# システムソフトウェア研究の最前線

#### 慶應義塾大学 河野健二

sslab

#### 自己紹介

sslab

- 河野健二(慶應義塾大学理工学部情報工学科)
  - 東大助手, 電通大講師, 慶應義塾大学准教授当を経て, 同教授
- 専門分野
  - オペレーティングシステムおよびシステムソフトウェアディペンダブル・コンピューティングにも少し浮気
- 主な論文
  - IEEE Trans. on Computers, USENIX ATC, ACM EuroSys, IEEE DSN, ACSAC, IEEE ICDCS, ACM VEE, ECOOP...
- 受賞

  - 2000年 情報処理学会山下記念研究賞2000年, 2009年, 2010年, 2013年 情報処理学会論文賞
  - 2014年 日本ソフトウェア科学会ソフトウェア論文賞

  - 2015年 IBM Faculty Award 2015年 ACM Service Award
  - 2016年 日本ソフトウェア科学会基礎研究賞

#### システムソフトウェアの研究トピック



- ACM SIGOPS == 国際的な研究コミュニティ
  - SOSP/OSDI には毎年 600 名程度集まる
  - アクティブな研究者は 150名ぐらい?
- 研究トピックは多岐に渡る
  - Many-core/Multi-core 向けのオペレーティングシステム ◆ GPU のスケジューリング、GPUの仮想化など
  - ライブラリOS, 仮想化支援ハードウェアによるその支援
  - NVM(不揮発性メモリ)を利用したストレージ, KVS, FS
  - システムソフトウェアのバグ解析
    - Linux や Apache, MySQL 等のバグ解析
  - などなど

# システムソフトウェア研究の楽なところ

sslab

- システムソフトウェアは中間管理職
  - "アプリケーション"という名のわがまま上司
  - "ハードウェア" という名の身勝手な部下
- どこが楽なの???

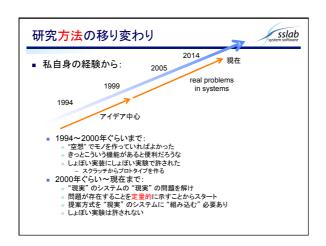


#### システムソフトウェア研究の楽なところ



- 研究テーマを考える必要がない
  - 上司のわがままを聞いてあげればよい
    - ◆ アプリケーションの進化が新しいテーマをもたらしてくれる
  - 部下の身勝手に付き合ってあげればよい
    - ハードウェアの進化が新しいテーマをもたらしてくれる
- ■「もうやることない」と思っても 向こうから研究テーマが やってくる

#### 研究テーマの移り変わり sslab 現在 ■ 私自身の経験から: 2005 仮想化 GPU 🥕 1999 Peer-to-Peer Non-Volatile 1994 Memory Web セキュリティに浮気 ■ 1994~2000年ぐらいまで: LANやPCクラスタを前提とした closed な分散システムの研究2000~2010年ぐらいまで: ▶ Peer-to-Peer などの Internet-scale の分散システムの研究 2005年~現在まで: ● 仮想化をベースとしたクラウドのための基盤技術の研究 ■ 2013年~現在まで: GPU や non-volatile memory など新しめのデバイスのための基盤



## システムソフトウェア研究の辛いところ

sslab

- "空想"の問題を解いても評価されない
- "正しさ" を証明する方法がない
- "良さ"を客観的に比較する方法がない
- 研究に対する評価が厳しい
  - 査読は減点法になってしまう
  - 文句をつけられる箇所があったら reject というノリ

#### どうやって辛さを乗り切るのか?(その1) system software



- "空想"の問題を解いても評価されない
- 実際のシステムで起きている問題を見つけろ!
- いろいろな視点から性能評価をやってみる
  - まずは人為的なベンチマークで極限状況を調べる
  - 直感に反する結果が出ていないかどうか調べる・考える
  - 直感に反する結果を見つけたら・・・
  - 何が起きているのか詳細に調べていく 原因がわかったら・・・
    - 同じ現象が起きる実アプリケーションを探す

Containers or Hypervisors, Which is Better for Database Consolidation?

