



Rethinking Tiered Storage: Talk to File Systems, Not Device Drivers

Jiyuan Zhang

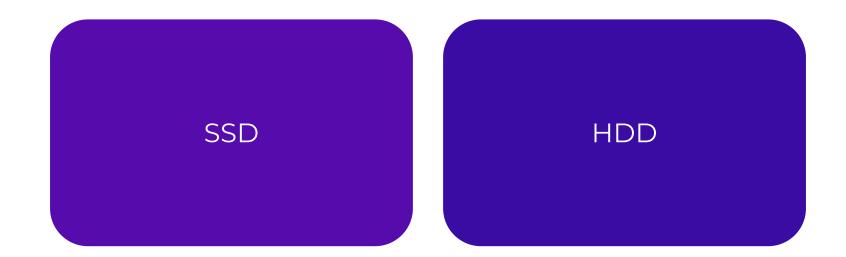
Jongyul Kim, Chloe Alverti, Peizhe Liu, Weiwei Jia, and Tianyin Xu





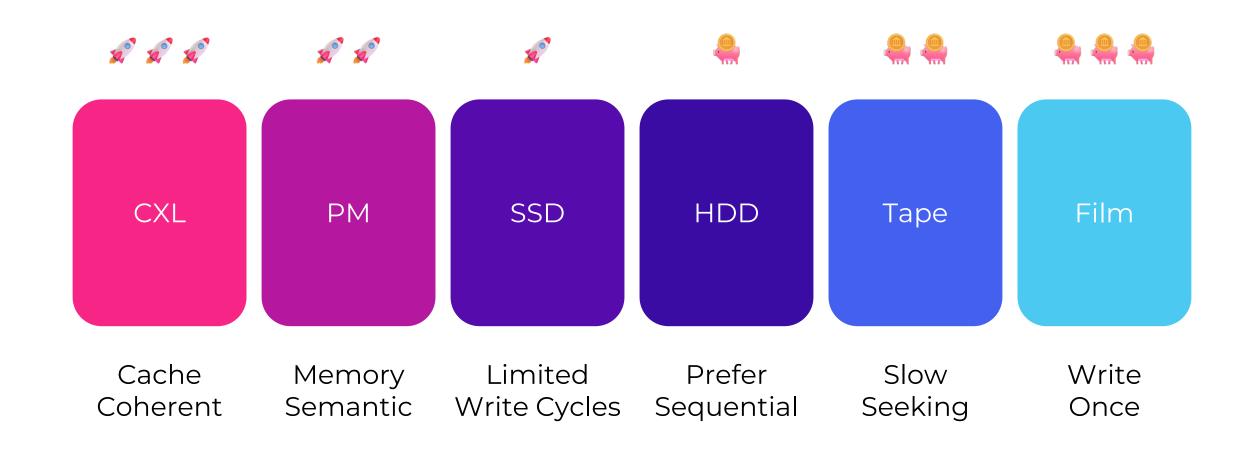


Storage devices are diversifying





Storage devices are diversifying





New devices lead to heterogeneous configurations

We need to use the new devices AND the old ones.



