

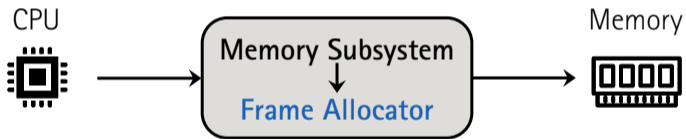
LLFree: Scalable and Optionally-Persistent Page-Frame Allocation

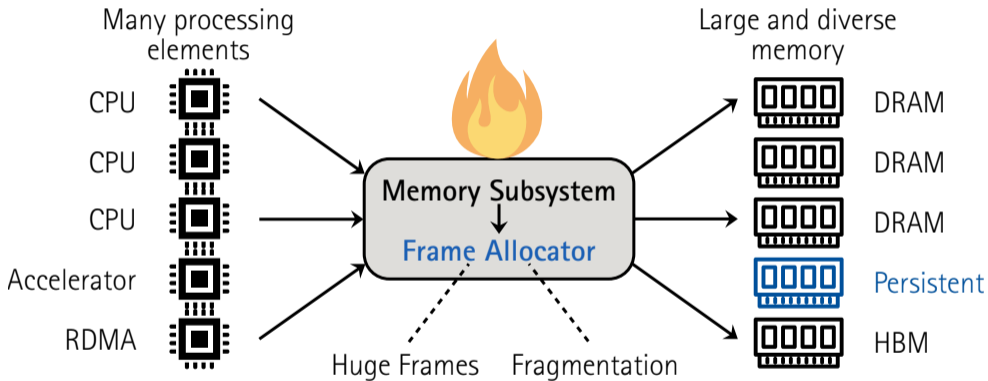
USENIX ATC'23

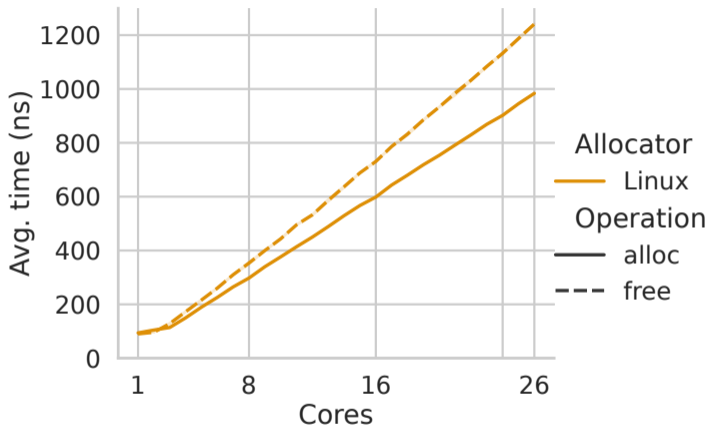
Lars Wrenger Florian Rommel Alexander Halbuer Christian Dietrich Daniel Lohmann

