

Системные баги под нагрузкой

Миша Жилин, Postgres Professional

C++ Russia, 17 мая 2026

Миша Жилин

- Физтех, 2008 год выпуска
- 18 лет занимаюсь
производительностью программ
- Больше исследователь
- Меньше разработчик
- Обожаю Open Source, PostgreSQL и
FreeBSD
- Работаю 5,5 лет в Postgres
Professional



О чём доклад

- Операционные системы работают как часы



О чём доклад

- Операционные системы работают как часы
- Базы данных всегда возвращают корректный результат



О чём доклад

- Операционные системы работают как часы
- Базы данных всегда возвращают корректный результат
- Системный софт - надёжный, качественный, проверенный



О чём доклад

- Операционные системы работают как часы
- Базы данных всегда возвращают корректный результат
- Системный софт - надёжный, качественный, проверенный
- Или нет?!



О чём доклад

- Операционные системы работают как часы
- Базы данных всегда возвращают корректный результат
- Системный софт - надёжный, качественный, проверенный
- Или нет?!
- Расскажу только то, что видел



О чём доклад

- Операционные системы работают как часы
- Базы данных всегда возвращают корректный результат
- Системный софт - надёжный, качественный, проверенный
- Или нет?!
- Расскажу только то, что видел
- Что не видел, то не расскажу



О чём доклад

- Поиск невозможного



О чём доклад

- Поиск невозможного
- Несколько историй из личного опыта



О чём доклад

- Поиск невозможного
- Несколько историй из личного опыта
- Включая много попыток поехать кукухой



О чём доклад

- Поиск невозможного
- Несколько историй из личного опыта
- Включая много попыток поехать кукухой
- Мало кода



О чём доклад

- Поиск невозможного
- Несколько историй из личного опыта
- Включая много попыток поехать кукухой
- Мало кода
- Много QR кодов (а может и немного)



Разгон

Intel SSD, 2019 год

Intel SSD, 2019 год

- Пробный период запуска мобильного оператора



Intel SSD, 2019 год

- Пробный период запуска мобильного оператора
- Отказ в распределенном дисковом пространстве CEPH



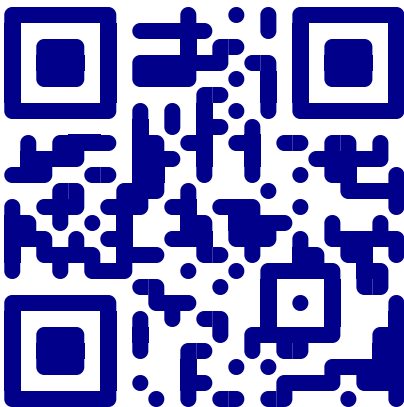
Intel SSD, 2019 год

- Пробный период запуска мобильного оператора
- Отказ в распределенном дисковом пространстве CEPH
- Зависание дисков Intel SSD S45210 и S4610



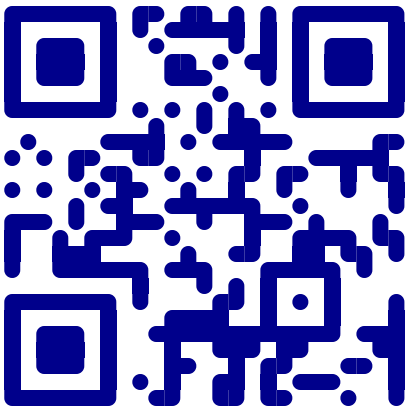
Intel SSD, 2019 год

- Intel SSD D3-S4510 and Intel SSD D3-S4610 series 1.92 TB and 3.84 TB drives become unresponsive after 1,700 cumulative idle power-on hours
- Currently addressed in the Maintenance Release 1 (MR1) of Intel firmware (version XCV10110) as of March 2019



Intel SSD, 2019 год

- Intel удалил все упоминания
- <https://forums.servethehome.com/index.php?threads/psa-errors-with-intel-r-solid-state-drive-dc-s4510-and-s4610-series.26196/>



WAL corruption, 2020 год

WAL corruption, 2020 год

- Смена работы
 - Инженер по отказоустойчивости в кластерной группе
 - 4ая неделя испытательного срока



WAL corruption, 2020 год

- Смена работы
 - Инженер по отказоустойчивости в кластерной группе
 - 4ая неделя испытательного срока
- Проблема от клиента



WAL corruption, 2020 год

- Смена работы
 - Инженер по отказоустойчивости в кластерной группе
 - 4ая неделя испытательного срока
- Проблема от клиента
- База данных на 2 серверах: мастер и реплика



WAL corruption, 2020 год

- Смена работы
 - Инженер по отказоустойчивости в кластерной группе
 - 4ая неделя испытательного срока
- Проблема от клиента
- База данных на 2 серверах: мастер и реплика
- Ошибка репликации
 - [2020-10-27 12:03:02 MSK u= d= h= p=14570 l=99]
 - LOG: incorrect resource manager data checksum in record at 61E7/ED365B20]

WAL corruption, 2020 год

- Смена работы
 - Инженер по отказоустойчивости в кластерной группе
 - 4ая неделя испытательного срока
- Проблема от клиента
- База данных на 2 серверах: мастер и реплика
- Ошибка репликации
 - [2020-10-27 12:03:02 MSK u= d= h= p=14570 l=99]
 - LOG: incorrect resource manager data checksum in record at 61E7/ED365B20]
- Неверная контрольная сумма - повреждение WAL журналов

WAL corruption, 2020 год

- Смена работы
 - Инженер по отказоустойчивости в кластерной группе
 - 4ая неделя испытательного срока
- Проблема от клиента
- База данных на 2 серверах: мастер и реплика
- Ошибка репликации
 - [2020-10-27 12:03:02 MSK u= d= h= p=14570 l=99]
 - LOG: incorrect resource manager data checksum in record at 61E7/ED365B20]
- Неверная контрольная сумма - повреждение WAL журналов
- Материал для анализа: record at 61E7/ED365B20

003100 D0 BA D0 B0 D0 B7 D0 B0 D0 BD D0 B8 D0 B5 20 D1 83 D1 81 D0
003120 B6 00 72 07 12 00 30 40 37 31 38 38 20 30 33 30 32 20 30 38
003140 20 34 20 30 33 20 39 30 30 34 39 20 32 34 32 13 30 30 31 31 4 03 90049 242!!0011
003160 30 36 33 31 00 00 00 00 01 00 2F 00 1B 00 30 40 37 31 38 38 0631 © / ← 0@7188
003180 20 30 33 30 32 D1 82 D0 B8 D0 B7 D0 B8 D1 80 D0 BE D0 B2 D0 0302
003200 B0 D0 BD D0 BD D1 8B D1 85 20 D1 80 D0 B0 D0 B1 D0 BE D1 87

003100 D0 BA D0 B0 D0 B7 D0 B0 D0 BD D0 B8 D0 B5 20 D1 83 D1 81 D0
003120 B6 D1 83 D0 B3 20 D0 B2 20 D1 81 D1 84 D0 B5 D1 80 D0 B5 20
003140 D0 98 D0 9A D0 A2 20 D0 BF D0 BE 20 D0 B0 D1 82 D1 82 D0 B5
003160 D1 81 D1 82 D0 B0 D1 86 D0 B8 D0 B8 20 D0 B0 D0 B2 D1 82 D0
003180 BE D0 BC D0 B0 D1 82 D0 B8 D0 B7 D0 B8 D1 80 D0 BE D0 B2 D0
003200 B0 D0 BD D0 BD D1 8B D1 85 20 D1 80 D0 B0 D0 B1 D0 BE D1 87

WAL corruption, 2020 год

- Ошибка диска/OS: и всё что над ним RAID контролер, драйвера, блочной системы, файловой системы



WAL corruption, 2020 год

- Ошибка диска/OS: и всё что над ним RAID контролер, драйвера, блочной системы, файловой системы
- Ошибка СУБД (не исключена вероятность “гуляющего” указателя)

WAL corruption, 2020 год

- Ошибка диска/OS: и всё что над ним RAID контролер, драйвера, блочной системы, файловой системы
- Ошибка СУБД (не исключена вероятность “гуляющего” указателя)
- **Ошибка памяти:**
 - адрес повреждения в WAL файле – 0x366500 (выровнен на 256 байт)
 - размер повреждения - 64 байта (размер кэшлайна процессора)

WAL corruption, 2020 год

- Через месяц ещё одно падение



WAL corruption, 2020 год

- Через месяц ещё одно падение
- Потом ещё несколько



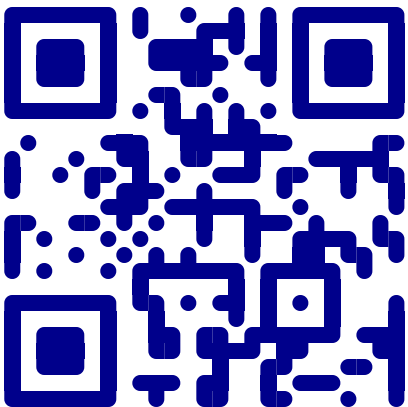
WAL corruption, 2020 год

- Через месяц ещё одно падение
- Потом ещё несколько
- Но всегда портился только WAL, данные проверяли - они все верные



WAL corruption, 2020 год

- Двойной WAL буфер - защита от disk/memory corruption
- Настройка `wal_sender_check_crc` в PostgresPro Enterprise



Наступил 2021 год

WAL corruption, 2021 год

- Изоляция проблемного сервера и доступ для экспериментов



WAL corruption, 2021 год

- Изоляция проблемного сервера и доступ для экспериментов
- BullSequana S series
 - Несколько материнских плат
 - 4 терабайта памяти
 - 8 процессоров



WAL corruption, 2021 год

- Изоляция проблемного сервера и доступ для экспериментов
- BullSequana S series
 - Несколько материнских плат
 - 4 терабайта памяти
 - 8 процессоров
- Недельные тесты
 - Ноль результатов



WAL corruption, 2021 год

- Изоляция проблемного сервера и доступ для экспериментов
- BullSequana S series
 - Несколько материнских плат
 - 4 терабайта памяти
 - 8 процессоров
- Недельные тесты
 - Ноль результатов
- MemTest
 - Ноль результатов



WAL corruption, 2021 год

- Полный ноль



WAL corruption, 2021 год

- Полный ноль
- Фрустрации



WAL corruption, 2021 год

- Полный ноль
- Фрустрации
 - А что ищем?



WAL corruption, 2021 год

- Полный ноль
- Фрустрации
 - А что ищем?
 - Как задать вопрос?



WAL corruption, 2021 год

- Полный ноль
- Фрустрации
 - А что ищем?
 - Как задать вопрос?
 - Что я не понимаю?



