

# Boosting Data Center Performance via Intelligently Managed Multi-backend Disaggregated Memory

Jing Wang, Hanzhang Yang, Chao Li\*, Yiming Zhuansun, Wang Yuan, Cheng Xu, Xiaofeng Hou, Minyi Guo, Yang Hu, Yaqian Zhao

Shanghai Jiao Tong University, Tsinghua University, IEIT SYSTEMS Co., Ltd

## Introduction

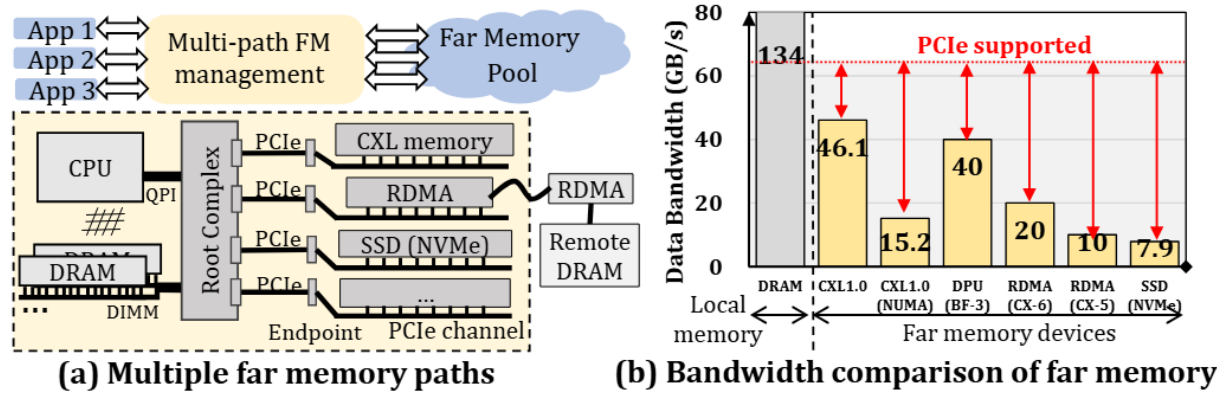


Figure 1: (a) Multi-backend far memory. (b) Bandwidth comparison of various far memory technologies.

In recent years, disaggregated memory (DM) architecture stands out as a way to enhance data center capabilities by offering highly flexible memory expansion. In this case, a memory-hungry monolithic server can access a PCI Express (PCIe) based secondary memory device (i.e., far memory) with low data access latency.

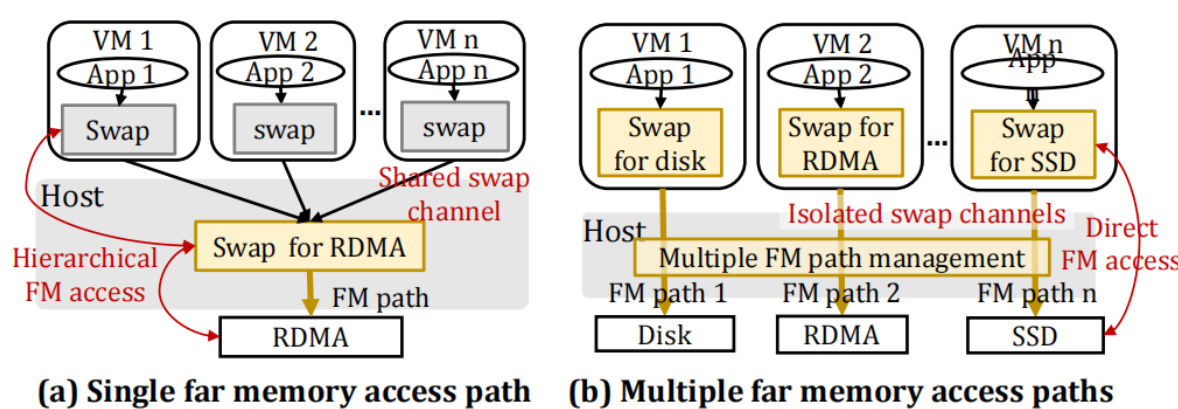
Unfortunately, prior works mostly limit their designs on a **single far memory (FM) device** based on PCIe, which cannot meet the needs of high data/task throughput in today's data center. There is a wide gap between the data bandwidth that a single FM device could support (from 7.9 to 46 GB/s) and the maximum available bandwidth that modern PCIe protocol could provide, as shown in Figure 1.