Leibniz Universität Hannover
wrenger@sra.uni-hannover.de

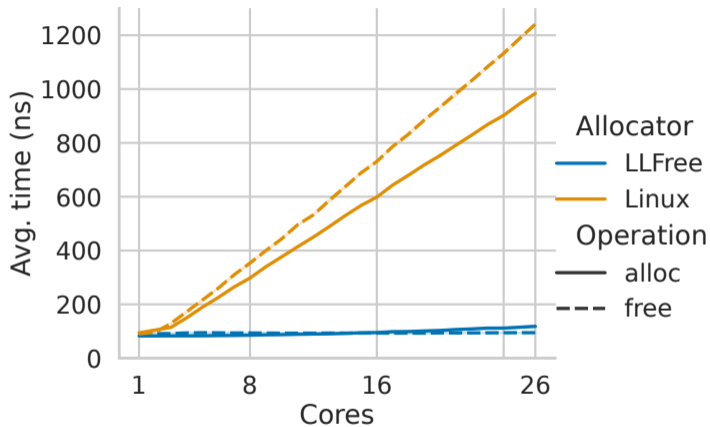
2023-07-12



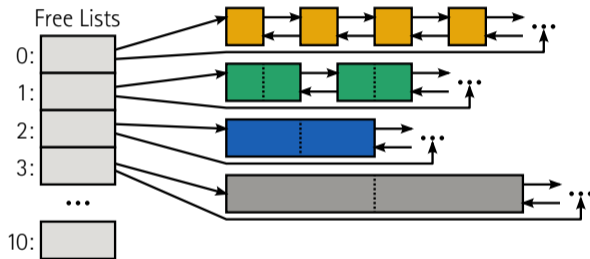


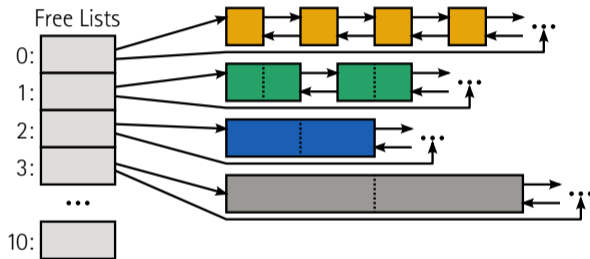


- Allocate 64 GiB as 4 KiB pages in parallel (Linux 6.0)

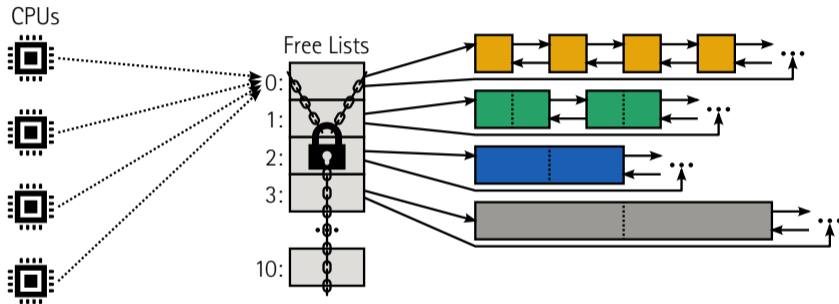


- Allocate 64 GiB as 4 KiB pages in parallel (Linux 6.0)