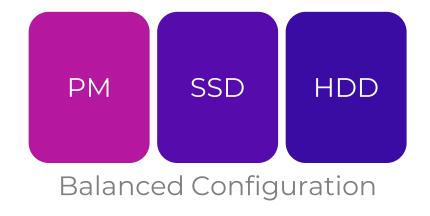
Latest Features

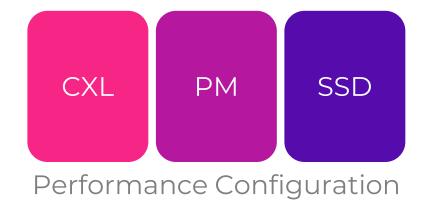
Performance Advantage

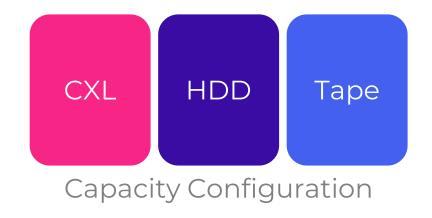
Data Reliability Backward Compatibility



New devices lead to heterogeneous configurations











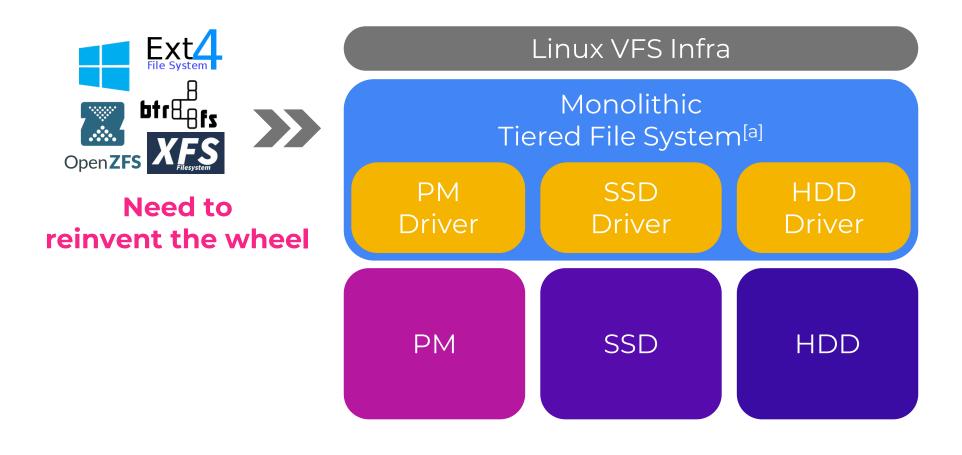
New devices lead to heterogeneous configurations

Research Question:

How to build **practical** file system for **diverse** storage configurations?



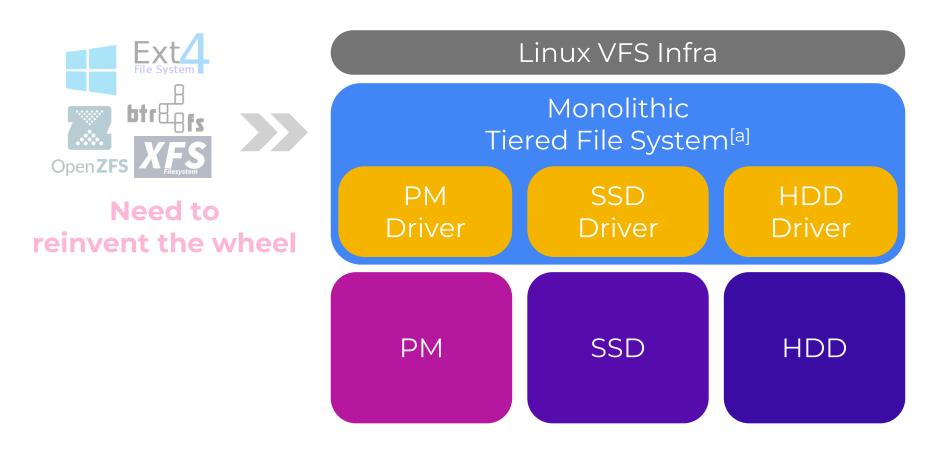
Existing solutions are falling behind

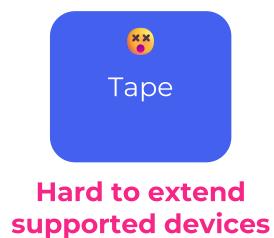


a. We observe the imperfections in the state-of-the-art monolithic tiered file system Strata