WAL corruption, 2021 год

- Полный ноль
- Фрустрации
 - А что ищем?
 - Как задать вопрос?
 - Что я не понимаю?
- Google it, опытные коллеги



WAL corruption, 2021 год

- Полный ноль
- Фрустрации
 - А что ищем?
 - Как задать вопрос?
 - Что я не понимаю?
- Google it, опытные коллеги
- Время менять подходы



WAL corruption, 2021 год

- Основная гипотеза - это память



WAL corruption, 2021 год

- Основная гипотеза - это память
- stress-ng - утилита нагрузки на железные ресурсы

- `stress-ng --vm-rw 1000 --vm-rw-bytes 2500G --verify --metrics-brief -t 60m`



WAL corruption, 2021 год

- Основная гипотеза - это память
- stress-ng - утилита нагрузки на железные ресурсы
 - `stress-ng --vm-rw 1000 --vm-rw-bytes 2500G --verify --metrics-brief -t 60m`
- Никаких ошибок программы



WAL corruption, 2021 год

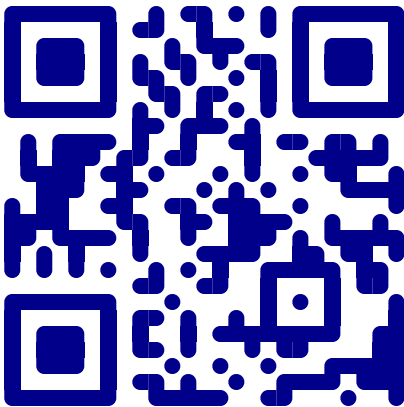
- Основная гипотеза - это память
- stress-ng - утилита нагрузки на железные ресурсы
 - `stress-ng --vm-rw 1000 --vm-rw-bytes 2500G --verify --metrics-brief -t 60m`
- Никаких ошибок программы
- Зато Fiber Channel начал мигать
 - `2021 Feb 15 23:43:57 dn-MDS-C9148S-1 %PORT-5-IF_DOWN_LINK_FAILURE: %$VSAN2221%$ Interface fc1/15 is down (Link failure loss of signal)`

WAL corruption, 2021 год

- Основная гипотеза - это память
- stress-ng - утилита нагрузки на железные ресурсы
 - `stress-ng --vm-rw 1000 --vm-rw-bytes 2500G --verify --metrics-brief -t 60m`
- Никаких ошибок программы
- Зато Fiber Channel начал мигать
 - `2021 Feb 15 23:43:57 dn-MDS-C9148S-1 %PORT-5-IF_DOWN_LINK_FAILURE: %$VSAN2221%$ Interface fc1/15 is down (Link failure loss of signal)`
- Время и идеи закончились

WAL corruption, 2021 год

- Через полгода статья от техподдержки производителя железа
 - <https://habr.com/ru/companies/atosservers/articles/562376/>
- Обновление драйвера FC линков (драйвер устройства операционной системы Linux).

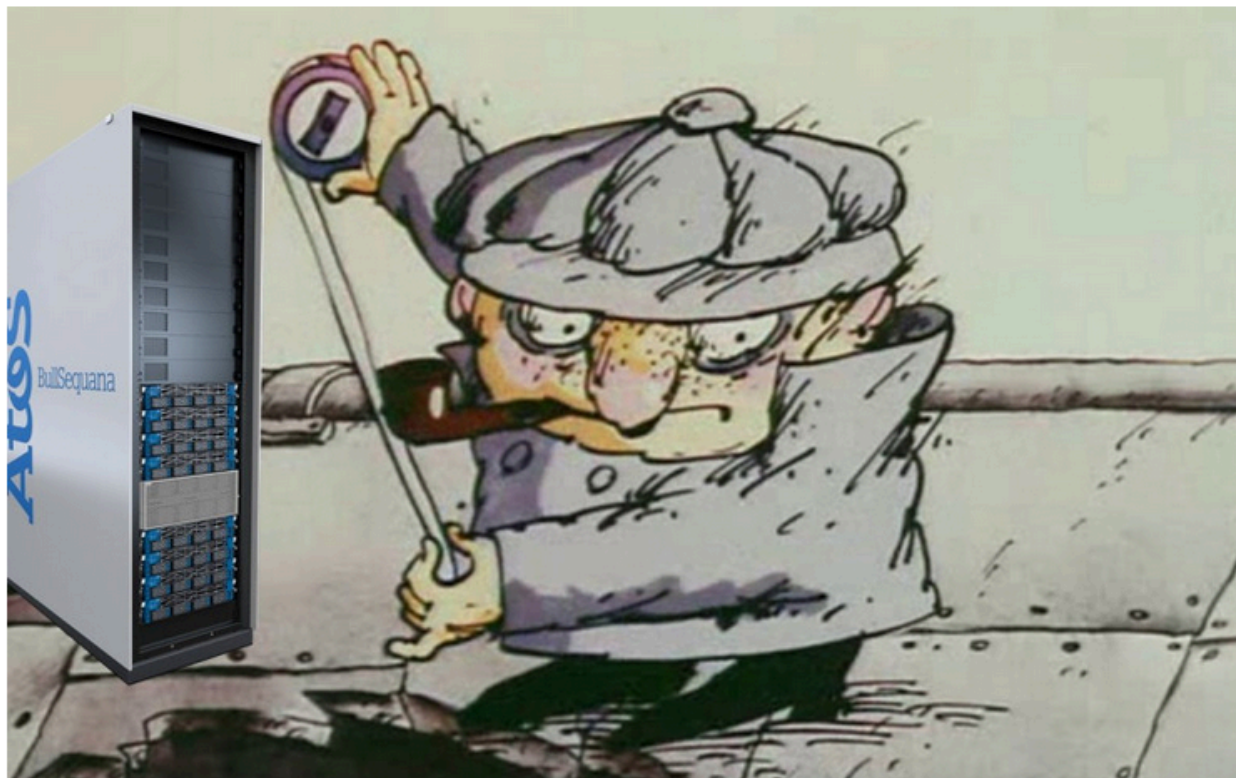


sharkattack 18 июн 2021 в 18:24

Следствие вели: пропажа FC-линков HBA Emulex на сервере Atos BullSequana S1600

🕒 6 мин 👁 3.8К

Блог компании Atos, Высоконагруженные системы*, Тестирование IT-систем*, IT-компании



Цитатник

- "За всё время этой проблемы пошатали и потрогали всё – саму ОС, FW, прошивку HW сервера, настройки HW сервера, параметры GRUB, настройки фабрики, свитчей и линков..."
- "Техподдержка врёт (ну или добросовестно заблуждается), поэтому приходится старательно все проверять самому"
- "Трогать и шатать при траблшутинге надо всё!"



Итоги

- Обновления драйвера и реализации двойного буфера WAL журналов
- Больше повреждений WAL не было
- Кто виноват - до сих пор загадка
- Это успех или нет?
- Негативный опыт - тоже опыт

Полезные моменты

- Нужен стенд. Люди, увы - не телепаты
- Если тест не выявляет дефект, то наверно не то нагружаем или не достаточно
- Полезно добавлять генераторы нагрузки на ресурсы ОС
- Пример - **stress-ng**



Прошли годы...

Прошли годы

- Меня выбросили из кластерной группы как непрофпригодного
 - "Иди делай группу нагрузочных исследователей" 😜

Прошли годы

- Меня выбросили из кластерной группы как непрофпригодного
 - "Иди делай группу нагрузочных исследователей" 🤪
- Коллаборация с нашей техподдержкой
 - Они нам задачи
 - Мы им помощь в исследованиях



Прошли годы

- Меня выбросили из кластерной группы как непрофпригодного
 - "Иди делай группу нагрузочных исследователей" 😜
- Коллаборация с нашей техподдержкой
 - Они нам задачи
 - Мы им помощь в исследованиях
- История с подвисанием сервера на Linux 5.15

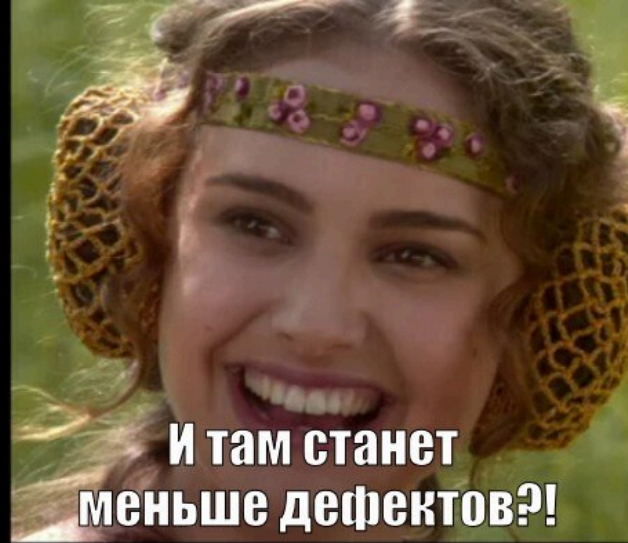
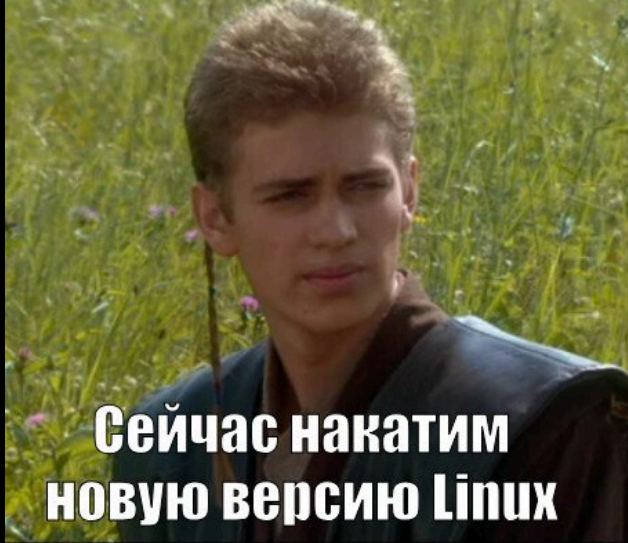


Прошли годы

- Меня выбросили из кластерной группы как непрофпригодного
 - "Иди делай группу нагрузочных исследователей" 😜
- Коллаборация с нашей техподдержкой
 - Они нам задачи
 - Мы им помощь в исследованиях
- История с подвисанием сервера на Linux 5.15
- Проблема масштабирования на lookup-а файлов в XFS

Прошли годы

- Меня выбросили из кластерной группы как непрофпригодного
 - "Иди делай группу нагрузочных исследователей" 😜
- Коллаборация с нашей техподдержкой
 - Они нам задачи
 - Мы им помощь в исследованиях
- История с подвисанием сервера на Linux 5.15
- Проблема масштабирования на lookup-а файлов в XFS
- **Общая рекомендация:** обновитесь до новой версии операционной системы



segfault, 2024

segfault, 2024

- 30 декабрь 2024 года



segfault, 2024

- 30 декабрь 2024 года
- Миграция данных



segfault, 2024

- 30 декабрь 2024 года
- Миграция данных
- Упала база данных по segfault



segfault, 2024

- 30 декабрь 2024 года
- Миграция данных
- Упала база данных по segfault
- ```
[Mon Dec 30 17:49:50 2024] postgres[4092660]: segfault at 55edf40f20d0 ip
000055edf40f20d0 sp 00007ffc7a7773c8 error 15 in
postgres[55edf3f67000+2a9000] likely on CPU 40 (core 16, socket 1)
```

# segfault, 2024

## ■ checkpointer (4092660)

```
#0 0x000055edf40f20d0 in ?? ()
#1 0x000055edf3f1a370 in hash_search (hashp=0x55edf5c1ada8, keyPtr=0x7ffc7a777480, action=HASH_ENTER, foundPtr=0x7ffc7a777480) at hashtable.c:131
#2 0x000055edf3dd2dd2 in RememberSyncRequest (ftag=0x7ffc7a777480, type=<optimized out>) at sync.c:565
#3 0x000055edf3dd2fb5 in RegisterSyncRequest (ftag=ftag@entry=0x7ffc7a777480, type=type@entry=SYNC_REQUEST, retryOnError=retryOnError@entry=0) at sync.c:612
```

