Traces & code
pdl.cmu.edu/CILES

# **Baleen:** ML Admission & Prefetching for Flash Caches

Center for Coastal Studies

## Daniel L.-K. Wong

wonglkd@cmu.edu

Hao Wu[†], Carson Molder[§], Sathya Gunasekar[†], Jimmy Lu[†], Snehal Khandkar[†]
Abhinav Sharma[†], Daniel S. Berger[‡], Nathan Beckmann, Gregory R. Ganger
[†]Meta, [‡]Microsoft/UW, [§]UT Austin

PARALLEL DATA LABORATORY
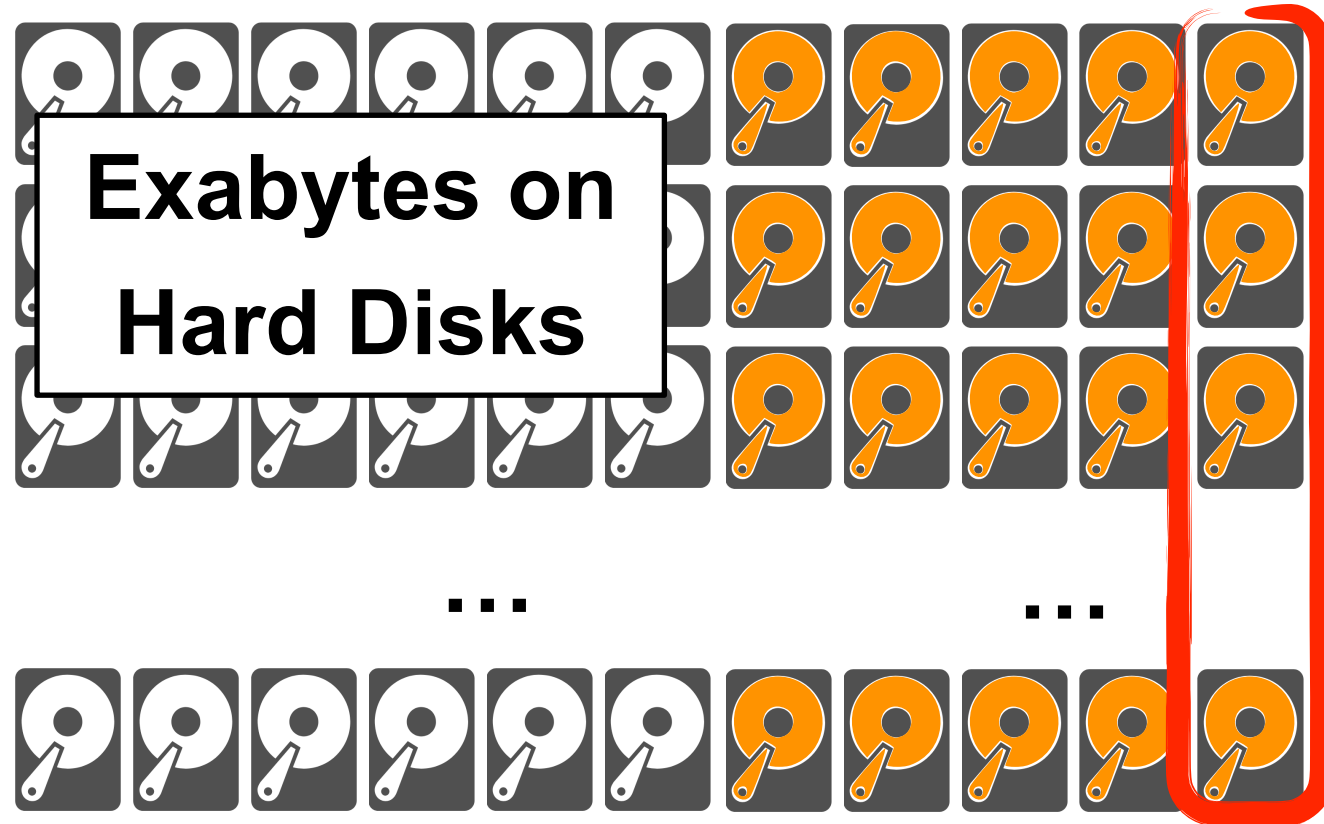Carnegie Mellon University

Carnegie Mellon
Parallel Data Laboratory
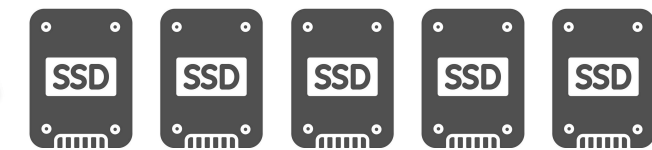
∞ Meta