Existing solutions are falling behind

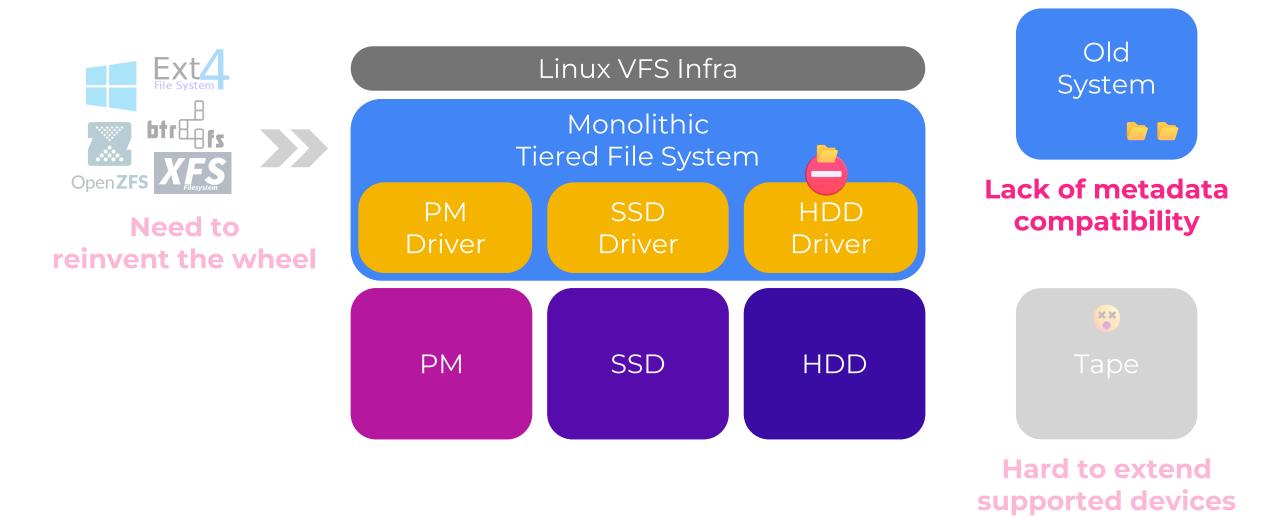




a. We observe the imperfections in the state-of-the-art monolithic tiered file system Strata