## ■ startup (4092662)

```
#0 0x000055edf40f20d0 in ?? ()
#1 0x000055edf3f1a370 in hash_search (hashp=0x55edf5c458f8, keyPtr=0x7ffc7a7781ec, action=HASH_FIND, foundPtr=0x0) at hashtable.c:131
#2 0x000055edf3b40603 in XLogPrefetcherIsFiltered (blockno=1523809, rnode=... , prefetcher=0x55edf5c375f8) at xlogprefetcher.c:101
#3 XLogPrefetcherNextBlock (pgsr_private=94480518837752, lsn=<optimized out>) at xlogprefetcher.c:693
#4 0x000055edf3b40f51 in lrq_prefetch (lrq=<optimized out>) at xlogprefetcher.c:256
#5 lrq_complete_lsn (lsn=67516543241680, lrq=0x55edf5c3aff8) at xlogprefetcher.c:294
```

# segfault, 2024

## ■ checkpointer (4092660)

```
#0 0x000055edf40f20d0 in ?? ()
#1 0x000055edf3f1a370 in hash_search (hashp=0x55edf5c1ada8, keyPtr=0x7ffc7a777480, action=HASH_ENTER, foundPtr=0x7ffc7a777480) at hashtable.c:131
#2 0x000055edf3dd2dd2 in RememberSyncRequest (ftag=0x7ffc7a777480, type=<optimized out>) at sync.c:565
#3 0x000055edf3dd2fb5 in RegisterSyncRequest (ftag=ftag@entry=0x7ffc7a777480, type=type@entry=SYNC_REQUEST, retryOnError=<optimized out>) at sync.c:615
```

## ■ startup (4092662)

```
#0 0x000055edf40f20d0 in ?? ()
#1 0x000055edf3f1a370 in hash_search (hashp=0x55edf5c458f8, keyPtr=0x7ffc7a7781ec, action=HASH_FIND, foundPtr=0x0) at hashtable.c:131
#2 0x000055edf3b40603 in XLogPrefetcherIsFiltered (blockno=1523809, rnode=... , prefetcher=0x55edf5c375f8) at xlogprefetcher.c:101
#3 XLogPrefetcherNextBlock (pgsr_private=94480518837752, lsn=<optimized out>) at xlogprefetcher.c:693
#4 0x000055edf3b40f51 in lrq_prefetch (lrq=<optimized out>) at xlogprefetcher.c:256
#5 lrq_complete_lsn (lsn=67516543241680, lrq=0x55edf5c3aff8) at xlogprefetcher.c:294
```

# segfault, 2024

## ■ checkpointer (4092660)

```
#0 0x000055edf40f20d0 in ?? ()
#1 0x000055edf3f1a370 in hash_search (hashp=0x55edf5c1ada8, keyPtr=0x7ffc7a777480, action=HASH_ENTER, foundPtr=0x7ffc7a777480) at hashtable.c:131
#2 0x000055edf3dd2dd2 in RememberSyncRequest (ftag=0x7ffc7a777480, type=<optimized out>) at sync.c:565
#3 0x000055edf3dd2fb5 in RegisterSyncRequest (ftag=ftag@entry=0x7ffc7a777480, type=type@entry=SYNC_REQUEST, retryOnError=retryOnError@entry=0) at sync.c:611
```

## ■ startup (4092662)

```
#0 0x000055edf40f20d0 in ?? ()
#1 0x000055edf3f1a370 in hash_search (hashp=0x55edf5c458f8, keyPtr=0x7ffc7a7781ec, action=HASH_FIND, foundPtr=0x0) at hashtable.c:131
#2 0x000055edf3b40603 in XLogPrefetcherIsFiltered (blockno=1523809, rnode=..., prefetcher=0x55edf5c375f8) at xlogprefetcher.c:101
#3 XLogPrefetcherNextBlock (pgsr_private=94480518837752, lsn=<optimized out>) at xlogprefetcher.c:693
#4 0x000055edf3b40f51 in lrq_prefetch (lrq=<optimized out>) at xlogprefetcher.c:256
#5 lrq_complete_lsn (lsn=67516543241680, lrq=0x55edf5c3aff8) at xlogprefetcher.c:294
```

# segfault, 2024

- Два процесса одновременно



# segfault, 2024

- Два процесса одновременно
- **hash\_search**
  - Фундаментальный код
  - Горячий код (часто выполняемый)
  - Вероятность найти ошибку почти нулевая



# segfault, 2024

- Два процесса одновременно
- **hash\_search**
  - Фундаментальный код
  - Горячий код (часто выполняемый)
  - Вероятность найти ошибку почти нулевая
- Трудность



# segfault, 2024

- Начинаем искать зацепки



# segfault, 2024

- Начинаем искать зацепки
- Вводные
  - Дамп памяти (coredump)
  - Конечное состояние регистров
  - Ассемблеровские инструкции



# segfault, 2024

- Начинаем искать зацепки
- Вводные
  - Дамп памяти (coredump)
  - Конечное состояние регистров
  - Ассемблеровские инструкции
- Поиск непротиворечивого сценария



# segfault, 2024

- Начинаем искать зацепки
- Вводные
  - Дамп памяти (coredump)
  - Конечное состояние регистров
  - Ассемблеровские инструкции
- Поиск непротиворечивого сценария
- Никаких вариантов, кроме одного, но нереального



# segfault, 2024

```
#0 0x000055edf40f20d0 in ?? ()
#1 0x000055edf3f1a370 in hash_search (hashp=0x55edf5c458f8, keyPtr=0x7ffc7a7781ec, action=HASH_FIND, foundPtr=0x0)
```

```
0x000055edf3f5870b <+123>: cmp $0xb,%r10d
0x000055edf3f5870f <+127>: ja 0x55edf3f586b0 <hash_bytes+32>
0x000055edf3f58711 <+129>: lea 0x1999b8(%rip),%rdi
0x000055edf3f58718 <+136>: movslq(%rdi,%r10,4),%rsi
0x000055edf3f5871c <+140>: add %rdi,%rsi
0x000055edf3f5871f <+143>: jmp *%rsi
0x000055edf3f58721 <+145>: nopl 0x0(%rax)
```

# segfault, 2024

```
#0 0x000055edf40f20d0 in ?? ()
#1 0x000055edf3f1a370 in hash_search (hashp=0x55edf5c458f8, keyPtr=0x7ffc7a7781ec, action=HASH_FIND, foundPtr=0x0)
```

```
0x000055edf3f5870b <+123>: cmp $0xb,%r10d
0x000055edf3f5870f <+127>: ja 0x55edf3f586b0 <hash_bytes+32>
0x000055edf3f58711 <+129>: lea 0x1999b8(%rip),%rdi
0x000055edf3f58718 <+136>: movslq(%rdi,%r10,4),%rsi
0x000055edf3f5871c <+140>: add %rdi,%rsi
0x000055edf3f5871f <+143>: jmp *%rsi
0x000055edf3f58721 <+145>: nopl 0x0(%rax)
```

# segfault, 2024

```
#0 0x000055edf40f20d0 in ?? ()
#1 0x000055edf3f1a370 in hash_search (hashp=0x55edf5c458f8, keyPtr=0x7ffc7a7781ec, action=HASH_FIND, foundPtr=0x0)
```

```
0x000055edf3f5870b <+123>: cmp $0xb,%r10d
0x000055edf3f5870f <+127>: ja 0x55edf3f586b0 <hash_bytes+32>
0x000055edf3f58711 <+129>: lea 0x1999b8(%rip),%rdi
0x000055edf3f58718 <+136>: movslq (%rdi,%r10,4),%rsi
0x000055edf3f5871c <+140>: add %rdi,%rsi
0x000055edf3f5871f <+143>: jmp *%rsi
0x000055edf3f58721 <+145>: nopl 0x0(%rax)
```

# segfault, 2024

```
#0 0x000055edf40f20d0 in ?? ()
#1 0x000055edf3f1a370 in hash_search (hashp=0x55edf5c458f8, keyPtr=0x7ffc7a7781ec, action=HASH_FIND, foundPtr=0x0)
```

```
0x000055edf3f5870b <+123>: cmp $0xb,%r10d
0x000055edf3f5870f <+127>: ja 0x55edf3f586b0 <hash_bytes+32>
0x000055edf3f58711 <+129>: lea 0x1999b8(%rip),%rdi
0x000055edf3f58718 <+136>: movslq (%rdi,%r10,4),%rsi
0x000055edf3f5871c <+140>: add %rdi,%rsi
0x000055edf3f5871f <+143>: jmp *%rsi
0x000055edf3f58721 <+145>: nopl 0x0(%rax)
```

# segfault, 2024

```
0x000055edf3f5870b <+123>: cmp $0xb,%r10d
0x000055edf3f5870f <+127>: ja 0x55edf3f586b0 <hash_bytes+32>
0x000055edf3f58711 <+129>: lea 0x1999b8(%rip),%rdi # %rdi = (%rip+0x1999b8) | %rdi = 0x55edf40f20d0
0x000055edf3f58718 <+136>: movslq(%rdi,%r10,4),%rsi # %rsi = mem(%rdi+(%r10*4)) | %r10 = 0
0x000055edf3f5871c <+140>: add %rdi,%rsi # %rsi = %rsi + %rdi | %rsi = 0x55edf40f20d0
0x000055edf3f5871f <+143>: jmp *%rsi # jump to %rsi | segfault
0x000055edf3f58721 <+145>: nopl 0x0(%rax)
```

# segfault, 2024

```
0x000055edf3f5870b <+123>: cmp $0xb,%r10d
0x000055edf3f5870f <+127>: ja 0x55edf3f586b0 <hash_bytes+32>
0x000055edf3f58711 <+129>: lea 0x1999b8(%rip),%rdi # %rdi = (%rip+0x1999b8) | %rdi = 0x55edf40f20d0
0x000055edf3f58718 <+136>: movslq (%rdi,%r10,4),%rsi # %rsi = mem(%rdi+(%r10*4)) | %r10 = 0
0x000055edf3f5871c <+140>: add %rdi,%rsi # %rsi = %rsi + %rdi | %rsi = 0x55edf40f20d0
0x000055edf3f5871f <+143>: jmp *%rsi # jump to %rsi | segfault
0x000055edf3f58721 <+145>: nopl 0x0(%rax)
```

# segfault, 2024

```
0x000055edf3f5870b <+123>: cmp $0xb,%r10d
0x000055edf3f5870f <+127>: ja 0x55edf3f586b0 <hash_bytes+32>
0x000055edf3f58711 <+129>: lea 0x1999b8(%rip),%rdi # %rdi = (%rip+0x1999b8) | %rdi = 0x55edf40f20d0
0x000055edf3f58718 <+136>: movslq (%rdi,%r10,4),%rsi # %rsi = mem(%rdi+(%r10*4)) | %r10 = 0
0x000055edf3f5871c <+140>: add %rdi,%rsi # %rsi = %rsi + %rdi | %rsi = 0x55edf40f20d0
0x000055edf3f5871f <+143>: jmp *%rsi # jump to %rsi | segfault
0x000055edf3f58721 <+145>: nopl 0x0(%rax)
```