Incorporating multiple FM devices and oversubscribing the PCIe subsystem allows one to push the limit of server data throughput. This approach, which we call **multi-backend disaggregated memory**, unfortunately cannot be implemented solely based on current data transfer protocols and device drivers. Existing FM management schemes are blind to the possible multiple physical FM channels: they logically only support one data exchange path to FM devices due to original swap mechanism limitation.

We introduce **xDM**, a novel FM management system engineered to support high-performance data analytic workloads running on the new multi-backend disaggregated memory architecture.

## Challenges and Motivations

- Far Memory Usage Bottleneck:** It is impractical to directly support multi-path far memory management in the OS. A naive VM-based FM system has limitations due to the hierarchical swap architecture and considerable data swap overhead. Removing FM usage bottleneck requires expanding data swap channels and tapping into light-weight swapping.
- Far Memory Usage Effectiveness:** Most of the prior works follow a simple idea: by offloading part of data to far memory based on workload behaviors, the local DRAM can retain more latency-sensitive tasks. Applications often exhibit more complex performance variation, which requires a detailed analysis of page behaviors. These works unfortunately overlook the importance of several key system parameters.



## Design

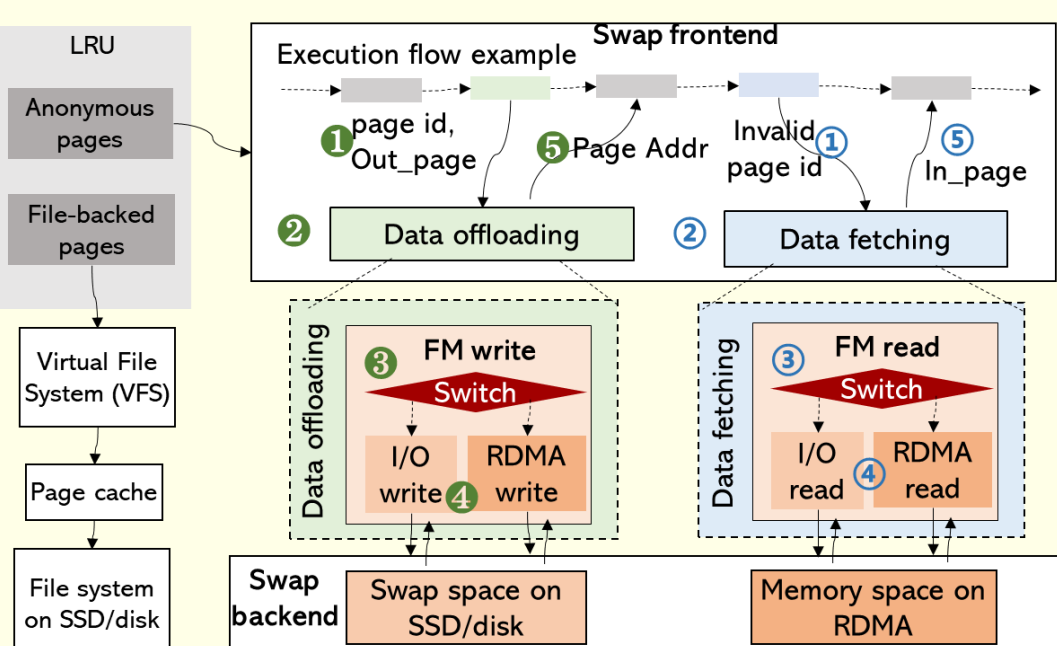
xDM is the first multi-path FM management system that not only allows *simultaneous* access of multiple FM devices but also provides *dynamic and implicit* adjustment of FM paths based on program behaviors.