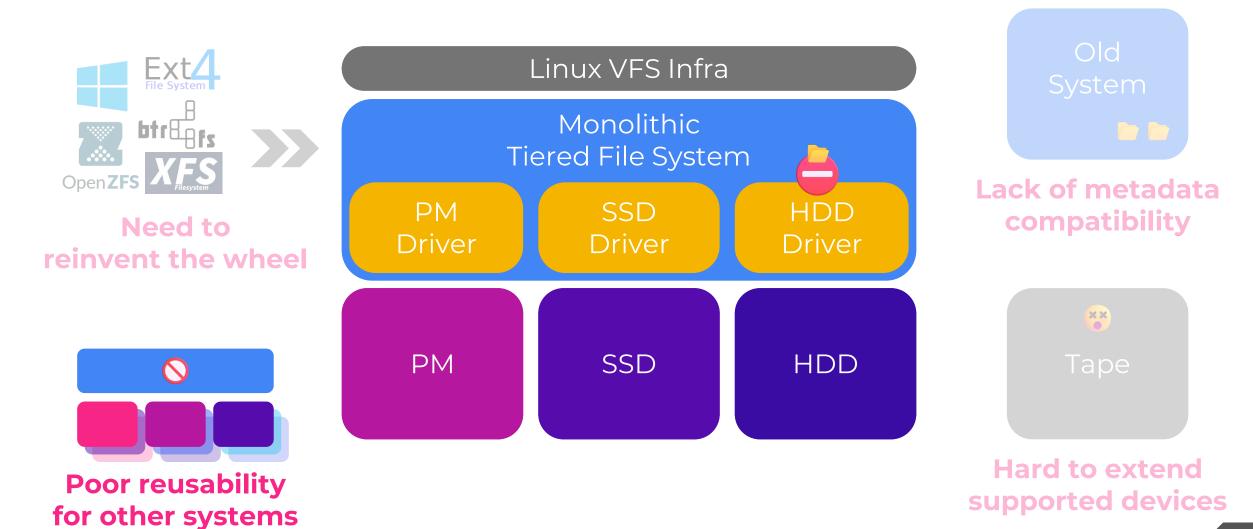


... and it introduces practical shortcomings



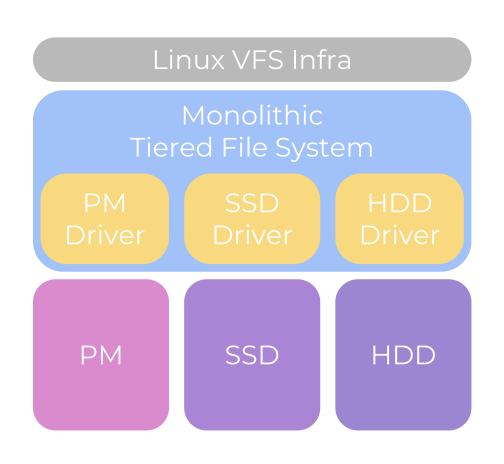


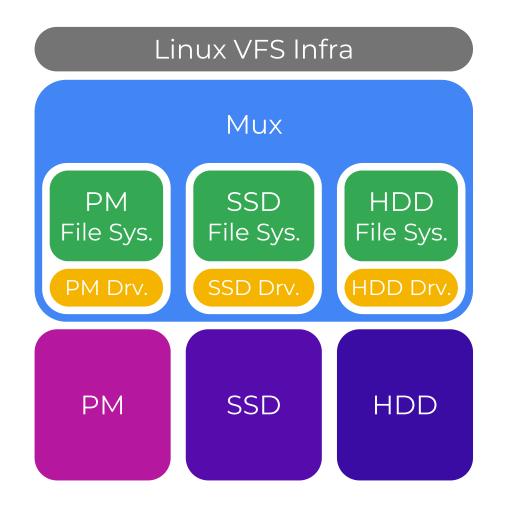
... and it introduces practical shortcomings