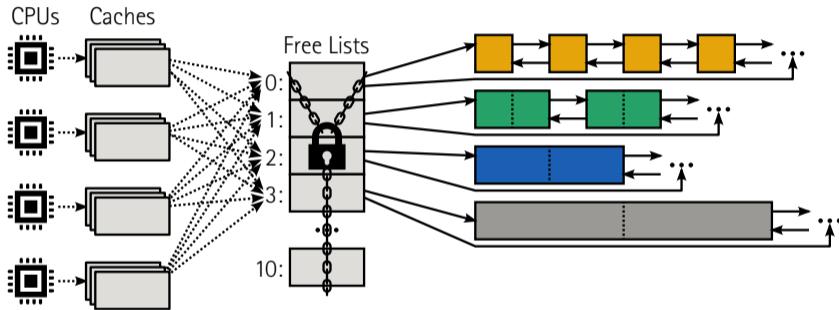




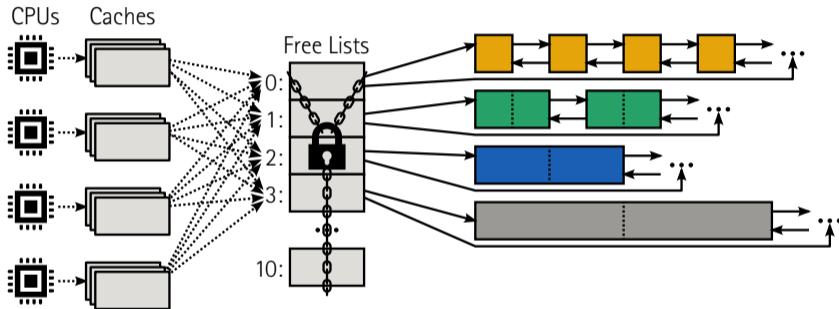
■ Split/Merge Costs



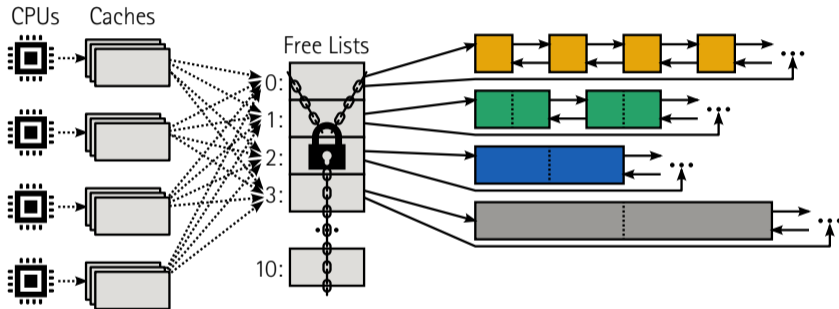
- Split/Merge Costs
- Scalability



- Split/Merge Costs
- Scalability



- Split/Merge Costs
- Scalability
- Fragmentation



- Split/Merge Costs
- Scalability
- Fragmentation
- Persistency

Design Principles

Respect Hardware: hardware-specific page sizes, cache-line granularity

Avoid Sharing: preventing memory sharing bottlenecks

Careful Redundancy: avoid synchronization problems for persistent state

Design Principles

Respect Hardware: hardware-specific page sizes, cache-line granularity

Avoid Sharing: preventing memory sharing bottlenecks

Careful Redundancy: avoid synchronization problems for persistent state

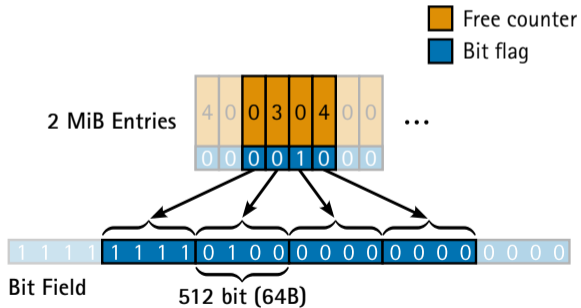


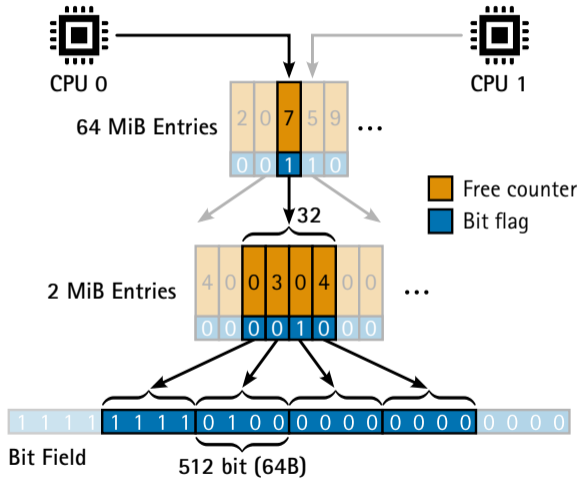
- **Lock- and log-free** → LLFree
 - Locks (Mutex): Not crash consistent!
 - Logging (ACID): Slow and intensifies write wearing!
- Atomic updates → "Jump between valid states"