# segfault, 2024

```
0x000055edf3f5870b <+123>: cmp $0xb,%r10d
0x000055edf3f5870f <+127>: ja 0x55edf3f586b0 <hash_bytes+32>
0x000055edf3f58711 <+129>: lea 0x1999b8(%rip),%rdi # %rdi = (%rip+0x1999b8) | %rdi = 0x55edf40f20d0
0x000055edf3f58718 <+136>: movslq(%rdi,%r10,4),%rsi # %rsi = mem(%rdi+(%r10*4)) | %r10 = 0
0x000055edf3f5871c <+140>: add %rdi,%rsi # %rsi = %rsi + %rdi | %rsi = 0x55edf40f20d0
0x000055edf3f5871f <+143>: jmp *%rsi # jump to %rsi | segfault
0x000055edf3f58721 <+145>: nopl 0x0(%rax)
```

# segfault, 2024

```
0x000055edf3f5870b <+123>: cmp $0xb,%r10d
0x000055edf3f5870f <+127>: ja 0x55edf3f586b0 <hash_bytes+32>
0x000055edf3f58711 <+129>: lea 0x1999b8(%rip),%rdi # %rdi = (%rip+0x1999b8) | %rdi = 0x55edf40f20d0
0x000055edf3f58718 <+136>: movslq (%rdi,%r10,4),%rsi # %rsi = mem(%rdi+(%r10*4)) | %r10 = 0
0x000055edf3f5871c <+140>: add %rdi,%rsi # %rsi = %rsi + %rdi | %rsi = 0x55edf40f20d0
0x000055edf3f5871f <+143>: jmp *%rsi # jump to %rsi | segfault
0x000055edf3f58721 <+145>: nopl 0x0(%rax)
```

# segfault, 2024

```
0x000055edf3f5870b <+123>: cmp $0xb,%r10d
0x000055edf3f5870f <+127>: ja 0x55edf3f586b0 <hash_bytes+32>
0x000055edf3f58711 <+129>: lea 0x1999b8(%rip),%rdi # %rdi = (%rip+0x1999b8) | %rdi = 0x55edf40f20d0
0x000055edf3f58718 <+136>: movslq (%rdi,%r10,4),%rsi # %rsi = mem(%rdi+(%r10*4)) | %r10 = 0
0x000055edf3f5871c <+140>: add %rdi,%rsi # %rsi = %rsi + %rdi | %rsi = 0x55edf40f20d0
0x000055edf3f5871f <+143>: jmp *%rsi # jump to %rsi | segfault
0x000055edf3f58721 <+145>: nopl 0x0(%rax)
```

# segfault, 2024

```
0x000055edf3f5870b <+123>: cmp $0xb,%r10d
0x000055edf3f5870f <+127>: ja 0x55edf3f586b0 <hash_bytes+32>
0x000055edf3f58711 <+129>: lea 0x1999b8(%rip),%rdi # %rdi = (%rip+0x1999b8) | %rdi = 0x55edf40f20d0
0x000055edf3f58718 <+136>: movslq(%rdi,%r10,4),%rsi # %rsi = mem(%rdi+(%r10*4)) | %r10 = 0
0x000055edf3f5871c <+140>: add %rdi,%rsi # %rsi = %rsi + %rdi | %rsi = 0x55edf40f20d0
0x000055edf3f5871f <+143>: jmp *%rsi # jump to %rsi | segfault
0x000055edf3f58721 <+145>: nopl 0x0(%rax)
```

# segfault, 2024

- Процессор зачитал значение равное **0** из адреса `%rdi+%r10*4`
- Должен был зачитать **0xffe667a7**

```
(gdb) p/x *(unsigned*)0x55edf40f20d0@12
$9 = {0xffe667a7, ... }
```

# segfault, 2024

- Процессор зачитал значение равное **0** из адреса `%rdi+%r10*4`
- Должен был зачитать **0xffe667a7**

```
(gdb) p/x *(unsigned*)0x55edf40f20d0@12
$9 = {0xffe667a7, ... }
```

# segfault, 2024

- 31 декабря, 2:30 ночи: документ с анализом



# segfault, 2024

- 31 декабря, 2:30 ночи: документ с анализом
- "Главное - мы видим аномалию в работе системы, по какой-то причине процессор зачитал нулевое значения из ненулевой памяти (которая постоянная, никогда не меняется)"



**Конечно никто не собирался  
нам поверить**

# segfault, 2024

- Документ ушёл в ТП вендора ОС



# segfault, 2024

- Документ ушёл в ТП вендора ОС
- Ответ



# segfault, 2024

- Документ ушёл в ТП вендора ОС
- Ответ
  - **"Причины SIGSEGV во всех случаях (внезапно!) не внешние, а со стороны userspace, вызванные процессом postgres"**

# segfault, 2024

- Документ ушёл в ТП вендора ОС
- Ответ
  - **"Причины SIGSEGV во всех случаях (внезапно!) не внешние, а со стороны userspace, вызванные процессом postgres"**
- Так заканчивался 2024 год...



segfault, 2025

# segfault, 2025

- Новые падения (6 января и т.д.) - случайный места, случайные программы (sshd)



# segfault, 2025

- Новые падения (6 января и т.д.) - случайный места, случайные программы (sshd)
- Больше доказательство о нулевой памяти



# segfault, 2025

- Новые падения (6 января и т.д.) - случайный места, случайные программы (sshd)
- Больше доказательство о нулевой памяти
- Гипотезы рождались как грибы

# Гипотезы

- Библиотеки сжатия



# Гипотезы

- Библиотеки сжатия
- Программные raid контроллеры



# Гипотезы

- Библиотеки сжатия
- Программные raid контроллеры
- Прошивки и микрокоды



# Гипотезы

- Библиотеки сжатия
- Программные raid контроллеры
- Прошивки и микрокоды
- FPU



# Гипотезы

- Библиотеки сжатия
- Программные raid контроллеры
- Прошивки и микрокоды
- FPU
- Гипотез много, фактов ноль



# Гипотезы

- Библиотеки сжатия
- Программные raid контроллеры
- Прошивки и микрокоды
- FPU
- Гипотез много, фактов ноль
- **"It's a capital mistake to theorize without facts"** (с) Артур Конан Дойль

# segfault, 2025

- Падения продолжались (28 января, 7 февраля)



# segfault, 2025

- Падения продолжались (28 января, 7 февраля)
- **rasdaemon**: RAS (Reliability, Availability and Serviceability) logging tool

# segfault, 2025

- Падения продолжались (28 января, 7 февраля)
- **rasdaemon**: RAS (Reliability, Availability and Serviceability) logging tool
- Тук-тук, у вас битые планки памяти!

# segfault, 2025

- Падения продолжались (28 января, 7 февраля)
- **rasdaemon**: RAS (Reliability, Availability and Serviceability) logging tool
- Тук-тук, у вас битые планки памяти!
- Эх, опять мимо

# segfault, 2025

- Падения продолжались (28 января, 7 февраля)
- **rasdaemon**: RAS (Reliability, Availability and Serviceability) logging tool
- Тук-тук, у вас битые планки памяти!
- Эх, опять мимо
- Падали разные сервера



# segfault, 2025

- Падения продолжались (28 января, 7 февраля)
- **rasdaemon**: RAS (Reliability, Availability and Serviceability) logging tool
- Тук-тук, у вас битые планки памяти!
- Эх, опять мимо
- Падали разные сервера
- Делали анализ, писали отчёты, всегда одно и тоже - **обнуление памяти!**



**Сложить  
лапки**



**Поигратся  
со стендом**

# segfault, 2025

- Стенд от производителя железа



# segfault, 2025

- Стенд от производителя железа
- **Спасибо** за оперативность



# segfault, 2025

- Стенд от производителя железа
- **Спасибо** за оперативность
- Одно из падений: перестроение индексов в базе данных (rebuild index)



# segfault, 2025

- Стенд от производителя железа
- **Спасибо** за оперативность
- Одно из падений: перестроение индексов в базе данных (rebuild index)
- Делаем большую синтетическую базу



# segfault, 2025

- Стенд от производителя железа
- **Спасибо** за оперативность
- Одно из падений: перестроение индексов в базе данных (rebuild index)
- Делаем большую синтетическую базу
- Добавляем проверку памяти (memtester)



## ■ 25 февраля

```
reindexdb -j 400
memtester 680G 1000
```

bash упал 11:18:44 ✓

Что сделал:

1) поднял базу которая лежала в /MBD\_DATA/pgdata. Там 100GB нужно было в huge pages, их туда добавил:

```
shell
echo 51543 > /proc/sys/vm/nr_hugepages
/opt/pgpro/ent-15/bin/pg_ctl -D /MBD_DATA/pgdata start
```

2) На базе поставил maintenance\_work\_mem в 96GB:

```
sql
alter system set maintenance_work_mem to '96GB';
```

3) В параллель запустил:

```
copy
reindexdb -j 400
```

и

```
copy
memtester 680G 1000
```

Прошло буквально несколько секунд-минут и reindex отвалился. В этот момент всё остановил и увидел в dmesg что был crash самого bash

11:22:07 ✓

- 25 февраля

```
reindexdb -j 400
memtester 680G 1000
```

- 400 процессов перестроения индексов

```
bash упал 11:18:44 ✓
```

Что сделал:

- 1) поднял базу которая лежала в /MBD\_DATA/pgdata. Там 100GB нужно было в huge pages, их туда добавил:

```
shell
echo 51543 > /proc/sys/vm/nr_hugepages
/opt/pgpro/ent-15/bin/pg_ctl -D /MBD_DATA/pgdata start
```

- 2) На базе поставил maintenance\_work\_mem в 96GB:

```
sql
alter system set maintenance_work_mem to '96GB';
```

- 3) В параллель запустил:

```
copy
reindexdb -j 400
```

и

```
copy
memtester 680G 1000
```

Прошло буквально несколько секунд-минут и reindex отвалился. В этот момент всё остановил и увидел в dmesg что был crash самого bash

```
11:22:07 ✓
```

- 25 февраля

```
reindexdb -j 400
memtester 680G 1000
```

- 400 процессов перестроения индексов
- 1000 итераций проверки 680 гигабайт памяти

bash упал 11:18:44 ✓

Что сделал:

1) поднял базу которая лежала в /MBD\_DATA/pgdata. Там 100GB нужно было в huge pages, их туда добавил:

```
shell
echo 51543 > /proc/sys/vm/nr_hugepages
/opt/pgpro/ent-15/bin/pg_ctl -D /MBD_DATA/pgdata start
```

2) На базе поставил maintenance\_work\_mem в 96GB:

```
sql
alter system set maintenance_work_mem to '96GB';
```

3) В параллель запустил:

```
copy
reindexdb -j 400
```

и

```
copy
memtester 680G 1000
```

Прошло буквально несколько секунд-минут и reindex отвалился. В этот момент всё остановил и увидел в dmesg что был crash самого bash

11:22:07 ✓

- 25 февраля

```
reindexdb -j 400
memtester 680G 1000
```

- 400 процессов перестроения индексов
- 1000 итераций проверки 680 гигабайт памяти
- Упал не postgres, а bash

bash упал 11:18:44 ✓

Что сделал:

1) поднял базу которая лежала в /MBD\_DATA/pgdata. Там 100GB нужно было в huge pages, их туда добавил:

```
shell
echo 51543 > /proc/sys/vm/nr_hugepages
/opt/pgpro/ent-15/bin/pg_ctl -D /MBD_DATA/pgdata start
```

2) На базе поставил maintenance\_work\_mem в 96GB:

```
sql
alter system set maintenance_work_mem to '96GB';
```

3) В параллель запустил:

```
copy
reindexdb -j 400
```

и

```
copy
memtester 680G 1000
```

Прошло буквально несколько секунд-минут и reindex отвалился. В этот момент всё остановил и увидел в dmesg что был crash самого bash

11:22:07 ✓

- 25 февраля

```
reindexdb -j 400
memtester 680G 1000
```

- 400 процессов перестроения индексов
- 1000 итераций проверки 680 гигабайт памяти
- Упал не postgres, а bash
- Воспроизводимый кейс

bash упал 11:18:44 ✓

Что сделал:

1) поднял базу которая лежала в /MBD\_DATA/pgdata. Там 100GB нужно было в huge pages, их туда добавил:

```
shell
echo 51543 > /proc/sys/vm/nr_hugepages
/opt/pgpro/ent-15/bin/pg_ctl -D /MBD_DATA/pgdata start
```

2) На базе поставил maintenance\_work\_mem в 96GB:

```
sql
alter system set maintenance_work_mem to '96GB';
```

3) В параллель запустил:

```
copy
reindexdb -j 400
```

и

```
copy
memtester 680G 1000
```

Прошло буквально несколько секунд-минут и reindex отвалился. В этот момент всё остановил и увидел в dmesg что был crash самого bash

11:22:07 ✓

- 25 февраля

```
reindexdb -j 400
memtester 680G 1000
```

- 400 процессов перестроения индексов
- 1000 итераций проверки 680 гигабайт памяти
- Упал не postgres, а bash
- Воспроизводимый кейс
- Маленькая победа!

bash упал 11:18:44 ✓

Что сделал:

- 1) поднял базу которая лежала в /MBD\_DATA/pgdata. Там 100GB нужно было в huge pages, их туда добавил:  

```
shell
echo 51543 > /proc/sys/vm/nr_hugepages
/opt/pgpro/ent-15/bin/pg_ctl -D /MBD_DATA/pgdata start
```
- 2) На базе поставил maintenance\_work\_mem в 96GB:  

```
sql
alter system set maintenance_work_mem to '96GB';
```
- 3) В параллель запустил:  

```
copy
reindexdb -j 400
```

и  

```
copy
memtester 680G 1000
```

Прошло буквально несколько секунд-минут и reindex отвалился. В этот момент всё остановил и увидел в dmesg что был crash самого bash

11:22:07 ✓

- 25 февраля

```
reindexdb -j 400
memtester 680G 1000
```

- 400 процессов перестроения индексов
- 1000 итераций проверки 680 гигабайт памяти
- Упал не postgres, а bash
- Воспроизводимый кейс
- Маленькая победа!
- Но что дальше?

bash упал 11:18:44 ✓

Что сделал:

- 1) поднял базу которая лежала в /MBD\_DATA/pgdata. Там 100GB нужно было в huge pages, их туда добавил:

```
shell
echo 51543 > /proc/sys/vm/nr_hugepages
/opt/pgpro/ent-15/bin/pg_ctl -D /MBD_DATA/pgdata start
```
- 2) На базе поставил maintenance\_work\_mem в 96GB:

```
sql
alter system set maintenance_work_mem to '96GB';
```
- 3) В параллель запустил:

```
copy
reindexdb -j 400
```

и

```
copy
memtester 680G 1000
```

Прошло буквально несколько секунд-минут и reindex отвалился. В этот момент всё остановил и увидел в dmesg что был crash самого bash

11:22:07 ✓

# segfault, 2025

- Инструкция для проверки нового оборудования



# segfault, 2025

- Инструкция для проверки нового оборудования
- Глобальное обновление железа (Апрель-Май)



# segfault, 2025

- Инструкция для проверки нового оборудования
- Глобальное обновление железа (Апрель-Май)
  - Смена ОС



# segfault, 2025

- Инструкция для проверки нового оборудования
- Глобальное обновление железа (Апрель-Май)
  - Смена ОС
  - Увеличение памяти с 1TiB до 4TiB



# segfault, 2025

- Инструкция для проверки нового оборудования
- Глобальное обновление железа (Апрель-Май)
  - Смена ОС
  - Увеличение памяти с 1TiB до 4TiB
  - Замена software RAID на hardware RAID



# segfault, 2025

- Инструкция для проверки нового оборудования
- Глобальное обновление железа (Апрель-Май)
  - Смена ОС
  - Увеличение памяти с 1TiB до 4TiB
  - Замена software RAID на hardware RAID
- Ура! Никаких больше падений



# segfault, 2025

- Инструкция для проверки нового оборудования
- Глобальное обновление железа (Апрель-Май)
  - Смена ОС
  - Увеличение памяти с 1TiB до 4TiB
  - Замена software RAID на hardware RAID
- Ура! Никаких больше падений
- Успешный выход проекта в мае-июне



# segfault, 2025

- Инструкция для проверки нового оборудования
- Глобальное обновление железа (Апрель-Май)
  - Смена ОС
  - Увеличение памяти с 1TiB до 4TiB
  - Замена software RAID на hardware RAID
- Ура! Никаких больше падений
- Успешный выход проекта в мае-июне
- **Увы, никто тогда не знал что причина была не в железе...**



КОНЕЦ

"Конец, конец... Концы в воду!"



**"Иди, и без ёлки не возвращайся"**

# Осень 2025, снова жалобы

- Различные клиенты



# Осень 2025, снова жалобы

- Различные клиенты
- Различное железо



# Осень 2025, снова жалобы

- Различные клиенты
- Различное железо
- Различные операционные системы



# Осень 2025, снова жалобы

- Различные клиенты
- Различное железо
- Различные операционные системы
- И не только PostgresPro, но PostgreSQL, Kafka и прочие



# Осень 2025, снова жалобы

- Различные клиенты
- Различное железо
- Различные операционные системы
- И не только PostgresPro, но PostgreSQL, Kafka и прочие
- **А точно ли дефект был в железе? Ответ: нет**



# Осень 2025, воспроизведение

- Linux kernel 6.1, XFS, 2 TiB RAM



# Осень 2025, воспроизведение

- Linux kernel 6.1, XFS, 2 TiB RAM
- Неожиданно воспроизвелось на наших серверах



# Осень 2025, воспроизведение

- Linux kernel 6.1, XFS, 2 TiB RAM
- Неожиданно воспроизвелось на наших серверах
- Подкрутили тест



# Осень 2025, воспроизведение

- Linux kernel 6.1, XFS, 2 TiB RAM
- Неожиданно воспроизвелось на наших серверах
- Подкрутили тест
- Воспроизвели на виртуалках с большим количеством vCPU и памяти




# Осень 2025, воспроизведение

- Linux kernel 6.1, XFS, 2 TiB RAM
- Неожиданно воспроизвелось на наших серверах
- Подкрутили тест
- Воспроизвели на виртуалках с большим количеством vCPU и памяти
- Ещё подкрутили тест



# Осень 2025, воспроизведение

- Linux kernel 6.1, XFS, 2 TiB RAM
  - Неожиданно воспроизвелось на наших серверах
  - Подкрутили тест
  - Воспроизвели на виртуалках с большим количеством vCPU и памяти
  - Ещё подкрутили тест
  - Воспроизвели на маленьких виртуалках (4-8 vCPU)
- 

# Осень 2025, воспроизведение

- Супер! Теперь можем что-нибудь поменять



# Осень 2025, воспроизведение

- Супер! Теперь можем что-нибудь поменять
- ext4 вместо XFS - **не** воспроизводится



# Осень 2025, воспроизведение

- Супер! Теперь можем что-нибудь поменять
- ext4 вместо XFS - **не** воспроизводится
- Debian вместо российской ОС - **не** воспроизводится



# Осень 2025, воспроизведение

- Супер! Теперь можем что-нибудь поменять
- ext4 вместо XFS - **не** воспроизводится
- Debian вместо российской ОС - **не** воспроизводится
- Магия!



# Осень 2025, воспроизведение

- Супер! Теперь можем что-нибудь поменять
- ext4 вместо XFS - **не** воспроизводится
- Debian вместо российской ОС - **не** воспроизводится
- Магия!
- Замена Linux kernel с вендорского на ванильное - **воспроизводится**

# Осень 2025, воспроизведение

- Давайте менять версии Linux kernel



# Осень 2025, воспроизведение

- Давайте менять версии Linux kernel
- 5.18 – 6.8

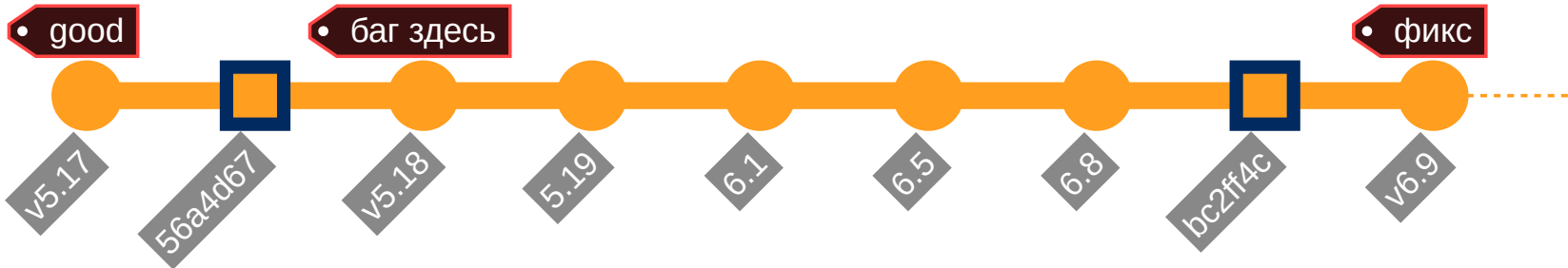


# Осень 2025, воспроизведение

- Давайте менять версии Linux kernel
- 5.18 – 6.8
- Поиск git коммита, который добавил дефект



main



# Осень 2025, воспроизведение

- Перебор коммитов



# Осень 2025, воспроизведение

- Перебор коммитов
- Проблемный коммит в драйвере клавиатуры!



# Осень 2025, воспроизведение

- Перебор коммитов
- Проблемный коммит в драйвере клавиатуры!
- Чё-то не то...