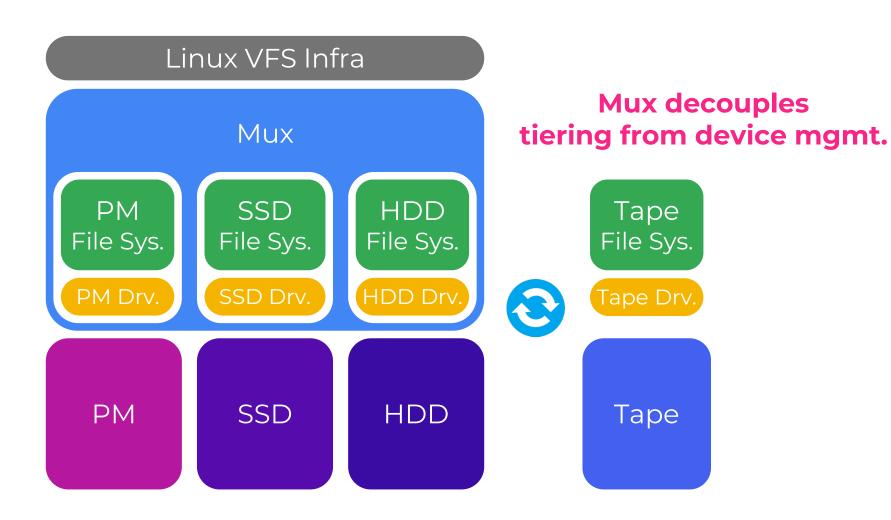
Mux: Talk to file systems, not device drivers







Our Contributions







Our Contributions

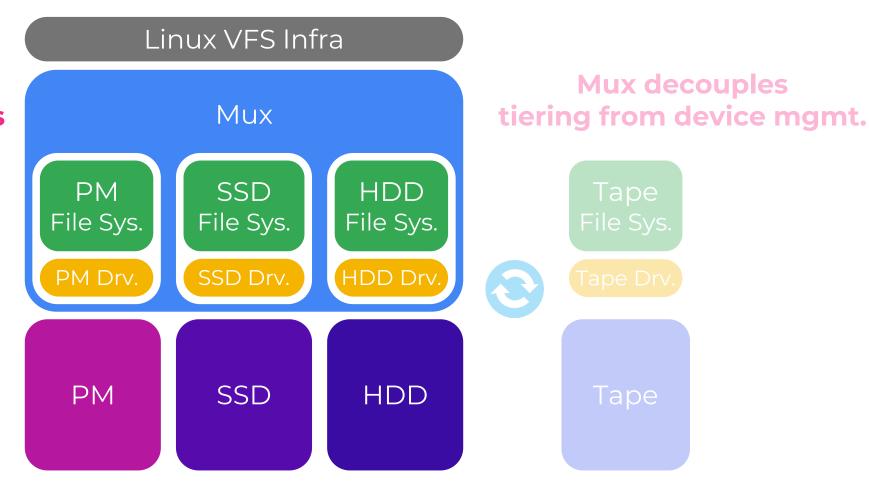
Mux can directly use latest file systems / devices