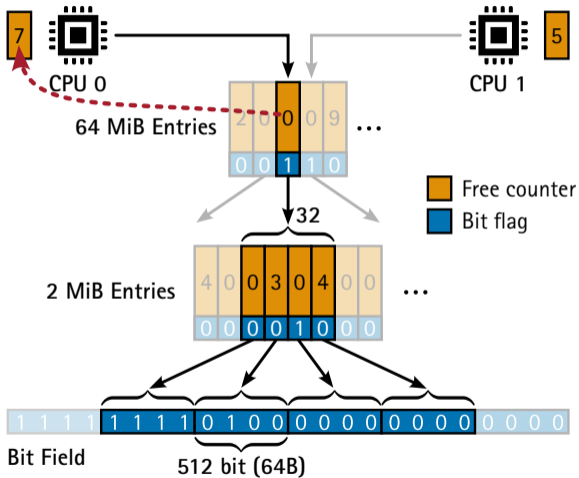
1 1 1 1 1 1 1 1 0 1 0 0 0 0 0 0 0 0 0 0 0 0

Bit Field

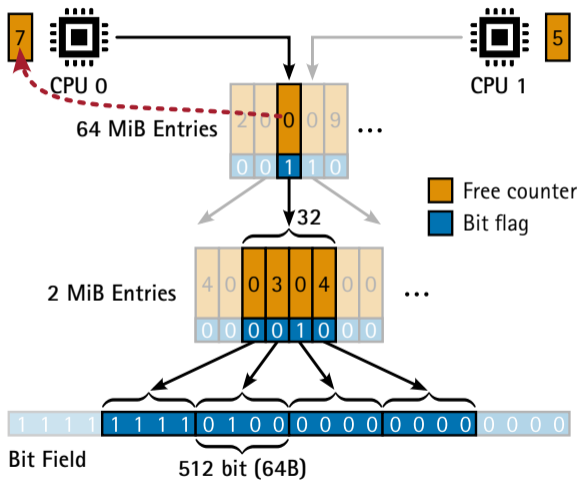




■ **Avoid Sharing** in 64 MiB Chunks



- **Avoid Sharing** in 64 MiB Chunks
 - **False-sharing** on the 64 MiB counters
⇒ Local split counters