# Осень 2025, воспроизведение

- Перебор коммитов
- Проблемный коммит в драйвере клавиатуры!
- Чё-то не то...
- GitHub vs git-bisect



# GITHUB

drivers: net: amd: lance: Remove this driver 


 lunn authored and kuba-moo committed 2 days ago

2fbd04d  

drivers: net: 3com: 3c589: Remove this driver 

 lunn authored and kuba-moo committed 2 days ago

4ff8d06  

drivers: net: 3com: 3c574: Remove this driver 

 lunn authored and kuba-moo committed 2 days ago

a7fbf27  

drivers: net: 3com: 3c515: Remove this driver 


 lunn authored and kuba-moo committed 2 days ago

082b2e0  

drivers: net: 3com: 3c509: Remove this driver 

 lunn authored and kuba-moo committed 2 days ago

91f3a27  

Merge patch series "eventpoll: fix ep\_remove() UAF and follow-up cleanup" 

 brauner committed 2 days ago

Verified

ac8777c  

eventpoll: drop vestigial epi->dying flag 

 brauner committed 2 days ago

Verified

07422c9  

eventpoll: drop dead bool return from ep\_remove\_epi() 

 brauner committed 2 days ago

Verified

3a4551e  

eventpoll: refresh eventpoll\_release() fast-path comment 

 brauner committed 2 days ago

Verified

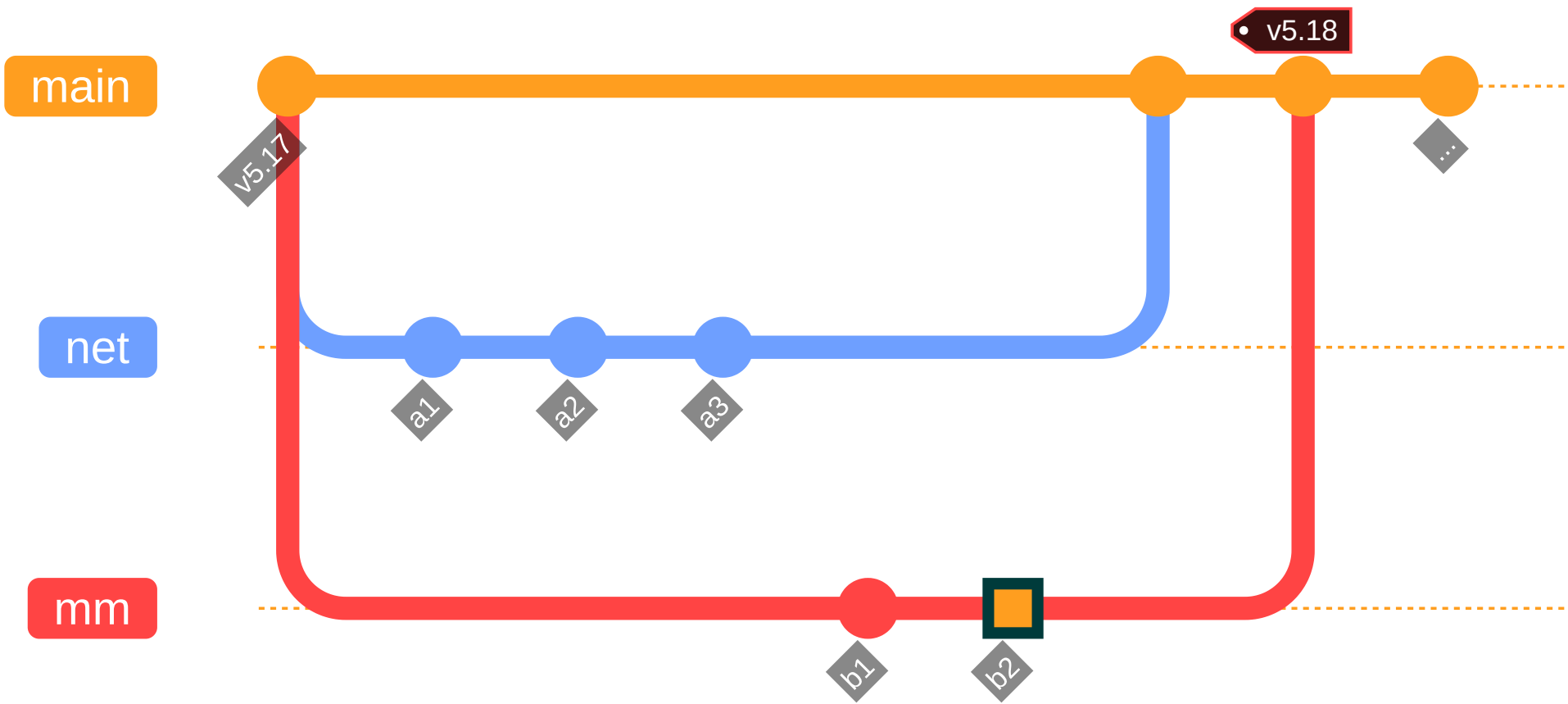
33e92e9  

eventpoll: move f\_lock acquisition into ep\_remove\_file() 

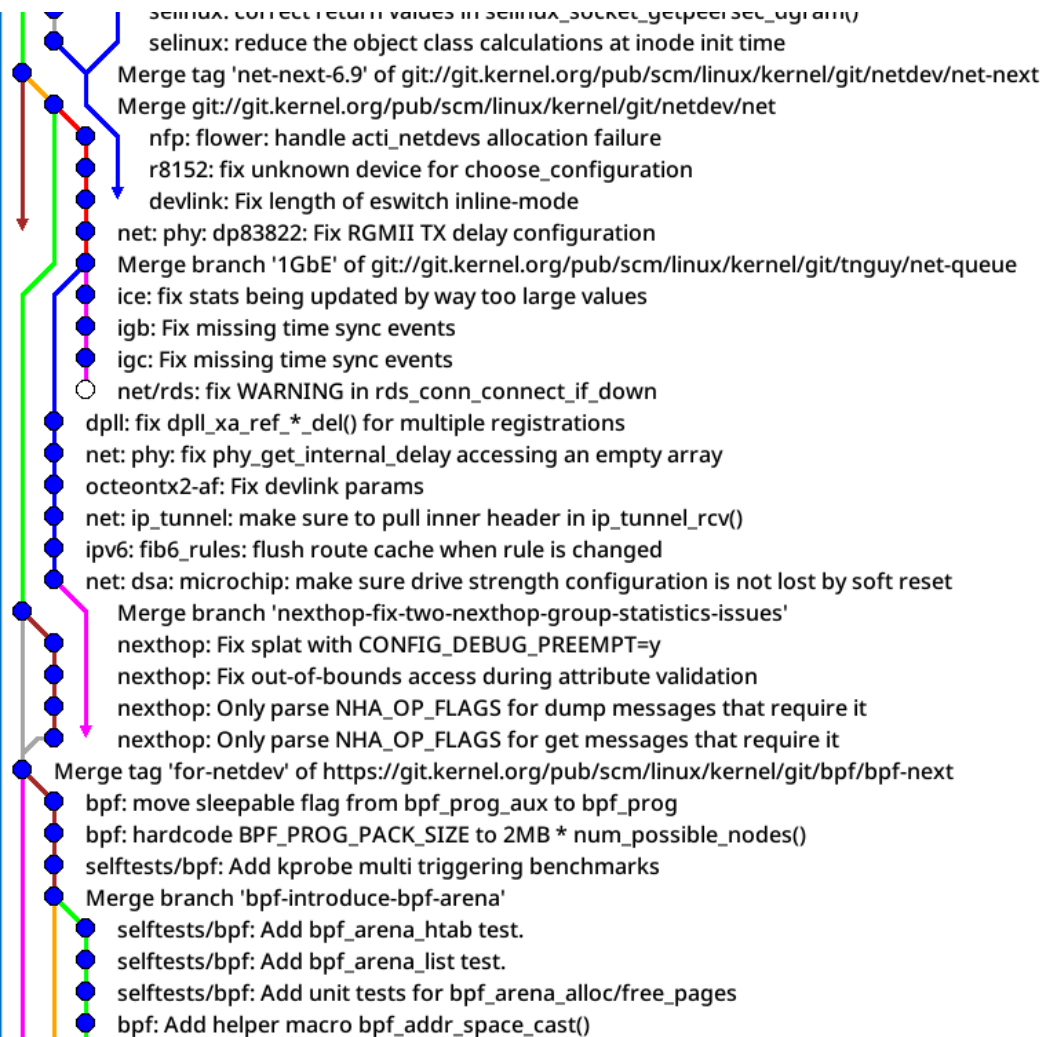
 brauner committed 2 days ago

Verified

d30deeb  







|                                      |                 |
|--------------------------------------|-----------------|
| Paul Moore <paul@paul-moore.com>     | 2024-01-19 00:4 |
| Paul Moore <paul@paul-moore.com>     | 2024-01-19 04:4 |
| Linus Torvalds <torvalds@linux-fou   | 2024-03-13 03:4 |
| Jakub Kicinski <kuba@kernel.org>     | 2024-03-12 06:3 |
| Duoming Zhou <duoming@zju.edu.c      | 2024-03-08 17:2 |
| Hayes Wang <hayeswang@realtek.c      | 2024-03-08 10:5 |
| William Tu <witu@nvidia.com>         | 2024-03-10 19:4 |
| Tim Pambor <tp@osasysteme.de>        | 2024-03-05 14:0 |
| Jakub Kicinski <kuba@kernel.org>     | 2024-03-08 22:3 |
| Przemek Kitszel <przemyslaw.kitsze   | 2024-02-27 17:3 |
| Vinicius Costa Gomes <vinicius.gom   | 2024-02-21 02:5 |
| Vinicius Costa Gomes <vinicius.gom   | 2024-02-21 02:5 |
| Edward Adam Davis <eadavis@qq.c      | 2024-03-05 03:1 |
| Jiri Pirko <jiri@nvidia.com>         | 2024-03-06 18:1 |
| Kévin L'hôpital <kevin.lhopital@savc | 2024-03-07 14:1 |
| Sunil Goutham <sgoutham@marvel       | 2024-03-07 13:5 |
| Eric Dumazet <edumazet@google.c      | 2024-03-07 13:0 |
| Shiming Cheng <shiming.cheng@m       | 2024-03-07 13:0 |
| Oleksij Rempel <o.rempel@pengutr     | 2024-03-04 16:5 |
| Jakub Kicinski <kuba@kernel.org>     | 2024-03-12 06:3 |
| Ido Schimmel <idosch@nvidia.com>     | 2024-03-11 19:2 |
| Ido Schimmel <idosch@nvidia.com>     | 2024-03-11 19:2 |
| Ido Schimmel <idosch@nvidia.com>     | 2024-03-11 19:2 |
| Ido Schimmel <idosch@nvidia.com>     | 2024-03-11 19:2 |
| Jakub Kicinski <kuba@kernel.org>     | 2024-03-12 04:0 |
| Andrii Nakryiko <andrii@kernel.org   | 2024-03-09 03:4 |
| Puranjay Mohan <puranjay12@gma       | 2024-03-11 15:2 |
| Jiri Olsa <jolsa@kernel.org>         | 2024-03-12 00:1 |
| Andrii Nakryiko <andrii@kernel.org   | 2024-03-12 01:3 |
| Alexei Starovoitov <ast@kernel.org   | 2024-03-08 04:0 |
| Alexei Starovoitov <ast@kernel.org   | 2024-03-08 04:0 |
| Alexei Starovoitov <ast@kernel.org   | 2024-03-08 04:0 |
| Alexei Starovoitov <ast@kernel.org   | 2024-03-08 04:0 |

# segfault, 2025

- git bisect: поиск бага по графу

```
git bisect start
git bisect good v5.17 # бага нет
git bisect bad v5.18 # баг есть
git bisect run ./test.sh # → 56a4d67c264e
```

# segfault, 2025

- Коммит, который сломал (21 марта, 2022): `mm/readahead: Switch to page_cache_ra_order`
  - Масштаб: 3 файла, +5 / -5 строк