- Flexibility
- Modularity





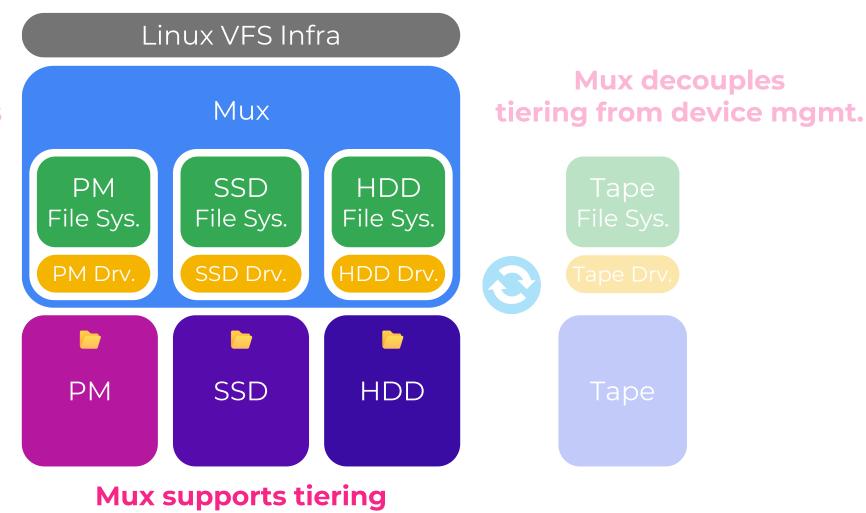
Our Contributions

Mux can directly use latest file systems / devices





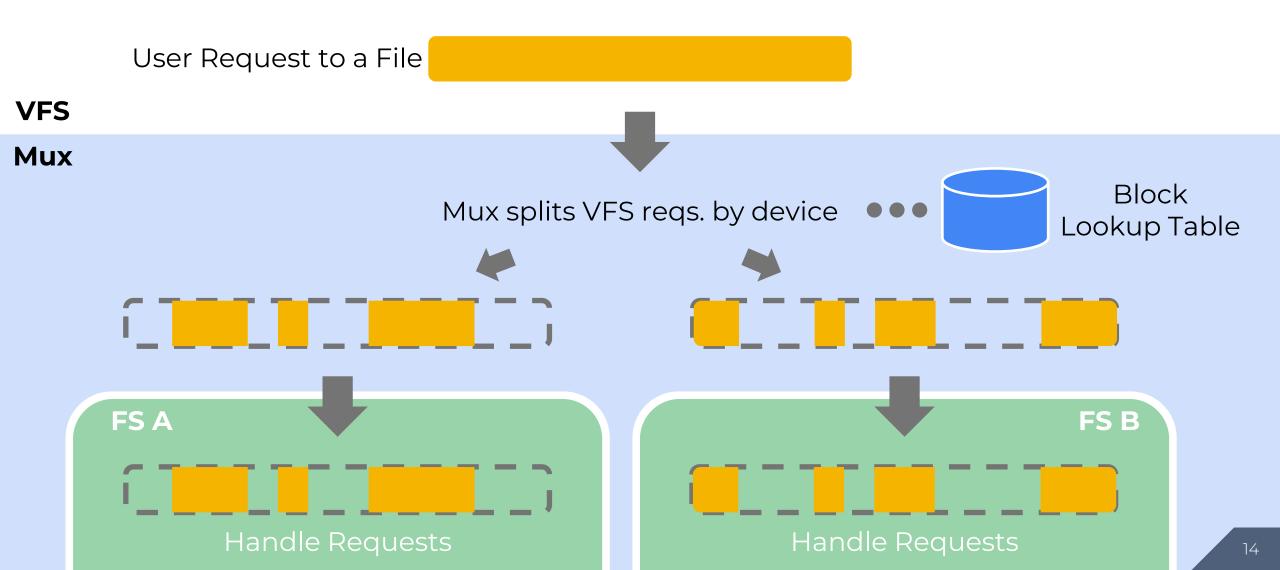
- Flexibility
- Modularity
- Compatibility