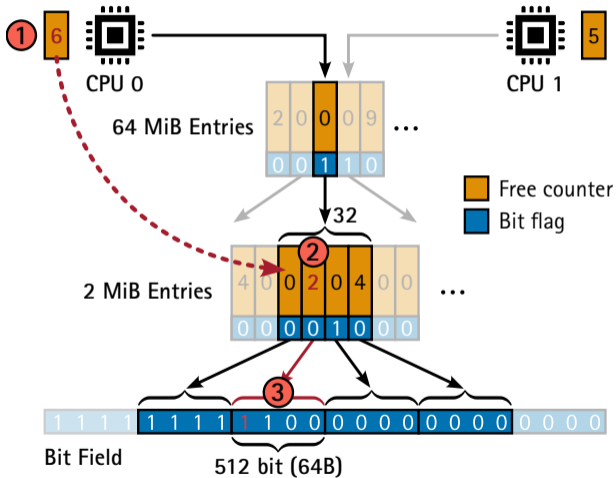
■ Avoid Sharing in 64 MiB Chunks

- **False-sharing** on the 64 MiB counters
⇒ Local split counters

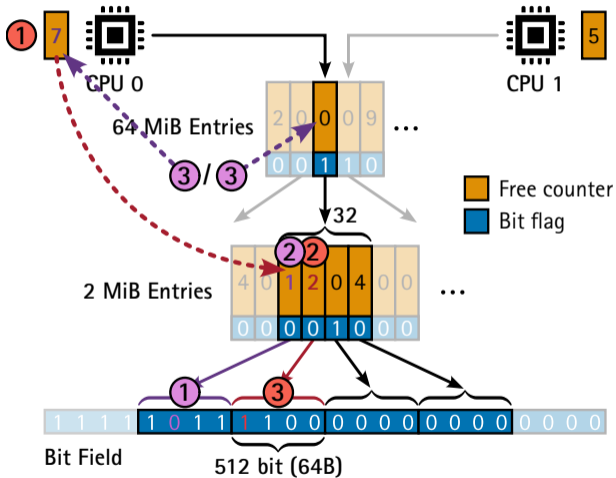
■ Reservation policy

→ Avoid Fragmentation

- Prioritize already fragmented
64 MiB Areas

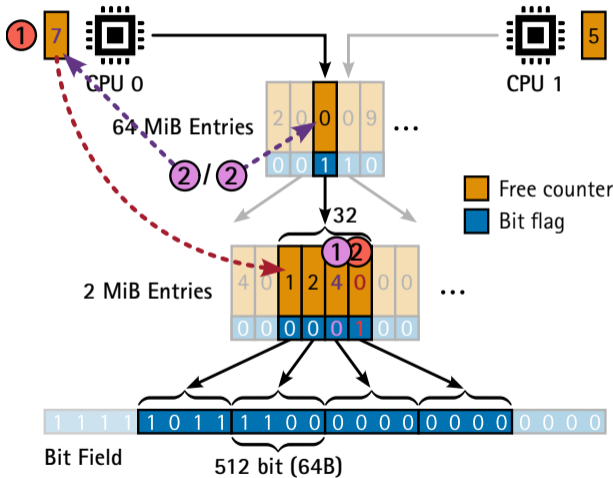


4 KiB Allocation: 3 cache lines



4 KiB Allocation: 3 cache lines