# segfault, 2025

- Коммит, который сломал (21 марта, 2022): mm/readahead: Switch to page\_cache\_ra\_order
  - Масштаб: 3 файла, +5 / -5 строк
- Коммит, который исправил (5 марта, 2024): mm: free folios in a batch in shrink\_folio\_list()
  - Масштаб: 1 файла, +9 / -11 строк

# segfault, 2025

- Коммит, который сломал (21 марта, 2022): mm/readahead: Switch to page\_cache\_ra\_order
  - Масштаб: 3 файла, +5 / -5 строк
- Коммит, который исправил (5 марта, 2024): mm: free folios in a batch in shrink\_folio\_list()
  - Масштаб: 1 файла, +9 / -11 строк
- Автор обоих коммитов: Matthew Wilcox, Oracle



# segfault, 2025

- Коммит, который сломал (21 марта, 2022): mm/readahead: Switch to page\_cache\_ra\_order
  - Масштаб: 3 файла, +5 / -5 строк
- Коммит, который исправил (5 марта, 2024): mm: free folios in a batch in shrink\_folio\_list()
  - Масштаб: 1 файла, +9 / -11 строк
- Автор обоих коммитов: Matthew Wilcox, Oracle
- Знаем коммиты, хотим ли мы знать причину? Ответ: **не хотим**

# segfault, 2025

- Коммит, который сломал (21 марта, 2022): mm/readahead: Switch to page\_cache\_ra\_order
  - Масштаб: 3 файла, +5 / -5 строк
- Коммит, который исправил (5 марта, 2024): mm: free folios in a batch in shrink\_folio\_list()
  - Масштаб: 1 файла, +9 / -11 строк
- Автор обоих коммитов: Matthew Wilcox, Oracle
- Знаем коммиты, хотим ли мы знать причину? Ответ: **не хотим**
- Знание по ядру есть на уровне исследователя, не программиста

# segfault, 2025

- Но хотим понять, почему на Debian не воспроизводится!



# segfault, 2025

- Но хотим понять, почему на Debian не воспроизводится!
- Ядра одинаковые, но есть и параметры загрузки ядра. Grub



# segfault, 2025

- Но хотим понять, почему на Debian не воспроизводится!
- Ядра одинаковые, но есть и параметры загрузки ядра. Grub
- `init_on_free = 0` во всех зарубежных



# segfault, 2025

- Но хотим понять, почему на Debian не воспроизводится!
- Ядра одинаковые, но есть и параметры загрузки ядра. Grub
- `init_on_free = 0` во всех зарубежных
- `init_on_free = 1` в российских



# Проверьте свои серверы

- Ядро 5.18 – 6.8 + XFS + init\_on\_free = обратитесь к вендору ОС
- Фикс в 6.9, но не помечен как исправление
- Бэкапорта нет: на LTS-ядрах фикс сам не придёт
- Нужен bug report в bugzilla
- А теперь прикол...

segfault или не segfault?

# segfault или не segfault?

- SIGSEGV - это сигнал
- Сигнал можно обработать
- Пишем свой обработчик



# segfault или не segfault?

```
sigaction_t sigact;

sigact.sa_sigaction = my_segfault_handler;
sigact.sa_flags = SA_RESETHAND | ...;

...

sigaction(SIGSEGV, &sigact, NULL);
```

# segfault или не segfault?

```
sigaction_t sigact;

sigact.sa_sigaction = my_segfault_handler;
sigact.sa_flags = SA_RESETHAND | ... ;

...

sigaction(SIGSEGV, &sigact, NULL);
```

# segfault или не segfault?

```
sigaction_t sigact;

sigact.sa_sigaction = my_segfault_handler;
sigact.sa_flags = SA_RESETHAND | ... ;

...

sigaction(SIGSEGV, &sigact, NULL);
```

# segfault или не segfault?

```
sigaction_t sigact;

sigact.sa_sigaction = my_segfault_handler;
sigact.sa_flags = SA_RESETHAND | ...;

...

sigaction(SIGSEGV, &sigact, NULL);
```

# segfault или не segfault?

- Ночью создавался индекс



# segfault или не segfault?

- Ночью создавался индекс
- С утра DBA находит 3 файла от нашего обработчика
  - `crash_1764605313879686_pid2130208.log`
  - `crash_1764605308648131_pid1867399.log`
  - `crash_1764605308641242_pid2130207.log`



# segfault или не segfault?

- Ночью создавался индекс
- С утра DBA находит 3 файла от нашего обработчика
  - `crash_1764605313879686_pid2130208.log`
  - `crash_1764605308648131_pid1867399.log`
  - `crash_1764605308641242_pid2130207.log`
- Упало создание индексов
  - `program_invocation_name: postgres: mydb: postgres myuser local CREATE INDEX`

# segfault или не segfault?

- ВНИМАНИЕ! Никто не пострадал



# segfault или не segfault?

- ВНИМАНИЕ! Никто не пострадал
- Индекс создан



# segfault или не segfault?

- ВНИМАНИЕ! Никто не пострадал
- Индекс создан
- В логах нет ошибок



# segfault или не segfault?

- ВНИМАНИЕ! Никто не пострадал
- Индекс создан
- В логах нет ошибок
- Процессы не падали, база не перезагружалась!

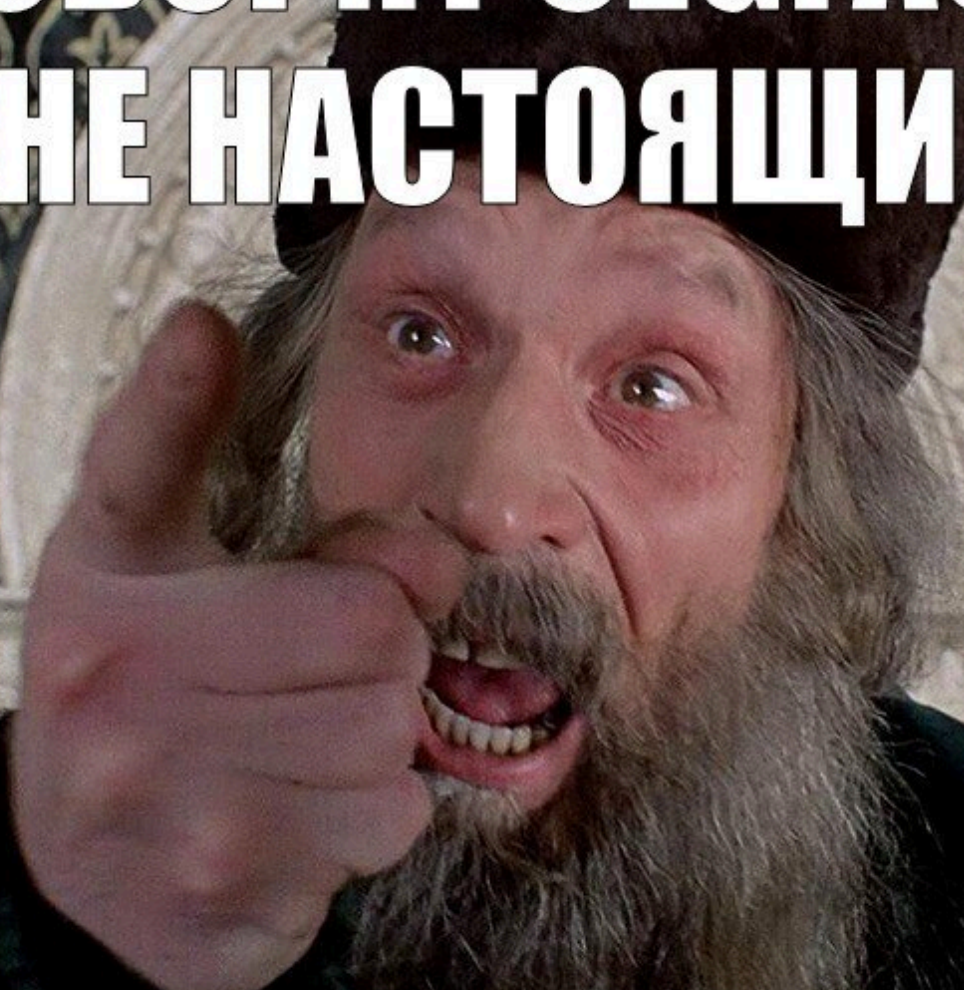


# segfault или не segfault?

- ВНИМАНИЕ! Никто не пострадал
- Индекс создан
- В логах нет ошибок
- Процессы не падали, база не перезагружалась!
- Индекс валидный!



**ГОВОРЯТ SEGFAULT  
НЕ НАСТОЯЩИЙ!**



# Без обработчика

- Случается SEGFAULT



# Без обработчика

- Случается SEGFAULT
- Инструкция процессора по адресу записанному в регистр RIP



# Без обработчика

- Случается SEGFAULT
- Инструкция процессора по адресу записанному в регистр RIP
- Ядерный обработчик логирует, делает coredump и убивается процесс

# С обработчиком

- Случается SEGFAULT



# С обработчиком

- Случается SEGFAULT
- Вызывается наш обработчик



# С обработчиком

- Случается SEGFAULT
- Вызывается наш обработчик
- Делаем полезное действие




# С обработчиком

- Случается SEGFAULT
- Вызывается наш обработчик
- Делаем полезное действие
- Потом возвращается обработчик по умолчанию (благодаря SA\_RESETHAND)



# С обработчиком

- Случается SEGFAULT
- Вызывается наш обработчик
- Делаем полезное действие
- Потом возвращается обработчик по умолчанию (благодаря SA\_RESETHAND)
- Операционная система заканчивает обработку сигнала

A close-up photograph of a man and a baby lying together. The man is on top, his eyes closed in a peaceful sleep. The baby is below him, looking upwards with a calm expression. The lighting is soft and intimate, highlighting their faces. The man is wearing a grey t-shirt, and the baby is in a white garment. The background is dark, making the subjects stand out.

