

# 针对用户态定长申请的内存分配器

## 星火十六期终审答辩

刘般若

清华大学

2022 年 5 月 22 日



## ① 课题背景

## ② 研究内容

## ③ 研究结果

## ④ 实现细节与创新

## ⑤ 未来目标

## ① 课题背景

## ② 研究内容

## ③ 研究结果

## ④ 实现细节与创新

## ⑤ 未来目标

# Motivation

- Modern volatile memory allocators over DRAM are quite optimized.(tcmalloc<sup>1</sup> jemalloc<sup>2</sup>)
- However, emerging NVM(non-volatile memory, or persistent memory) applications require allocators to be crash-consistent.
- To achieve that, general-purpose NVM allocator has to use transaction-like method, which devastates its performance.
- For low-level applications, a fast, scalable persistent allocator is in dire need.

---

<sup>1</sup><https://github.com/google/tcmalloc>

<sup>2</sup><https://github.com/jemalloc/jemalloc>

# Observation

- Since programmers endeavour to cater hardware characteristics, most of alloc requests are confined in a few size types.
- For indexes, such requests are fixed and highly skewed, usually related with values below.
  - 64 Bytes, typical CPU cacheline size
  - 256 Bytes, Optane DIMM write buffer size
  - 4K Bytes, page size for most X86 systems
- Thus, we only need to deal with a few fixed sizes.

## ① 课题背景

## ② 研究内容

## ③ 研究结果

## ④ 实现细节与创新

## ⑤ 未来目标

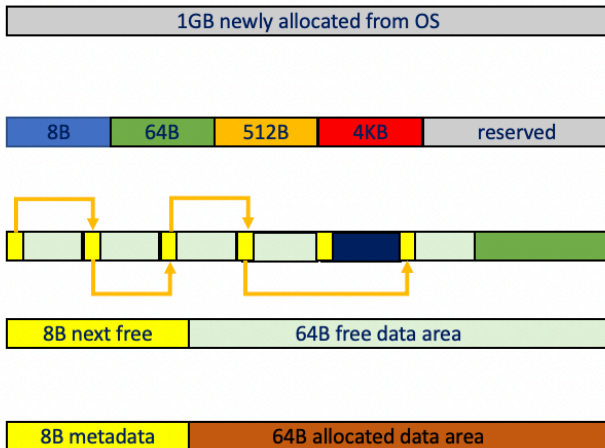
# User-awared Fixed Size Memory Allocator

- In this project, we present a memory pool, which accepts a size vector from user-space, provides alloc utility for such sizes, achieving:
- High throughput, low latency
- Near-linear scalability with thread number
- Low internal fragmentation
- Crash consistency (non-volatile)
- Outperform current mainstream persistent allocator. (pool in PMDK, Persistent Memory Development Kit <sup>3</sup>)

---

<sup>3</sup><https://pmem.io/pmdk/>

# Overview





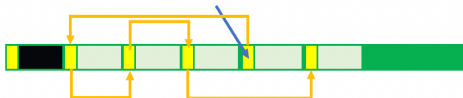
# Alloc and Free



• initial state



• alloc(64)



• free(addr), size is inferred from header

## ① 课题背景

## ② 研究内容

## ③ 研究结果

## ④ 实现细节与创新

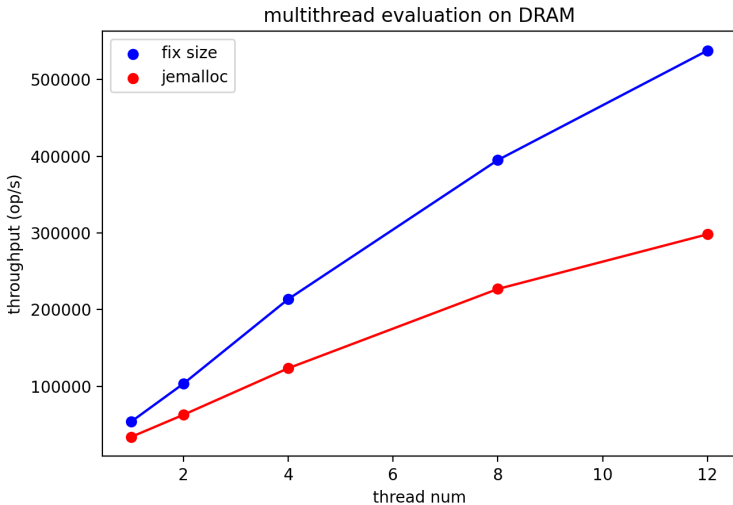
## ⑤ 未来目标

---

# Performance as Volatile Allocator

- We compared our allocator with current industrial-level general purpose allocator: jemalloc from Facebook
- Data are randomly generated alloc-free sequence of size in  $\{8, 64, 512, 4K, 2M\}$   
About 50000 operations per thread. Borrowing rarely occurs.