> in collaboration with Asraa Abdulrazak Ali Mardan

#### Introduction

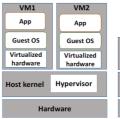


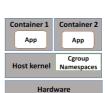
- Virtualization becomes an essential building block in today datacenters
  - Meeting growing demands of resources by providing resources sharing and consolidation
  - Brings significant cost savings
- Major virtualization technologies
  - Hypervisor-based solution like KVM , Vmware ..etc
  - OS-level Virtualization or containers such as LXC , Docker..etc

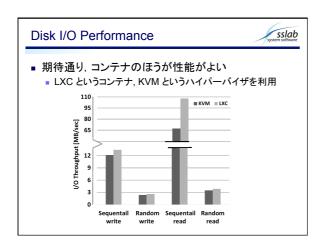
# Containers vs Hypervisors

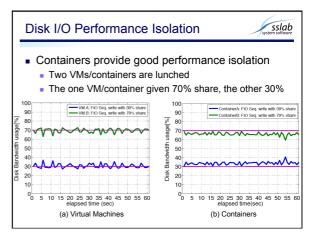


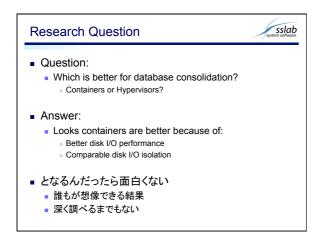
- Trade off between performance and isolation
  - Container provides near native performance with no virtualization overhead
  - Hypervisor provides strong isolation

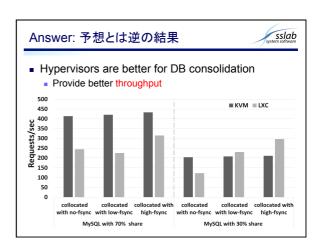


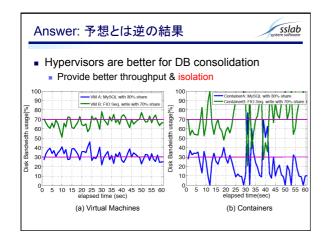


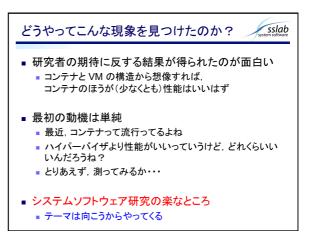












#### 実際には・・・



- 数百通りのパラメータの組み合わせで実験
  - (動きの解析しやすい)簡単なベンチマークを使う
  - とにかく brute-force で頑張る
  - 結果は出ると無邪気に信じる学生は半ベソでした・・・
- fsync の頻度が高いとコンテナの性能が低下する

  - fsync: ファイルの更新をディスクに反映する処理コンテナではジャーナリングの競合が起きるのが原因
- この辺りの分析は experience と expertise の勝利
- fsync を多用する実アプリケーションを探すデータベースなんかどうだろう?

  - We made it!!!

# どうやって辛さを乗り切るのか?(その2) sslab

- "空想"の問題を解いても評価されない
- 実際のシステムで起きている問題を見つける!
- 報告されている問題を全部調べる
  - 先入観は排除せよ
    - フィルタリングしたり絞り込んだりしては<mark>いけない</mark>
  - 先入観があると"想定範囲内"の問題しかみつからない
    - 驚きを伴ういい仕事になりにくい

# Linux のバグから学ぶ

~バグのフィールド調査からコード検査器の実現まで~

慶應義塾大学 河野健二

sslab

#### 本日の話題



- Linuxってバグがないわけではない・・・
  - でも Linux って社会のインフラになってます
- Linux のバグを駆逐しよう
  - OS 固有の知識を活かしたコード検査
    - OS に典型的なバグ・パターンを静的に検出する
      - ヌルポインタかどうかの検査忘れunlock や free などの呼び出し忘れ
      - 暗黙の約束事を忘れている
        - » ブロックしてはいけない関数の中でブロックする関数を呼び出す
- Linux に固有のバグって?
  - これまでは Linux 開発者の経験と勘が頼り

    - Linux には "○○" というバグが多い
       その "○○" っていうバグをとるコード検査器を作りました

#### 本日の話題

sslab

- コード検査器の実現は古典芸能の職人芸
  - ワシの経験と勘が頼りじゃ
- 経験と勘にたよらないコード検査器
  - 膨大なバグ修正記録を活用
    - ソフトウェア工学におけるビッグデータ
  - 37 万件を超える Linux のコード修正記録を分析
    - ◆ 自然言語処理によるパワープレイ
    - Amazon EC2. 100 インスタンスで 1ヶ月.
  - Linux に典型的なバグ・パターンの抽出
  - 割込み関連のバグ・パターンについてコード検査器を実現
  - 最新の Linux においても複数のバグを発見

