

# MADV\_LAZYFREE

“rcu\_free” for user memory  
Rik van Riel

# MADV\_DONTNEED overhead

TLB flush CPU overhead is ... bad

- In MADV\_DONTNEED calling thread: 50% of CPU time spent on TLB flush
- Add in TLB flushing on other CPUs: **90-97%**
- >90% of TLB flushes from MADV\_DONTNEED & MADV\_FREE
- Numbers with 10s of CPUs, newer systems have 100s of CPUs

RCU provides inspiration

- Make data inaccessible
- Users go away over time
- Freeing is cheap without users left

# Userspace API

API for use by malloc libraries and runtimes. API was simplified in collaboration with Jemalloc developers. Most user programs never need to touch this directly.

## MADV\_LAZYFREE

- Mark a memory range (start, end) for lazy freeing
- Program needs to ask the kernel before reusing the range

## MADV\_REUSE

- Make a previously MADV\_LAZYFREE'd virtual memory range (start, end) available for reuse
- Return values (open to better ideas)
  - 0: memory is available for reuse, the process has other MADV\_LAZYFREE ranges
  - ENOMEM: memory is available for reuse, no more MADV\_LAZYFREE ranges present

# MADV\_LAZYFREE pseudo-code

advise(..., MADV\_LAZYFREE)

- Kernel can discard data in this memory range, on its own schedule
- Advance mm\_struct tlb\_gen number
  - Currently in arch/x86, move to generic code?
- Store MADV\_LAZYFREE memory range
  - start, end, tlb\_gen
  - Some tree attached to mm\_struct
  - Link mm into MADV\_LAZYFREE range shrinker
- Mark PTEs non-present
  - PMD\_NONE, or some swap type like migration?
  - Keeps CPUs from loading valid page mappings for this address range

# MADV\_REUSE pseudo-code

`advise(..., MADV_REUSE)`

- Data in this range will be preserved after `MADV_REUSE` returns
- Mark PTEs present again
  - Pages may have been reclaimed already (nothing to do)
  - Pages might still be around if memory pressure is low
- Remove `MADV_LAZYFREE` memory range from mm tree
- No need to flush TLBs
  - If CPU `tlb_gen`  $\geq$  `MADV_LAZYFREE tlb_gen`, it has no mapping for this range
  - If CPU `tlb_gen`  $<$  `MADV_LAZYFREE tlb_gen`, there was no reclaim, and some of the same pages might still be mapped

# MADV\_LAZYFREE TLB flush requirements

No TLB flushes in MADV\_LAZYFREE or MADV\_REUSE themselves.

The TLB needs to be flushed on reuse, if:

- The virtual memory range is reused for something else
- A physical memory page is freed

If reuse is “far enough” into the future, the TLB flushes can happen naturally at context switch time, and we maybe able to avoid sending IPIs.

# Context switch pseudo-code

Piggyback on existing lazy TLB code

- Check whether CPU `tlb_gen` < `mm_struct tlb_gen`
  - If CPU `tlb_gen` is behind, flush the TLB
- The lazy TLB code on x86 already does this today
- `MADV_LAZYFREE` expands lazy TLB mode, minor change needed

# Reclaim & migration pseudo-code

- Check whether page is part of MADV\_LAZYFREE range
- Look up tlb\_gen of MADV\_LAZYFREE range
- Clear (already non-present) PTE
- Iterate over CPUs with mm loaded
  - CPU tlb\_gen < MADV\_LAZYFREE tlb\_gen? Send TLB flush IPI
  - CPU tlb\_gen >= MADV\_LAZYFREE tlb\_gen? Nothing to do
- Free page

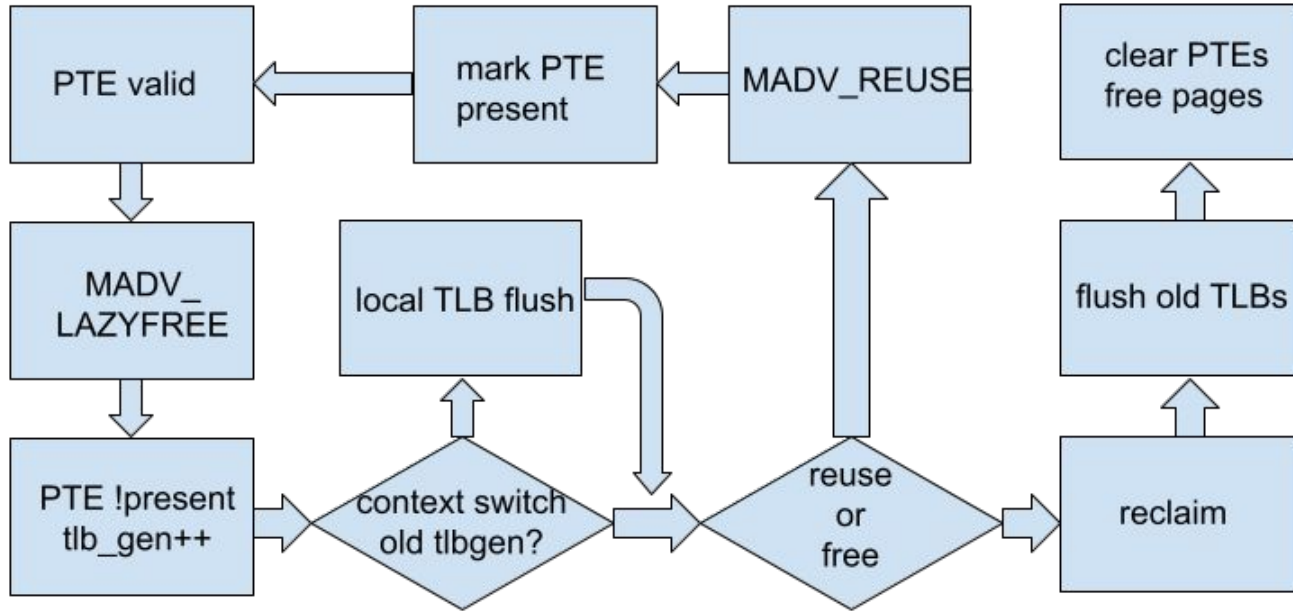


# MADV\_LAZYFREE potential optimizations

## Potential optimizations

- Keep a small buffer of flushed ranges
  - Do partial flushes at context switch time, instead of full flushes
  - Gives 1-2% IPC improvement in tests (without MADV\_LAZYFREE)
  - Effect reduced when system has less idle time
- Use lazy freeing for munmap() too?
  - Reduce mmap\_lock contention
- Integrate with Byungchul Park's Lazy Unmap Flush project?

# MADV\_LAZYFREE lifetime rules



**MADV\_LAZYFREE**

**Questions?**

**Suggestions?**

**Ideas?**