# Bulk storage systems depend on flash caches

**Bulk Storage**

*(Tectonic, Colossus, …)*
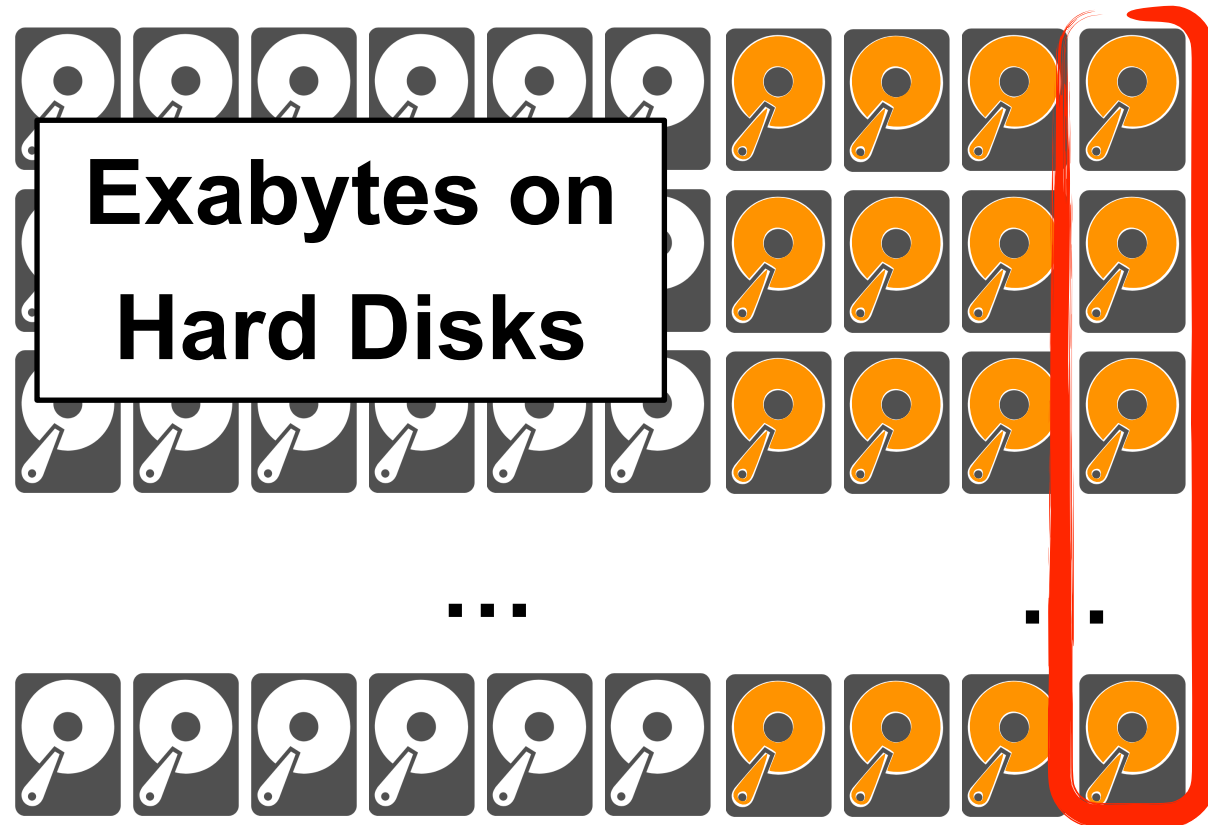
**Flash caches absorb HDD load**

*(CacheLib, …)*

**Exabytes on Hard Disks**

**Extra HDDs needed for IOPS & bandwidth**

# Better flash caches save more HDDs

**Bulk Storage**

*(Tectonic, Colossus, …)*



**Exabytes on Hard Disks**

**…**

…

**Better cache?**

**Extra HDDs needed for IOPS & bandwidth**

**Flash caches absorb HDD load**

*(CacheLib, …)*

Carnegie Mellon
**Parallel Data Laboratory**

# Flash caches are write-heavy

**Bulk Storage**

*(Tectonic, Colossus, …)*

**Flash caches absorb HDD load**

*(CacheLib, …)*

**Exabytes on Hard Disks**

**Better cache?**

**Extra HDDs needed for IOPS & bandwidth**

***Problem:* Limited write endurance**

Carnegie Mellon
**Parallel Data Laboratory**

# Costs dominated by #HDDs & #SSDs

**Baleen reduces costs by 17% on 7 traces**

**Reduce HDD load → less #HDDs**

**Reduce flash writes → less #SSDs**

**Trend: Less IOs/TB**

**Trend: Lower flash endurance**

**Even more important with denser storage!**

Carnegie
**Parallel Data Laboratory**

# How does Baleen reduce costs by 17%?

- 3 key ideas

  - Exploit a new cache residency model (**episodes**)

  - Train ML admission & ML prefetching policies

  - Optimize an end-to-end metric (disk-head time)

- *Why ML over heuristics?*

  - *More savings, more adaptive*

# Bulk storage clients access byte ranges within blocks

**Access**

1. Block ID
2. Byte Range

**Block**



8 MB

**Host**

36 hard disks



**Cluster**

Data center

1000s of hosts

# Fetching bytes from backend causes disk IO

**Client request**

**Backend (HDD)**

IO

*Byte range*

*8MB block*

**Carnegie Mellon**
**Parallel Data Laboratory**

# Cache stores segments (subset of block)

**Client request**

**Flash cache**

**Backend (HDD)**

...