#### Example of Linux Bugs (1)

sslab system software

- 初歩的なバグでさえ、まだまだ残っている
  - 例: "ヌルポインタ参照" のバグ
    - tun/tap: Fix crashes if open() /dev/net/tun and then poll() it.
    - Author: Mariusz Kozlowski <m.kozlowski@tuxland.pl>

ポインタ変数 tun を参照・tun はヌルではない

struct sock \*sk = tun->s unsigned int mask = 0; return POLLERR:

tun がヌルかどうかを検査 ・これは矛盾

struct sock \*sk; unsigned int mask = 0; if (!tun) return POLLERR;

tun がヌルかどうかを検査 してから tun->sk が正しい

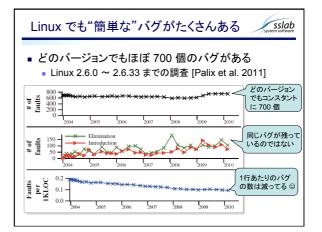
#### Linux における"簡単な"バグ

sslab

■ 簡単な静的解析で見つかるバグを調査

[Palix et al. 2011]

- Null (ヌル検査忘れ):
  - Null を返すかもしれない関数の返値のチェック忘れ
- Inull (Inconsistent null check):
  - ポインタ参照をした後にヌル・チェック
    - さきほどのバグの例
- Block (Calling blocking func in non-blocking context)
  - ブロックしていはいけないコンテキストでブロックする関数を呼ぶ
    - 例:スピンロックを保持したまま、ロックを獲得する
    - 例:ファイルシステムから・・・・
- など 12 種類のバグを検査



## Example of Linux Bug (2)



- "簡単"ではないバグの例
- 静的なコード検査では見つけるのが難しいもの
  - 関数にまたがった解析が必要なケース
  - ◆ 割込みなどの非同期的な振る舞いが絡むケース
  - などなど
- 割込みのタイミングに依存するバグ
  - デバイスの取り外し処理
    - 1. デバイスの割込みを解除する
    - 割込みを受け付けないようにする
  - 2. デバイス管理のためのデータ構造を開放する
  - 複雑なコードでは・・・
    - ♦ 1. と 2. の処理の順番がひっくり返ってしまうことがある

# Example of Linux Bug (2)



■ OS らしいバグ: 割込みのタイミングに依存するバグ



#### Example of Linux Bugs (3)



- Many Linux device drivers assume device perfection [Kadav et al. 2009]
- Example: Infinite polling
  - Driver waiting for device to enter particular state
  - If device not working correctly, the loop never ends. Hang

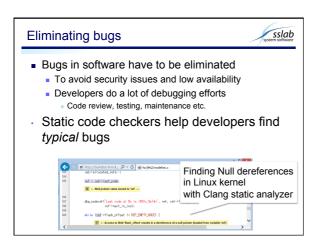
```
static int amd8111e_read_phy(......)
{
...
    reg_val = readl(mmio + PHY_ACCESS);
    while (reg_val & PHY_CMD_ACTIVE)
        reg_val = readl(mmio + PHY_ACCESS)
.
}
AMD 8111e network driver(amd8111e.c)
```

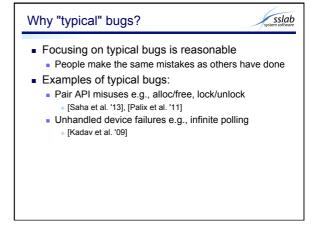
#### Example of Linux Bugs (3)

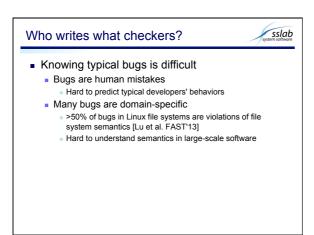


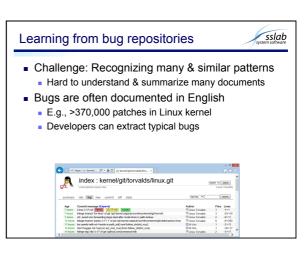
- Solution: add the code for timeout
  - If timeout occurs, recover code is invoked

