LINUX PLUMBERS CONFERENCE 2023, RICHMOND



## Kernel Livedump

Lukáš Hruška < lhruska@suse.cz>

Kernel Live Patching Developer

## Livedump

- Create a consistent image of memory without stopping or restarting the machine
- Debugging complex issues found only on machines which run HA services \_
- Initially introduced by Yoshida Maasanori in 2012 -



## Feature details

- restart not required =)
- copy-on-write
- consistent state
- x86 arch
- multiple target storage support
- output format supported by crash, drgn, ...



## Implementation

- decomposition into subproblems
  - wrprotect
    - make memory r/o
    - make sure the handler is called only once per PFN
  - livedump-core
    - generic module for livedump
    - sysfs
    - list of possible output destinations
  - memdump
    - Dumps into a raw block device



## Implementation

#### wrprotect diagram



-0-> not\_processed

-----

-0-> not\_processed

## Queue size problem

- memdump's handler currently saves content of a page on queue
- minimal required size of PF queue
- during stop machine state there are sensitive pages added to the queue
- dumping method also use some data structures which triggers another PF
  - might fail the whole dumping process
- starting calibration of min. required pages \* const (?)
  - wrprotect  $\rightarrow$  write to wrprotected var  $\rightarrow$  tracing write finish on selected storage
    - memdump tp: block:block\_rq\_complete (wrprotected var's PFN == tp.pfn)
- unlock as many as possible known PFNs having data needed to perform the write op (?)



## PF handling variants

- sometimes small inconsistency in dumped data might be OK \_
- two possible versions:
  - 100% consistent variant with possible performance impact
    - active waiting for space on queue (?)
    - dynamically resizing the PF queue till -ENOMEM (?)
    - sweep and PFs sharing same queue with sweep heuristic % to use (?) •
  - possible inconsistent variant with minimal performance impact
- 100% consistent variant with minimal performance impact possible to not complete (?)



### TODOS

- vmap area support (only correct synchronization needs to be solved)
- interrupt instead of stop-machine state
- analysis on impact of 4kB split of pages on performance (TLB overhead)
  - huge page support / PT restoration
- other target storage support (file, network, ...)





# Thank you

© SUSE LLC. All Rights Reserved. SUSE and the SUSE logo are registered trademarks of SUSE LLC in the United States and other countries. All third-party trademarks are the property of their respective owners.

For more information, contact SUSE at: +1 800 796 3700 (U.S./Canada)

(-)

Frankenstrasse 146

90461 Nürnberg

www.suse.com