4 KiB Free: 3 cache lines

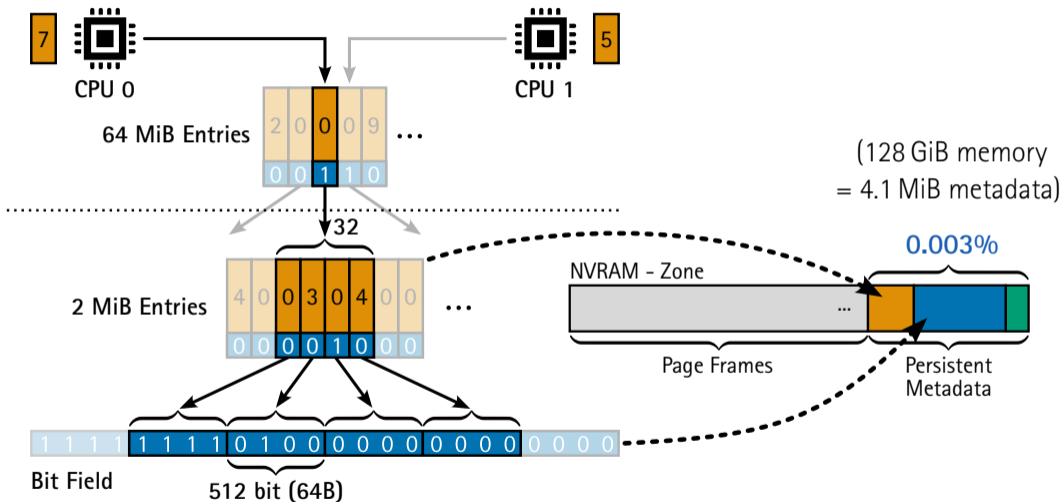


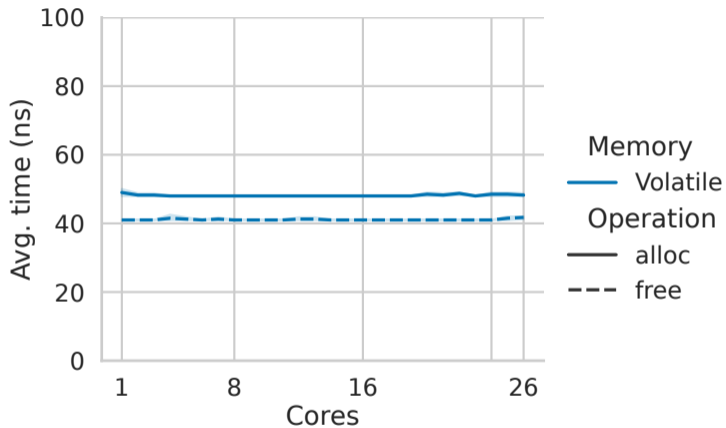
4 KiB Allocation: 3 cache lines

4 KiB Free: 3 cache lines

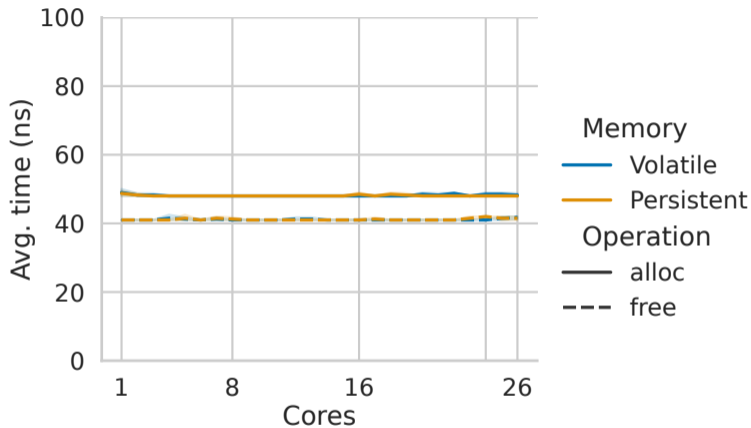
2 MiB Allocation: 2 cache lines

2 MiB Free: 2 cache lines

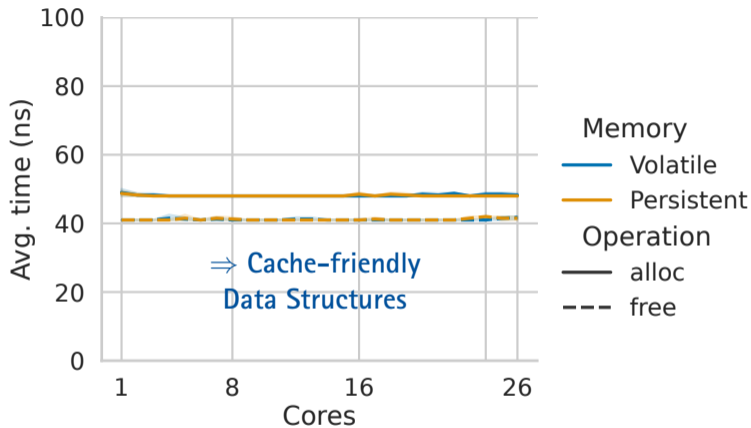




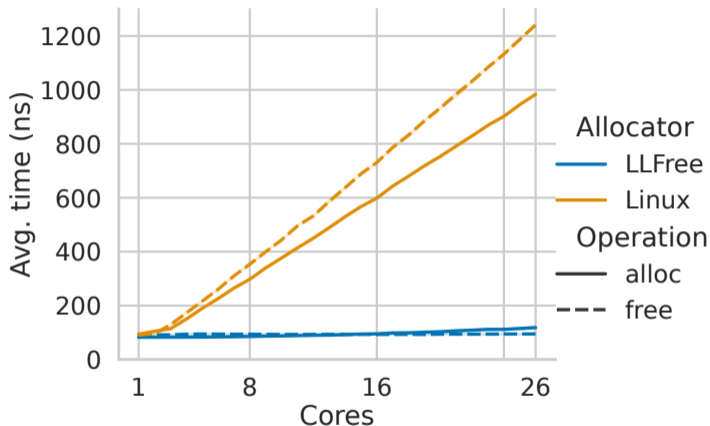
- Allocating 64 GiB as 4 KiB frames in parallel