### (1) Dynamic FM Switching Mechanism:

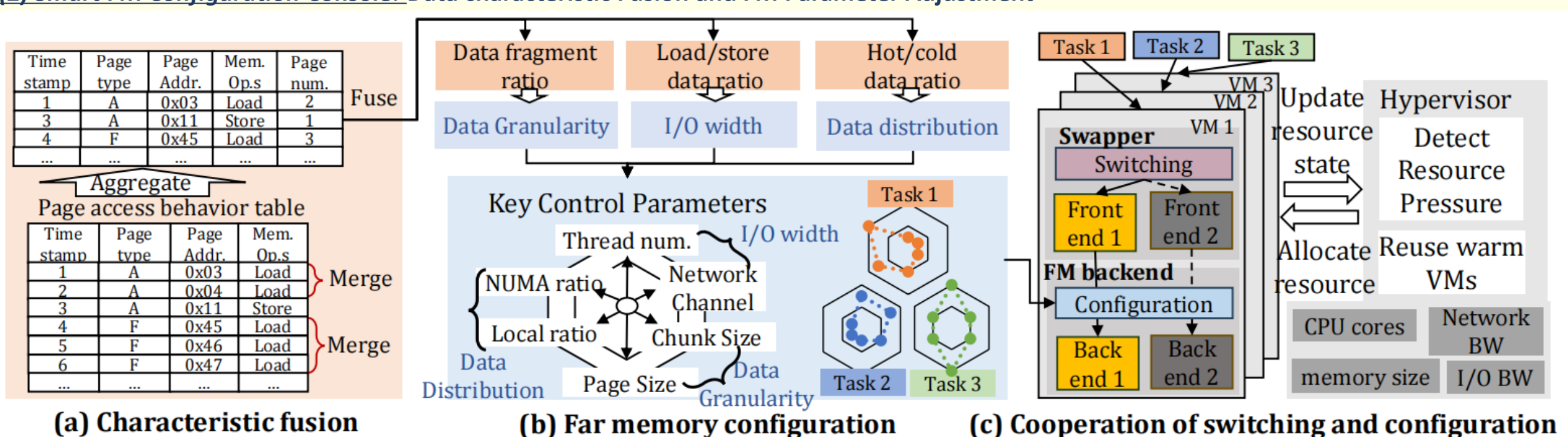
**Switchable FM Swapper:** The switchable swapper includes a modified swap frontend plus a variety of adaptive FM swap backends. Leveraging the data swap interfaces provided by Frontswap, we redirect page swapping to our customized FM read and write functions.

**Implicit FM Switching Strategy:** We design a far memory switching strategy to choose FM backends based on the page distribution statistics and application's sensitivity to different FM paths.

→ Data flow → Control flow ①~⑥ Offloading steps ⑦~⑨ Fetching steps

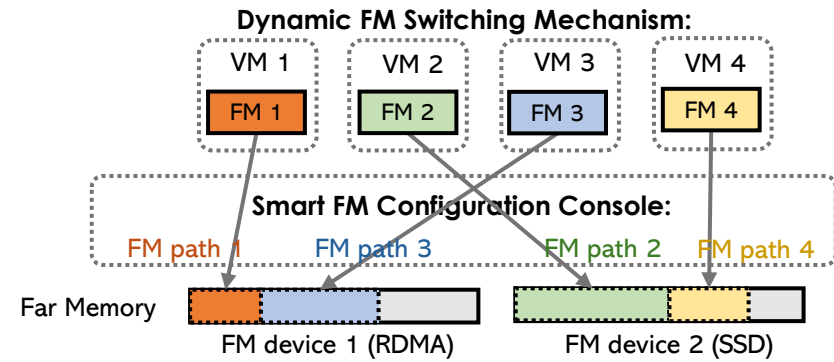


### (2) Smart FM Configuration Console: Data characteristic Fusion and FM Parameter Adjustment



xDM's adaptive FM parameter configuration. The characteristic fusion module provides a synthesis of information from page access traces. xDM can configure and fine-tune multi-backend FM from various aspects.

