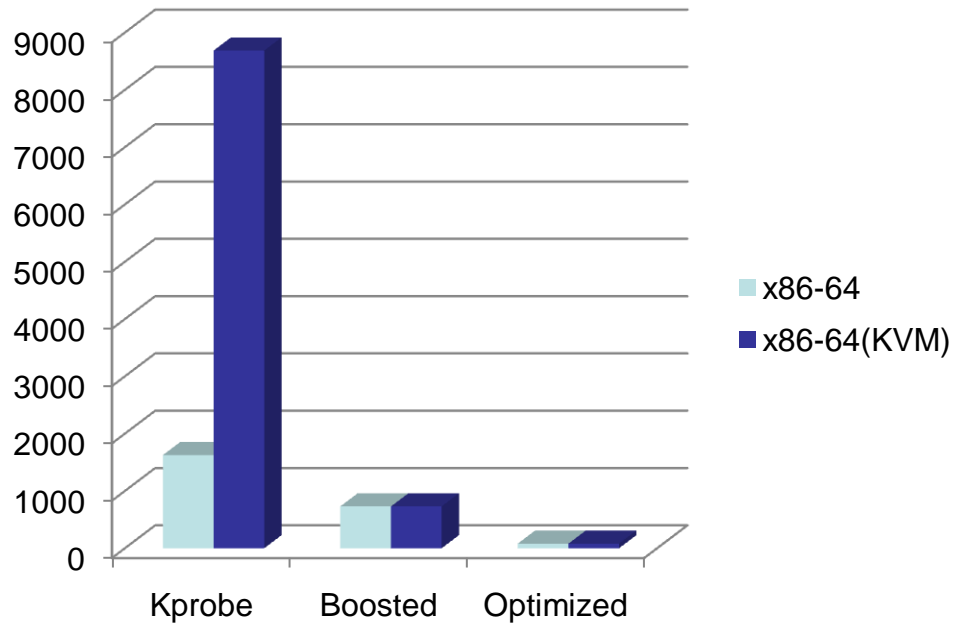- Performance results (unit is cycles)

**[cycles]**

**Intel® Core™ i7 920 (2.67Ghz)**



- kprobe
- kretprobe

Normal    Boosted    Optimized

•Optimization can reduce the overhead to ~100cycles
•Kretprobe is also optimized

# • Performance results on KVM

– On KVM, kprobes is much heavier, because trap is emulated
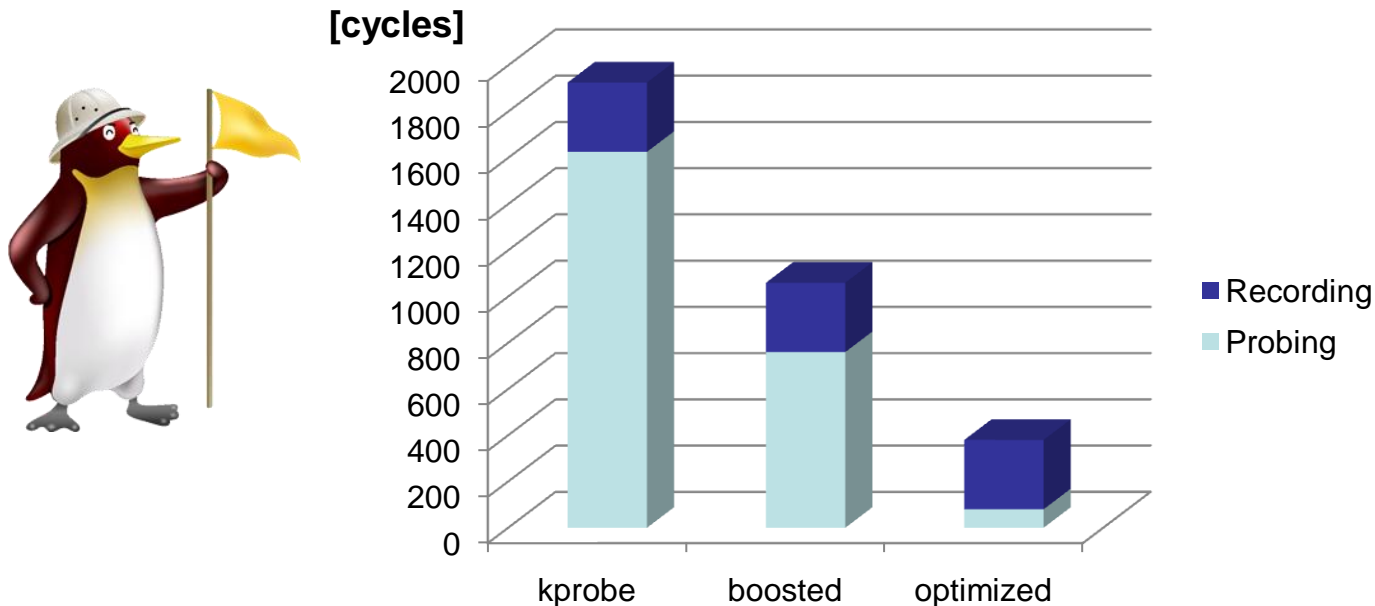
**[cycles]**

**Intel® Core™ i7 920 (2.67Ghz)**



**Optimized and boosted probes can run inside guest.**

*22*

- ## Lower overhead allows us to trace more events
  - Tracing overhead breakdown
    - Probing overhead (depends on optimization)
    - Recording overhead (~300 cycles)



- Total ~400cycles overhead/event allows us to trace **100K** events/sec with just **1~2%** overhead on 3GHz CPU

- # Kprobes
  - – Dynamic/Flexible in-kernel probing function
  - – But heavy, especially with Virtualization

- # Kprobe jump optimization
  - – Drastically reduce overhead of kprobes
  - – Some limitations
  - – Transparent optimization
    - • User need nothing to change
  - – Good performance with Virtualization

- Long history of kprobes jump optimization

- 2005 May: Got an idea for jump optimization
- 2005 Jul: First Prototype Release
- 2005 Aug: 1st Upstream Try
- 2006 Oct: 2nd Upstream Try
- 2007 Jul: 1st Presentation of "djprobe" in OLS
- 2008 – silent but things going forward…
- 2009 Jun: x86 instruction decoder Release
- 2009 Jun: Revised "Optprobe" Release
- 2010 Feb: Optprobe is merged!

- Minimizing instrumentation impacts (kprobes jump optimization)
  - http://lwn.net/Articles/365833/
- Kernel documents
  - Documents/kprobes.txt

**Thank you!**

- Linux is a trademark of Linus Torvalds in the United States, other countries, or both.

- Intel and Core are trademarks or registered trademarks of Intel Corporation in the United States, other countries, or both.

- Other company, product, or service names may be trademarks or service marks of others.