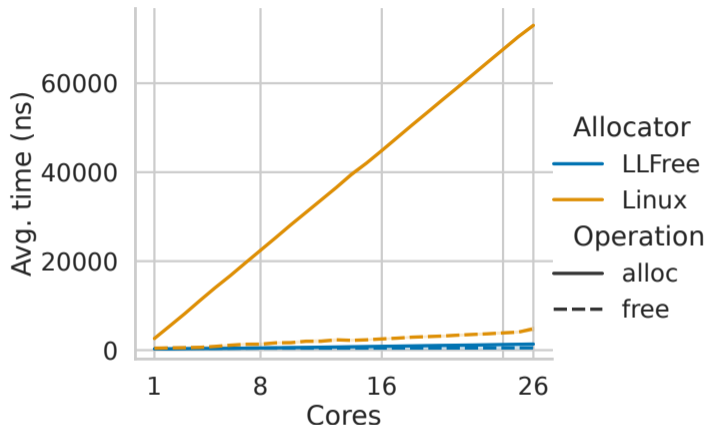
- Allocating 64 GiB as 4 KiB frames in parallel



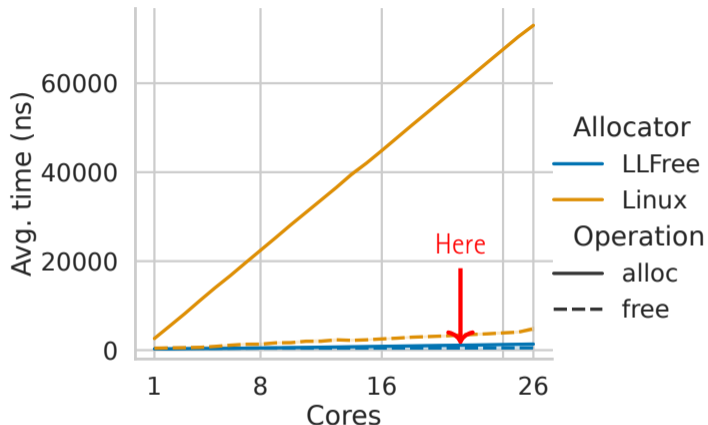
- Allocating 64 GiB as 4 KiB frames in parallel



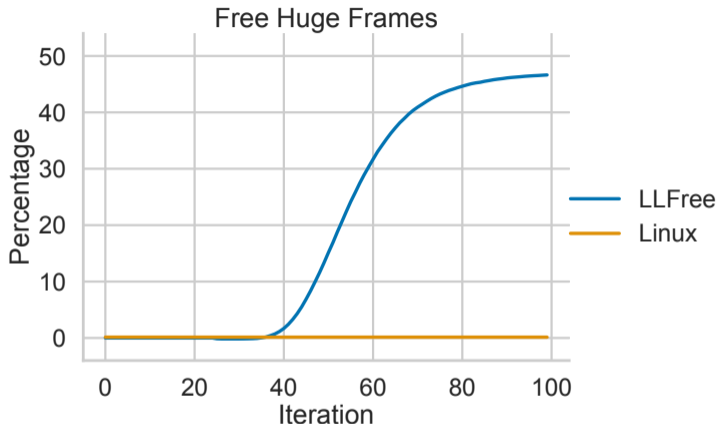
- Allocating 64 GiB as 4 KiB frames in the Linux Kernel (6.0)



- Allocating 64 GiB as **2 MiB** frames in the Linux Kernel (6.0)
- Linux has a severe performance problem (patched in 6.4)



- Allocating 64 GiB as **2 MiB** frames in the Linux Kernel (6.0)
- Linux has a severe performance problem (patched in 6.4)



- Initially, half of the memory is allocated and entirely fragmented
- In every iteration, randomly reallocate 10% of the allocated pages

- The Linux page allocator
 - Does not scale well
 - Prone to fragmentation
 - Not persistent
- LLFree
 - Lock- and log-free
 - Excellent multicore scaling
 - Over-time defragmentation
 - Persistent on NVM



Try LLFree

Open Source at
github.com/luhsra/llfree-rs



- Email: wrenger@sra.uni-hannover.de