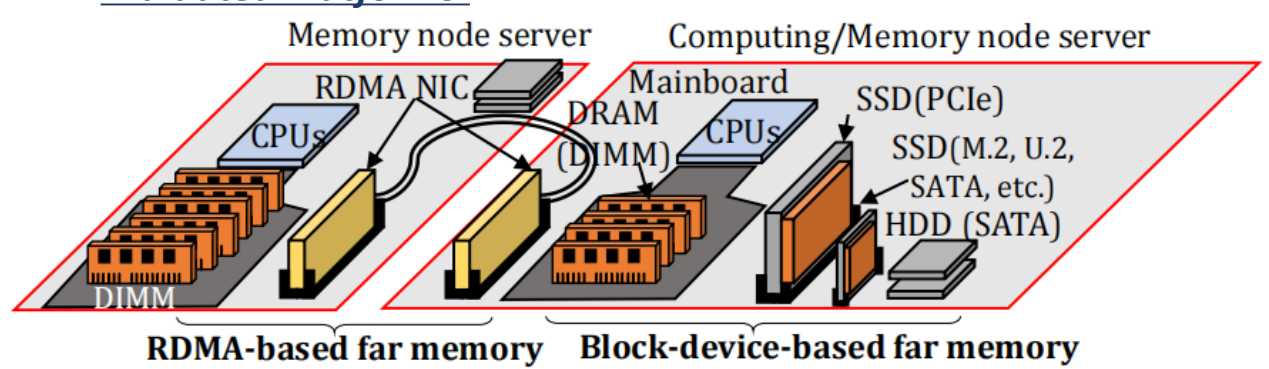
## Conclusion



- In this work, we design and implement xDM, a novel **multi-backend far memory system** with high bandwidth utilization and application performance.
- By turning the conventional swap mechanism into a **switchable data swap module**, we successfully realize simultaneous multi-path FM access.
- Based on a rich synthesis of application **page data**, we tailor the far memory path configurations to the needs of various applications.
- Our work shows up to **3.9x data swap** performance speedup, **2.8x data throughput** increase, and **5.1x data center task serving throughput** improvement compared with state-of-the-art works.
- Our design provides a **flexible solution** to scale out far memory access paths and an efficient way to manage them on monolithic servers. We expect that our design can improve the execution performance and **memory usage effectiveness** of memory-hungry tasks in cloud and near-edge micro data centers.