over existing storage in production

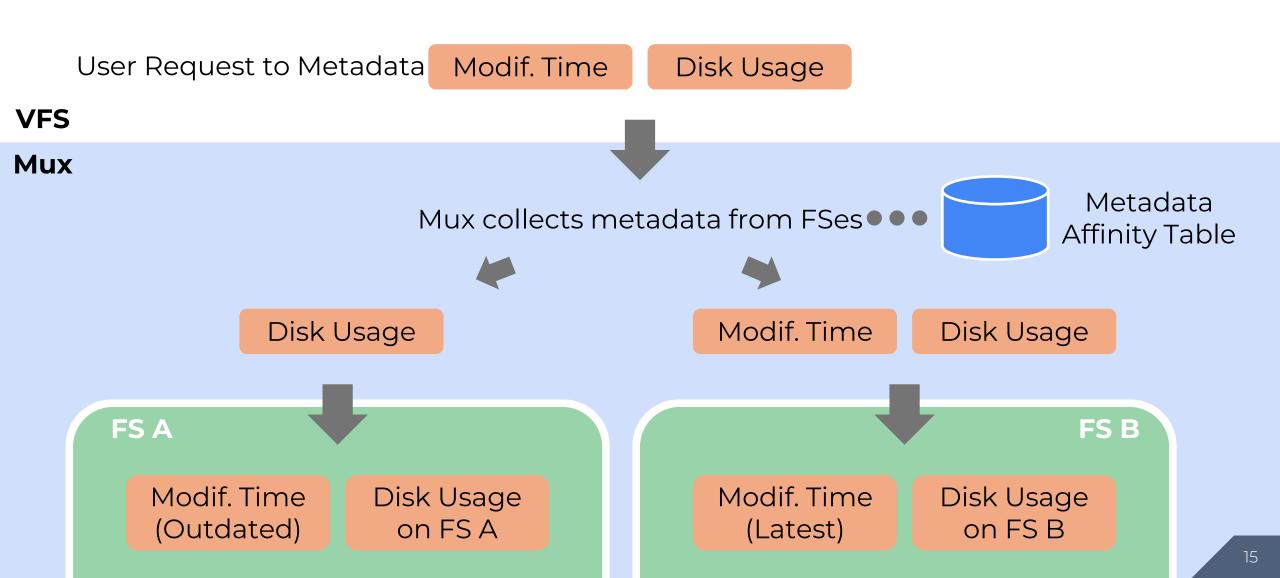


Dispatch user requests, not block I/Os





Track metadata location, not its content





VFS

Mux

File systems' lock scope only applies to themselves internally

FS A

Lock Scope

Lock Scope

FS B



VFS

Mux

Mux cannot utilize file system locks due to potential deadlock with user locks





VFS

Mux

Mux uses optimistic concurrency control for both correctness and performance

FS A

Version Counter

1

FS B

Version Counter

5



VFS

Mux

Mux uses optimistic concurrency control for both correctness and performance

FS A

Version Counter

2

FS B

Version Counter

6



VFS

Mux

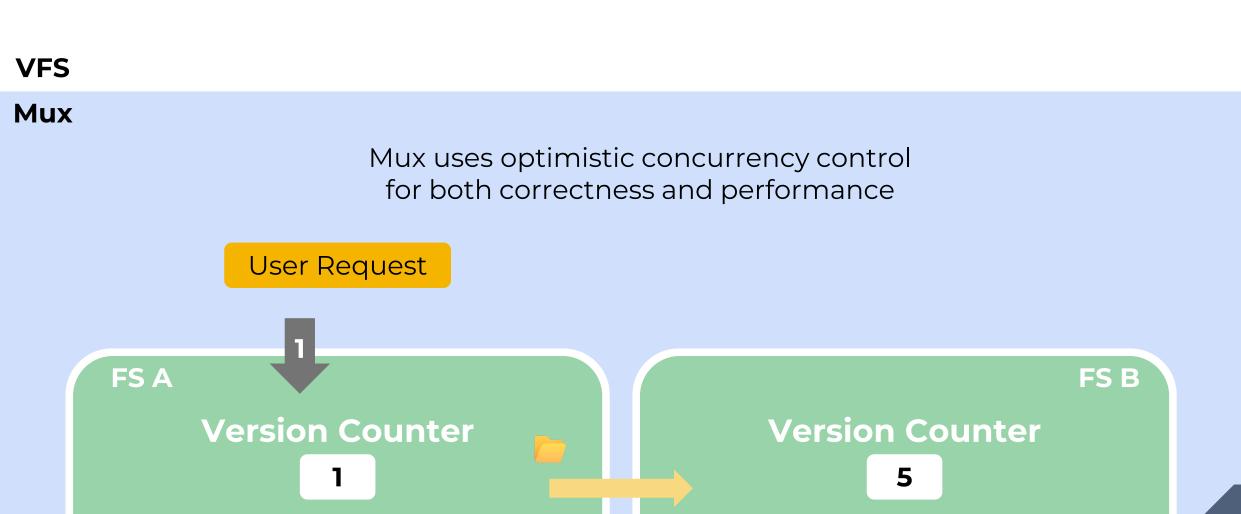
Mux uses optimistic concurrency control for both correctness and performance

FS A

Version Counter

7







VFS Mux Mux uses optimistic concurrency control for both correctness and performance User Request FS A FS B **Version Counter Version Counter** 3



VFS Mux Mux uses optimistic concurrency control for both correctness and performance User Request User Request FS A FS B **Version Counter Version Counter** 2 6