**ГОРОД ЗАСЫПАЕТ,  
ПРОСЫПАЕТСЯ  
МАФИЯ...**

# С обработчиком

- Повторяется инструкция процессора в тех же самых условиях



# С обработчиком

- Повторяется инструкция процессора в тех же самых условиях
- Значит, повторяется и SIGSEGV / SEGV



# С обработчиком

- Повторяется инструкция процессора в тех же самых условиях
- Значит, повторяется и SIGSEGV / SEGV
- А иногда и не повторяется



# С обработчиком

- Повторяется инструкция процессора в тех же самых условиях
- Значит, повторяется и SIGSEGV / SEGV
- А иногда и не повторяется
- Вторая попытка позже первой (разница до 1 секунды)



# С обработчиком

- Страница пустая только некоторый короткий интервал времени



# С обработчиком

- Страница пустая только некоторый короткий интервал времени
- Структуры памяти восстанавливаются



# С обработчиком

- Страница пустая только некоторый короткий интервал времени
- Структуры памяти восстанавливаются
- Ядерный обработчик не вызывается



# А почему раньше не помогало?

- Обнулялись страницы из сегментов отличных от .text
- "0x00 0x00" валидная инструкция

```
00 00 = add [rax], al
```

- Если RAX содержит невалидный адрес, то падает сразу
- Если RAX - валидный адрес, то... падает, но потом

Виртуальные страницы процесса

Физическая память (RAM)

Виртуальные страницы процесса

Page 0

Frame 0

Page 0

Page 1

Frame 1

Page 1

Page 2

Frame 2

Page 2

Page 3

Frame 3

Page 3

Диск (XFS)

Page 1 (на диске)

Page 2 (на диске)

Виртуальные страницы процесса

Физическая память (RAM)

Виртуальные страницы процесса

Page 0

Page 1

Page 2

Page 3

Frame 0

Frame 1 (free -> zeroing)

Frame 2

Frame 3

Page 0

Page 1

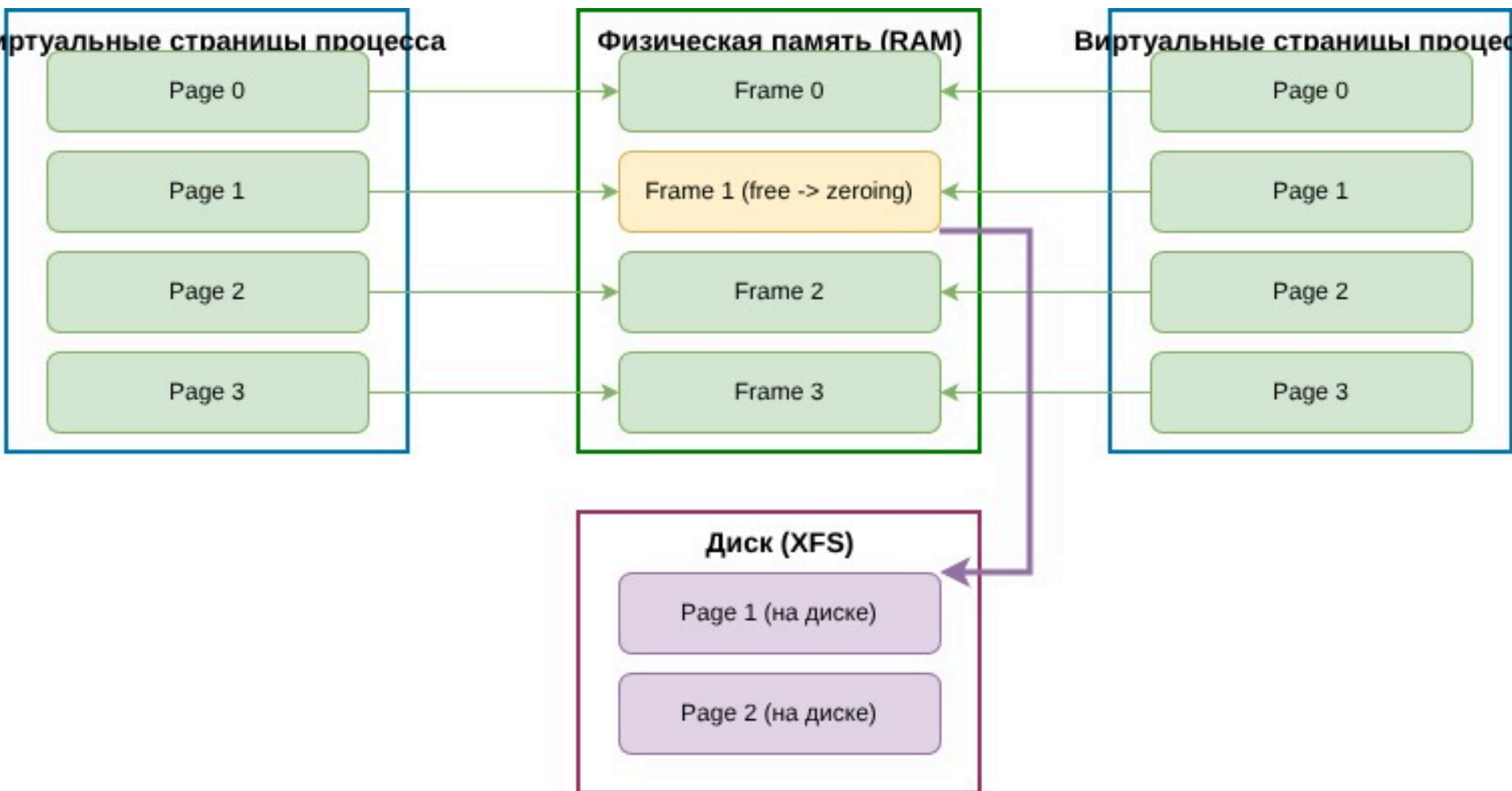
Page 2

Page 3

Диск (XFS)

Page 1 (на диске)

Page 2 (на диске)



Виртуальные страницы процесса

Page 0

Page 1

Page 2

Page 3

Физическая память (RAM)

Frame 0

PAGE FAULT

Frame 2

Frame 3

Виртуальные страницы процесса

Page 0

Page 1

Page 2

Page 3

Диск (XFS)

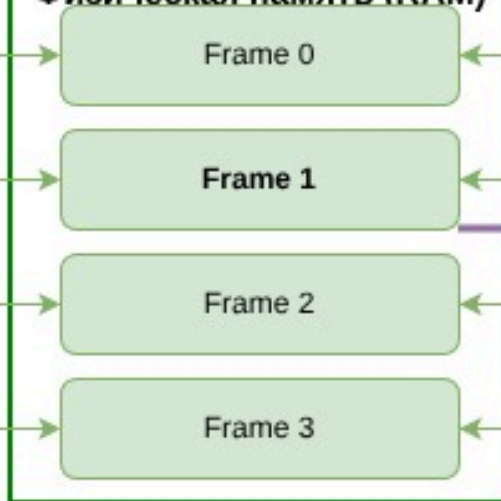
Page 1 (на диске)

Page 2 (на диске)

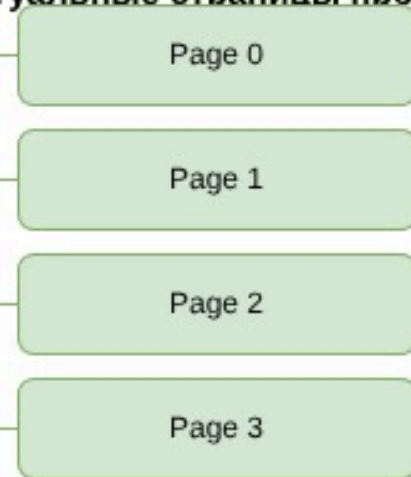
Виртуальные страницы процесса



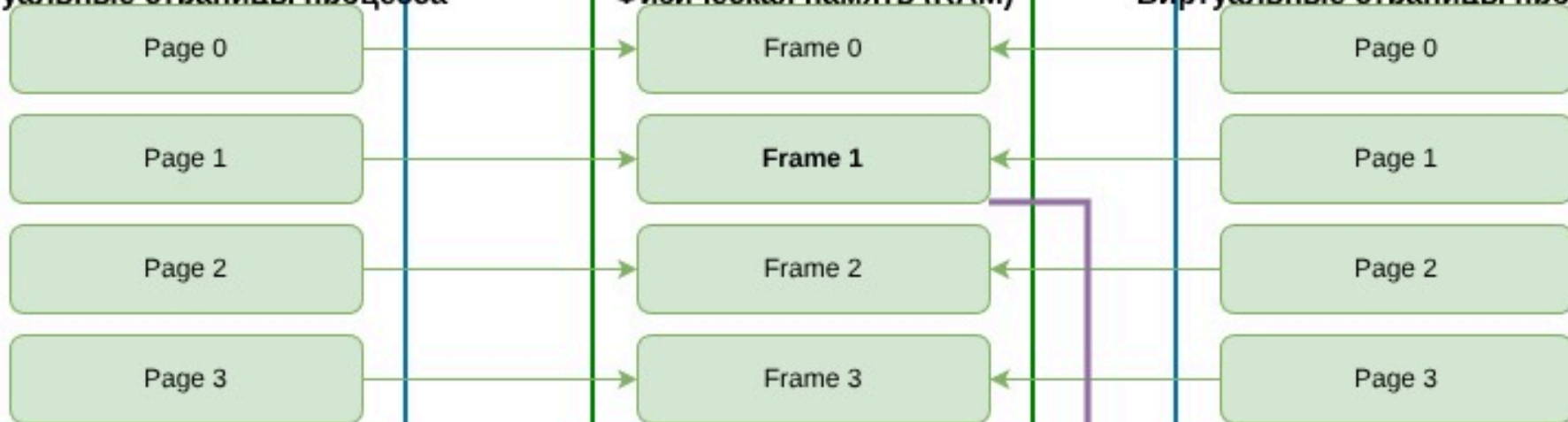
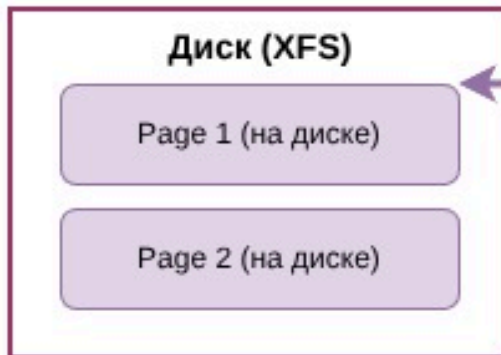
Физическая память (RAM)



Виртуальные страницы процесса



Диск (XFS)



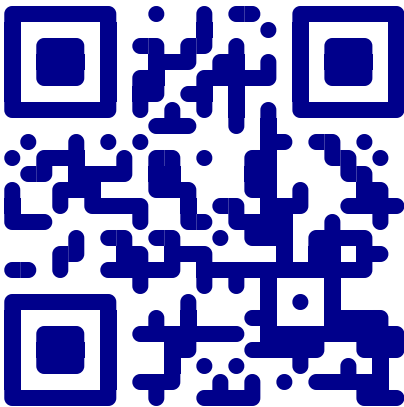
# Открытые вопросы

- Падаем до Page Fault или после?
- Дефект или недоделанная фича?
- Сможет ли AI агент помочь в поиске похожих проблем?
- Какие инструменты использовать для диагностики?

# Заключение

# Заключение

- Декабрь 2025 года
- Технический бюллетень о проблеме
- Передано российским вендорам и ИСП РАН
- Статья на Хабр, где куда больше деталей



# Guideline

- Нужен стенд. Люди, увы - не телепаты
- Детальный анализ условий падений
  - что за процессы работали в этот момент, какие запросы/действия выполняли
- Если тест не выявляет дефект → не то нагружается или не достаточно
- Гипер-усиливайте тест, чтобы сократить время воспроизведения
- Уменьшайте ресурсы постепенно, переводите на виртуалки
- Меняйте версии софта и параметры
- Полезно добавлять генераторы нагрузки на ресурсы ОС
- `git-bisect` / `memtester` / `stress-ng` / `rasdaemon`

# Ну а это конец...?

- Ребята из одного банка воспроизвели SEGFAULT и на `init_on_free=0`
- Но ушло не 1 минута, а 50 часов!
- Тикет в Linux bugzilla - я так и не завёл

# Благодарности

- Своим коллегам - Алексею Махмутову, Юре Соколову и Тимофею Мелко
- Коллегам из других компаний - Александру Бурцеву и Ко (product owner МБД Скала^p)
- И многих других компаниям и сочувствовавшим!
- Славе Смирнову из ПК Heisenbug за **ОГРОМНУЮ ПОМОЩЬ** в подготовке доклада!



Спасибо, что выслушали!

А данный доклад закончился

**Вопросы?**

Миша Жилин

Telegram: @mizhka

E-mail: mizhka@[gmail.com, freebsd.org]