## Evaluation

### Evaluated Platforms:

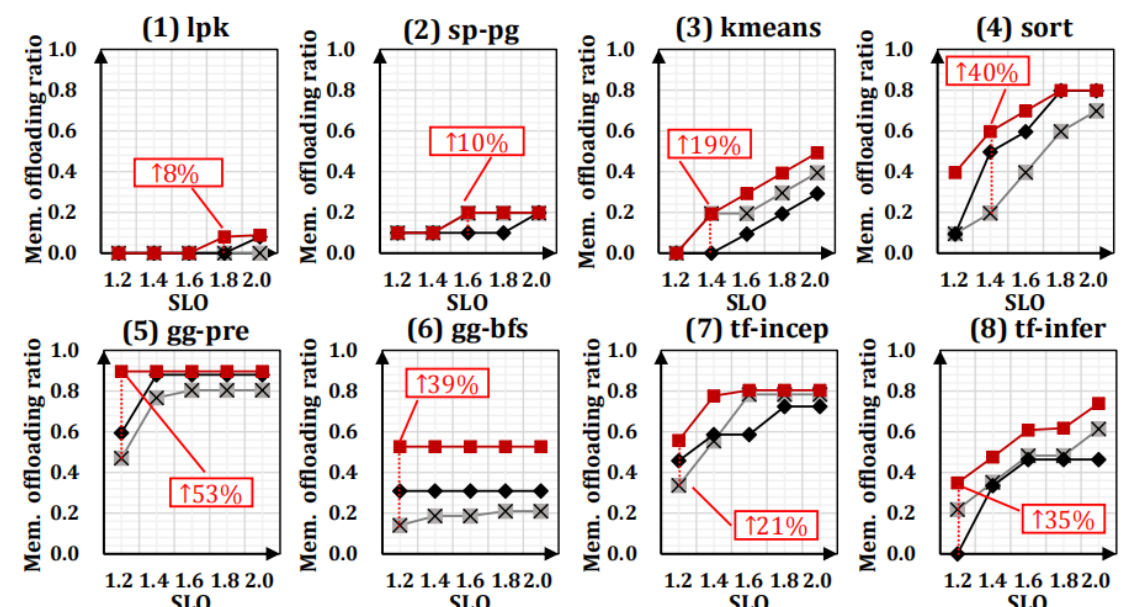


We show the swap performance speedup (Sp.) of our xDM compared with baselines on different backends. We categorize application swap features into two types: swap-sensitive (S, average Sp.  $\leq 1.5x$ ) and swap-friendly (F, average Sp.  $> 1.5x$ ).

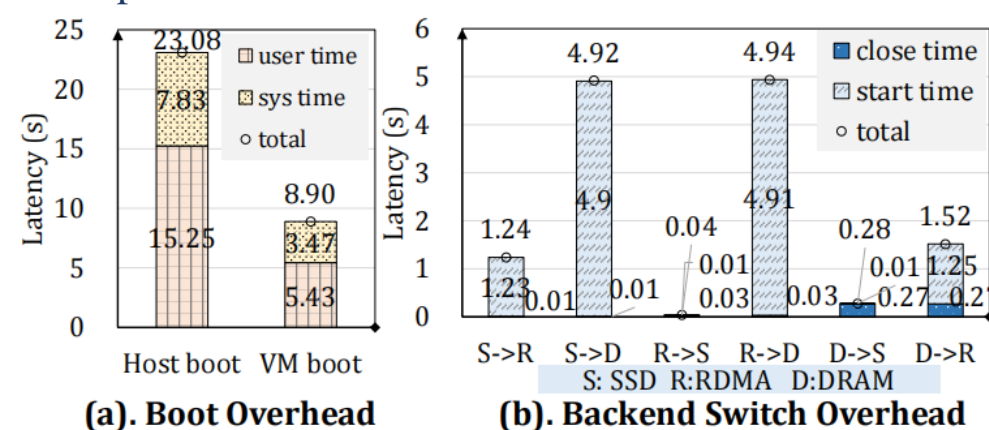
Evaluated Workload	stream	lpk	kmeans	sort	s-pg	gg-pre	gg-bfs	lg-bfs
Swap Feature	S	S	S	S	S	F	S	F
Sp. on DRAM	1.32x	1.18x	1.64x	1.05x	1.44x	2.24x	1.29x	2.00x
Sp. on SSD	1.01x	1.52x	0.88x	0.86x	1.01x	1.02x	1.18x	1.40x
Sp. on RDMA	1.25x	1.09x	1.40x	1.40x	1.37x	2.06x	1.19x	2.24x
Average Speedup	1.19x	1.26x	1.31x	1.11x	1.28x	1.77x	1.22x	1.88x

Evaluated Workload	lg-comp	lg-mis	tf-infer	tf-incep	clip	tf-te	chat-int	bert
Swap Feature	F	F	F	F	S	F	F	S
Sp. on DRAM	2.43x	2.17x	1.88x	1.72x	0.82x	1.28x	1.15x	1.03x
Sp. on SSD	1.52x	1.36x	1.51x	1.34x	0.91x	2.16x	1.92x	1.75x
Sp. on RDMA	2.22x	2.07x	2.70x	2.53x	2.46x	2.55x	3.89x	1.10x
Average Speedup	2.05x	1.86x	2.03x	1.86x	1.40x	2.00x	2.32x	1.29x

Our design shows **larger memory offloading ratios** than baselines on evaluated workloads under different SLOs.



The acceptable **overhead** of virtualization and backend switching.



For more details, please see our full paper in SC'24.

## Acknowledgement

We sincerely thank all the anonymous reviewers for their valuable comments. This work is supported by the National Key R&D Program of China (No. 2022YFB4501702), and the National Natural Science Foundation of China (No. 62122053). The corresponding author is Chao Li. Thanks for the support of Emerging Parallel Computing Center (EPCC).