# Performance as Volatile Allocator



## Performance as Persistent Allocator

- We compared with PMDK::pool on real workloads.
- We changed allocator of non-volatile index from PMDK to fixed size pool. And test the load phase duration, which consists about  $2 * 10^6$  insert operation.<sup>4</sup>
- We tested our project on three indexes, FastFair<sup>5</sup>, p-masstree<sup>6</sup> and CCEH<sup>7</sup>

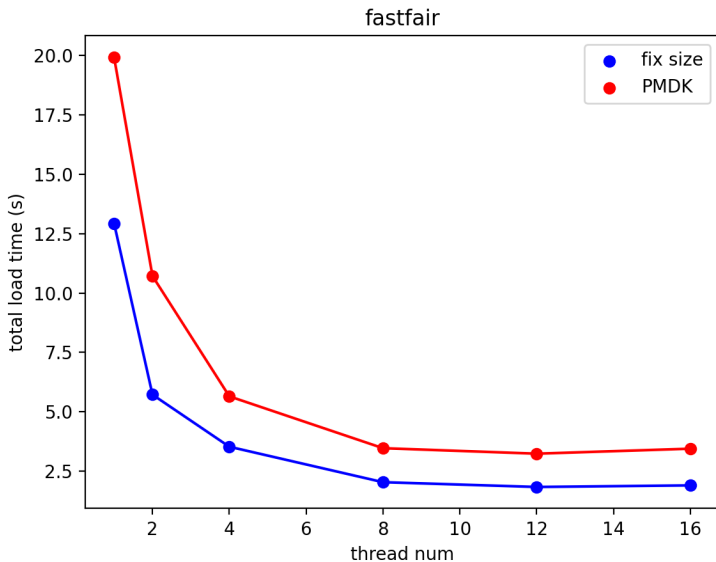
---

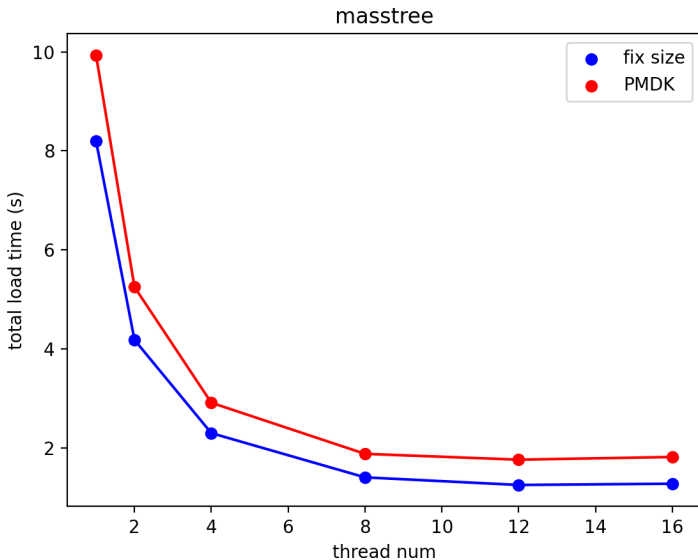
<sup>4</sup>Wang, Qing, et al. "Nap: A Black-Box Approach to NUMA-Aware Persistent Memory Indexes."

<sup>5</sup>Hwang, Deukyeon, et al. "Endurable Transient Inconsistency in Byte-Addressable Persistent B+-Tree."