23



Cache on DAX devices, not only DRAM

DRAM Memory

VFS

Mux

FS A

Use Page Cache

FS B



Cache on DAX devices, not only DRAM

DRAM Memory

VFS

Mux



Mux can support page cache with both DRAM and DAX pages





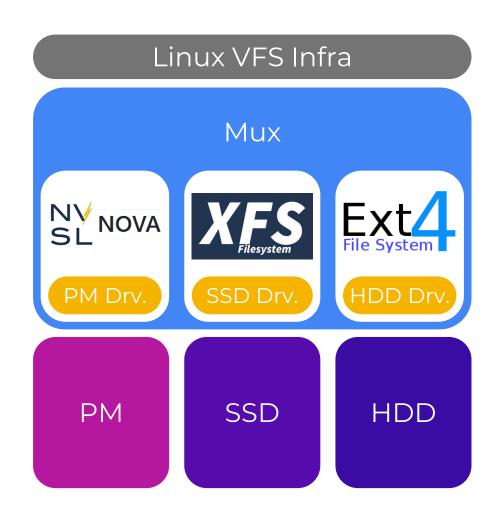
DAX Mem Pool

DAX Enabled





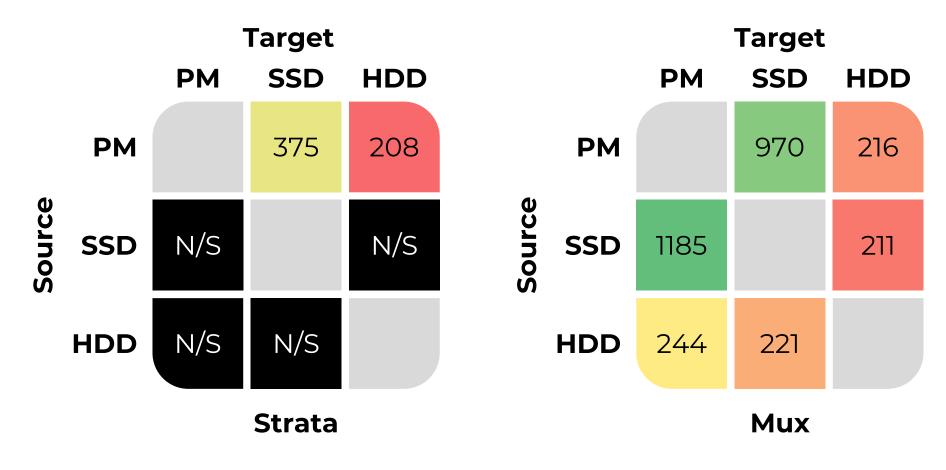
Evaluating a prototype of Mux



- Extensibility = 5 LoC / added device
- Env = Ubuntu 20.04 on Linux 5.15
- Tool = microbench from Strata
- PM = Intel Optane PM 200
- SSD = Intel Optane SSD DC P4800X
- HDD = Seagate Exos X18



Mux brings more tiering possibilities



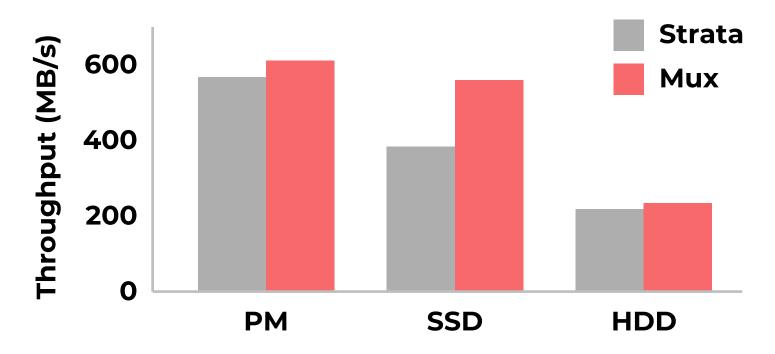
Mux supports all migration directions

Mux enables 2.59x faster migration via FS optimizations

Unit: MB/s N/S: Not Supported



Mux makes high-performance achievable



Mux shows 1.46x higher throughput for device I/O

Mux adds a worst-case read latency overhead of 6.6% to 87.3% compared to non-tiered FSes

Mux adds a worst-case write throughput overhead of 1.6% to 3.5% compared to non-tiered FSes



Future work: Mux can be further improved

- Metadata could be better handled to reduce overhead, e.g., read latency.
- Crash consistency can be enhanced beyond the lowest of all member file systems.
- Feature imparity among file systems could be preserved for better compatibility.

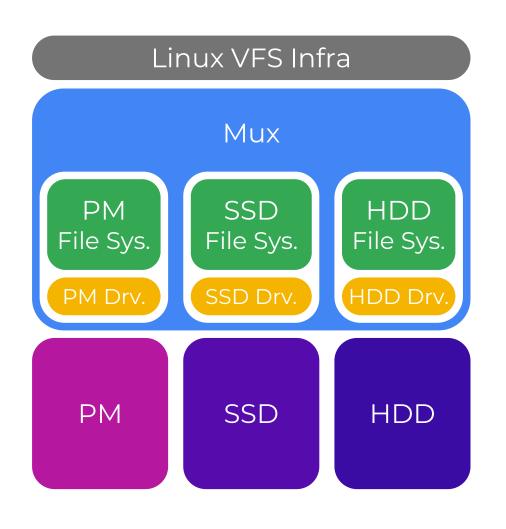


Future work: Mux enables new storage research

- **Scheduling** of file I/O could now be done at VFS level with Mux, rather than block level.
- Configuration of file systems for a given workload or a given set of storage devices can now be studied with Mux.
- **Distributed Mux** can be designed to perform tiering among an array of machines over network.



Conclusion



- Rethinking Tiered Storage: Talk to File Systems, Not Device Drivers
- Mux is a new tiered file system that accesses devices through file systems, rather than device drivers.
 - **Extensibility** as a first-class principle to decouple tiering from device mgmt.
- 1.46x higher throughput for device I/O;
 2.59x faster data migration over Strata.