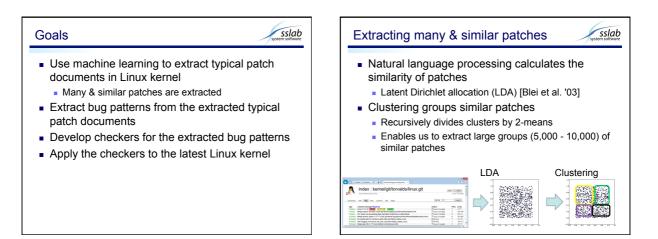
```
static int amd8111e_read_phy(......)
{
...
timeout = 0;
reg_val = read1(mmio + PHY_ACCESS);
while (reg_val & PHY_CMD_ACTIVE) {
    reg_val = read1(mmio + PHY_ACCESS)
    if (timeout++ >= 200)
        __shadow_recover();
}
.
AMD 8111e network driver(amd8111e.c)
```











#### LDA in short



- LDA infers latent topics in documents
  - A document is regarded as probability sets of topics
    - The similarity of two documents is the distance between the probability sets for them
  - Keywords characterizing patches can be obtained

In function devkmsg\_read/writev/llseek/poll/open()..., the function raw\_spin\_lock/unlock is used, there is potential deadlock case happening. CPU1: thread1 doing the cat /dev/kmsg: raw\_spin\_lock(&logbuf\_lock); while (user->seq == log\_next\_seq) { when thread1 run here, a this time one interrupt is coming on CPU1 and running based on this thread,if the interrupt handle called the printk which need the log buf\_lock spin also, it will cause deadlock.

topic A: 0.113, topic B: 0.055, topic C: 0.04 topic D: 0.038, topic E: 0.038, topic F: . Keywords:lock, unlock, spin, lock, protect,

## Analyzing Linux patch documents



- 370,403 patch documents are analyzed
  - Linux 2.6-rc2 ~ 3.12-rc5 (April 2005 October 2013)
  - Merge commits are excluded
  - Analyzer has been implemented on Hadoop MapReduce
  - LDA from Apache Mahout
- Result: 66 clusters
  - Linux had topics for general software bugs, OSs, devices, CPU platforms, etc.

#### Result 1/3: common bugs



- Clusters for bugs in general software
- Null dereferences, memory leaks, lock/unlock
  - Typical software bugs in the Linux kernel
- Example document: memory leak
  - Topic keywords: memory, leak, cpu,hotplug
  - Misuse of kmalloc() and kfree()

commit 003615301, 2nd paragraph

USB: io\_ti: fix port-data memory leak
Fix port-data memory leak by moving port data allocation
and deallocation to port probe and port remove.

### Result 2/3: hardware



- Clusters for CPU platforms and devices
  - ARM, X86, GPU, USB, NIC, etc.
  - Major hardware-speicifc issues
- Example document: ARM
  - Topic keywords: arm, mach, h, asm
  - Problems deriving from ARM features

#### commit 9cff337, 3rd paragraph

So far as I am aware this problem is ARM specific, because only ARM supports software change of the CPU (memory system) byte sex, however the partition table parsing is in generic MTD code.

#### Result 3/3: common OS features



- Clusters for common OS features
  - Interrupt handling, buffer cache, DMA, etc
  - Typical implementation issues in OSs
- Example document: interrupt handling
  - Topic keywords: irq, interrupt, msi
  - Problems around masking interrupts

commit ea6dedd, 2nd paragraph

The current OMAP GPIO IRQ framework doesn't use the do\_edge\_IRQ, do\_level\_IRQ handlers, but instead calls do\_simple\_IRQ. This doesn't handle disabled interrupts properly, so drivers will still get interrupts after calling disable\_irq. ....

#### Extracting bug patterns from a cluster

sslab

- The cluster for interrupts are expected to typical bugs in OSs
  - The cluster remains too large (5,334 patches)
- Topic keywords help us extract interesting subclusters
  - A sub-cluster with keyword "free" is expected to have bugs around free and interrupts
- The sub-cluster with keyword "free" contains 364 patches
  - 160 are identified as bugs