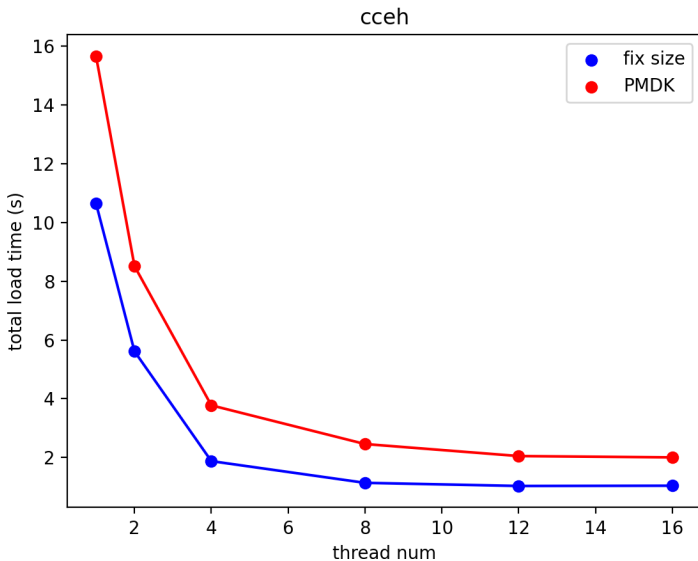
<sup>6</sup>Lee, Se Kwon, et al. "Recipe: Converting concurrent DRAM indexes to persistent-memory indexes."

<sup>7</sup>Nam, Moohyeon, et al. "Write-Optimized Dynamic Hashing for Persistent Memory."

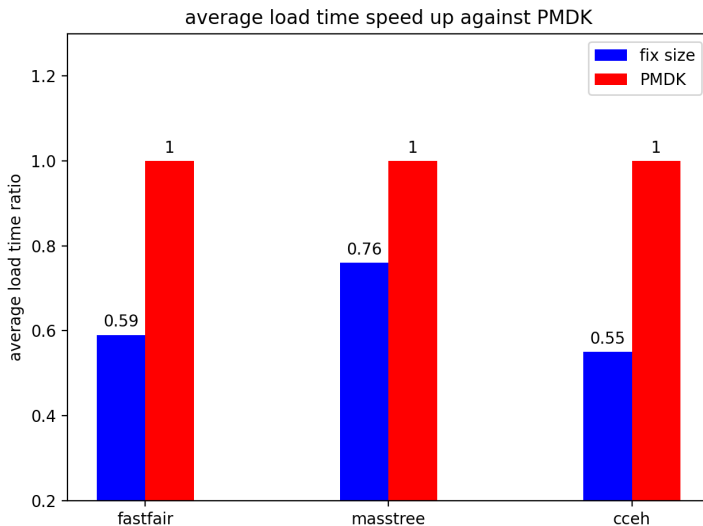








# Result



① 课题背景

② 研究内容

③ 研究结果

④ 实现细节与创新

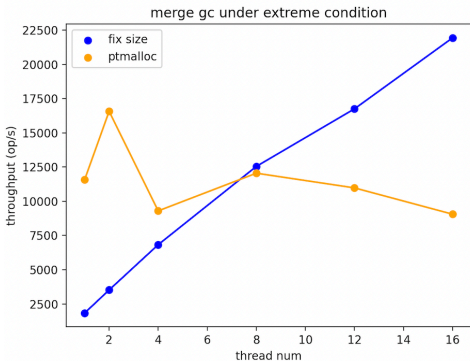
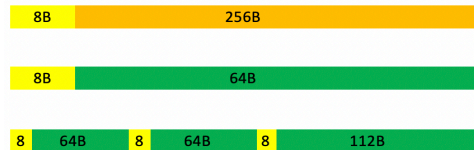
⑤ 未来目标

# Consistency

For a general purpose allocator, it uses refined tricks and techniques to deal with arbitrary sizes. As a result, it is much more complex to persist its structure on NVM.

However, since our goal is inherently simple, our design is much concise, and we can use techniques like write-ahead logging to guarantee consistency.

# Reduce Internal Fragmentation



- Through rarely happens, we have invented practical plan for borrowing and lending to ensure correctness. We use header-division or redundant allocation.
- We also use heuristic segmented merge for efficiency and parallelizability under extreme workloads.

## Small Size Optimization

- For small allocations i.e. 8 Bytes, using header will introduce about 50% internal fragmentation.
- A self-decode pointer is easy to implement and has decent performance. But linked-list is not cache-favorable. Performance is compromised when there is successive requests for small size.
- Using multi-level bitmap and random magic number checksum. We can fit the available slots in L1 cache, which improves the performance.

## ① 课题背景

## ② 研究内容

## ③ 研究结果

## ④ 实现细节与创新

## ⑤ 未来目标

## Future Plan

- Do more experiments on other NVM applications, i.e. key-value storage, filesystem...
- Optimize recovery strategy, using techniques like checkpoint to make the allocator more durable.
- Hopefully, make it a library that can actually be integrated into other projects.



*Thanks!*