Erasing Belady's Limitations: In Search of Flash Cache Offline Optimality

Yue Cheng^{1,2} Fred Douglis² Philip Shilane² Michael Trachtman² Grant Wallace² Peter Desnoyers³ Kai Li⁴

¹Virginia Tech ²EMC ³Northeastern University ⁴Princeton University



Background

- Unlike magnetic HDDs, flash SSDs erase data in large blocks
- Limited number of erasures before errors occur

... Want to make effective use of SSDs while being cognizant of erasure penalty







Flash Caching

- Using SSDs for an HDD cache
 - SSD cache as a tier btw. DRAM and HDD
 - Goal: to balance performance against endurance
 - Nitro [USENIX ATC '14], CacheDedup [USENIX FAST '16]
 - RIPQ [USENIX FAST '15]
 - Pannier [ACM Middleware '15]
- How do we know if we're doing well?
 - Comparison against an offline optimal "best case"
 - But what is the offline optimal for flash caches?

\mathbf{W} VirginiaTech \mathbf{EMC}^2



Clients

DRAM

HDDs

Flash cache

0

Offline Optimality

- Belady's MIN: A simple offline caching algorithm when the next access is known
 - Inserts a new entry into the cache, removing the entry that will be used furthest in the future
 - Yields the optimal read hit ratio (RHR)
- MIN is not able to provide the optimal erasures in the context of flash caching

UirginiaTech EMC⁴

MIN inserts data that won't actually be read



Flash Cache Offline Optimality

- Objectives:
 - Õ1: Minimize erasures s.t. maximal RHR
 - Never insert items if it does not increase RHR
 - O2: Maximize RHR s.t. an endurance limit
 - O3: Maximize combination of RHR and erasures
- True optimal
 - How to compute the offline optimal?
 - What is the complexity?
 - Heuristics

UirginiaTech **EMC**²

How can we approximate?

The focus of our work





MIN

Blocks

- Minimum unit of access to the cache (e.g., 4KB)
- MIN priority queue
 (ranked based on next ref.
 Head
 timestamp)
 - Don't insert blocks w/ furthest next ref.

UirginiaTech **EMC**





MIN

Blocks

- Minimum unit of access to the cache (e.g., 4KB)
- MIN priority queue (ranked based on next ref.^H timestamp)
 - Don't insert blocks w/ furthest next ref.
- 2. A new block is inserted WirginiaTech \mathbf{EMC}^2





Key Components

- MIN priority queue
 - In-RAM structure tracking runtime status at block granularity
- Containers
 - Unit of insertion & eviction for flash cache (e.g., 2MB)
- Write buffer

UirginiaTech **EMC**²

- Packs blocks into in-RAM write buffer (e.g., 8MB)
- Container priority queue
 - In-RAM structure tracking flash containers

Container Block Block Block Block Block Block Write buffer Container PQ MIN PQ Northeastern

RAM

Flash

PRINCETON

UNIVERSITY

Block_▲

- Insert a block
 - Into the write buffer (sorted by eviction timestamp)
 - Into the MIN queue (sorted by next ref. timestamp)



B_A: Eviction timestamp: **15**

Next ref timestamp: 10

 \mathbf{W} VirginiaTech \mathbf{EMC}^2

Insert a block

UirginiaTech **EMC**²

- Into the write buffer (sorted by eviction timestamp)
- Into the MIN queue (sorted by next ref. timestamp)



- Write buffer is dispersed into containers when it is FULL
- The containers are written to the flash cache



- Update block status on invalidation or when MIN would evict block as furthest in future ("evictpending")
 - Remains in container until that is GC'd
- Rank container queue by # valid blocks
- Evict the tail container to make room for new data
- Copy forward valid blocks to the write buffer

 \mathbf{W} VirginiaTech \mathbf{EMC}^2



Optimizations

 \mathbf{W} VirginiaTech \mathbf{EMC}^2

- R₁: Only insert blocks read at least once before eviction
- TRIM: Skips dead blocks during GC
 Dead: overwritten or never reaccessed
- Copy-forwarding reduction: Eliminates wasted CF blocks
- E: Segregates blocks by eviction timestamp



Experimental Methodology

- Storage traces (34): limit to large enough datasets
 - EMC VMAX: 25
 - MSR Cambridge: 3
 - MS Production Servers: 6
- Implementation
 - Full-system flash cache simulator
 - Vary cache size as function of unique data accessed in trace: 1-10%
- Metrics
 - Performance: Read hit ratio (RHR)
 - Endurance: Erasures per block per day (EPBPD)
 - Function of RHR and EPBPD: Weighted flash usage effectiveness (WFUE)

UirginiaTech \mathbf{EMC}^2



Comparing Algorithms

All 3 offline algorithms achieve the same (optimal) RHR



X

Comparing Algorithms





Evaluating Heuristic Techniques



Evaluating Heuristic Techniques

10% cache size



Conclusion

- Important to have a baseline for the offline optimal considering both RHR and endurance
- Our container-optimized heuristic maintains the optimal RHR while reducing erasures by up to 67%



 Additional optimizations may be possible to move this heuristic to the *true* optimal



Erasing Belady's Limitations: In Search of Flash Cache Offline Optimality

Thank you! Q & A

Acknowledgments:

Dan Arnon, Cheng Li, Stephen Manley, Darren Sawyer, Dan Tsafrir, Kevin Xu, and Sanjeev Arora and his students from Princeton University