2

3

4

5

IO

...

*Byte range*

*128KB segments*

*8MB block*

# Cache hits save disk IO

**Client request**

**Flash cache**

...

| 2 |
| 3 |
| 4 |
| 5 |

...

**Byte range**

**Backend (HDD)**

**Legend**

| Hit |
| Miss |

***128KB segments*** SSD

***8MB block***

**Carnegie Mellon**
**Parallel Data Laboratory**

# Cache miss causes disk IO

**Client request**

**Flash cache**
...

2
3
4
5
...

**Backend (HDD)**

IO

*Byte range*

*Legend*   Hit   Miss

*128KB segments* SSD

*8MB block*

**Carnegie Mellon**
**Parallel Data Laboratory**

# Decompose flash caching into 3 decisions

**Goal: Reduce HDD load without excessive flash writes**

**Policy Decisions:**

**Flash cache**

...

(a) **Admit misses?** | 3 | 4 | **Baleen**

(b) **Prefetch?** | 5 | **Baleen**

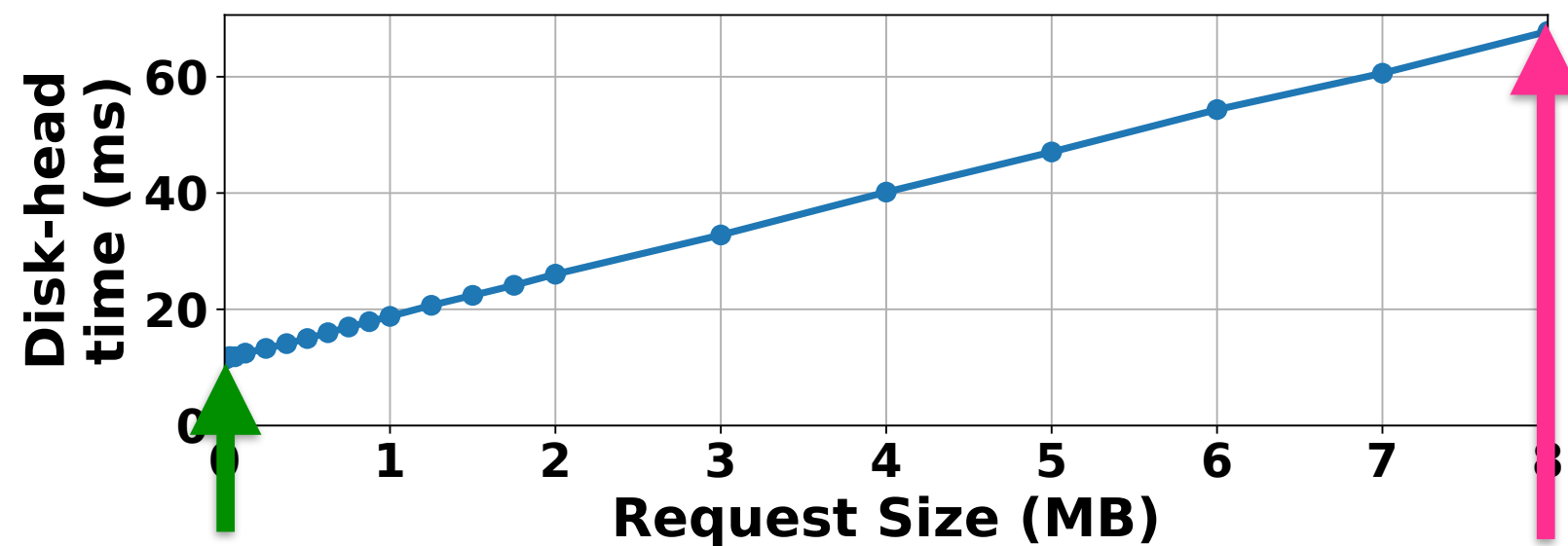(c) **When to evict?** *LRU*

2
3
4
5

...

# Our metric: Disk-head time (DT)

- *Q:* Why DT instead of miss rates?
  - *A:* **Variable size IOs** (reducing **#IOs** & **Size of IOs** both important)
  - Using only **IO hit rate** or **byte miss rate** is an easy misstep (we did!)
- ***DT = Positioning time + Read time***



constrained by **IOPS**                    constrained by **bandwidth**

- *Intuition:* DT is weighted sum of **#IOs** & **#Bytes**

# Design
## Episodes model

**Carnegie Mellon**
**Parallel Data Laboratory**

# ML for caching not straightforward

**Typical supervised learning**

- e.g., "Is this picture a cat?"

**ML for Caching**

- Data: trace of accesses

- *Multiple related decisions: Admit now? Later? Never?*

  - *Depend on AND affect cache contents, future decisions*

- **Tend to overfit on easy decisions**

- **Underfit on examples at margin that distinguish policies**

**Training on accesses non-trivial**

Carnegie Mellon
Parallel Data Laboratory

# What is an episode?

**Episode:**

Group of accesses corresponding to
the block's residency in flash
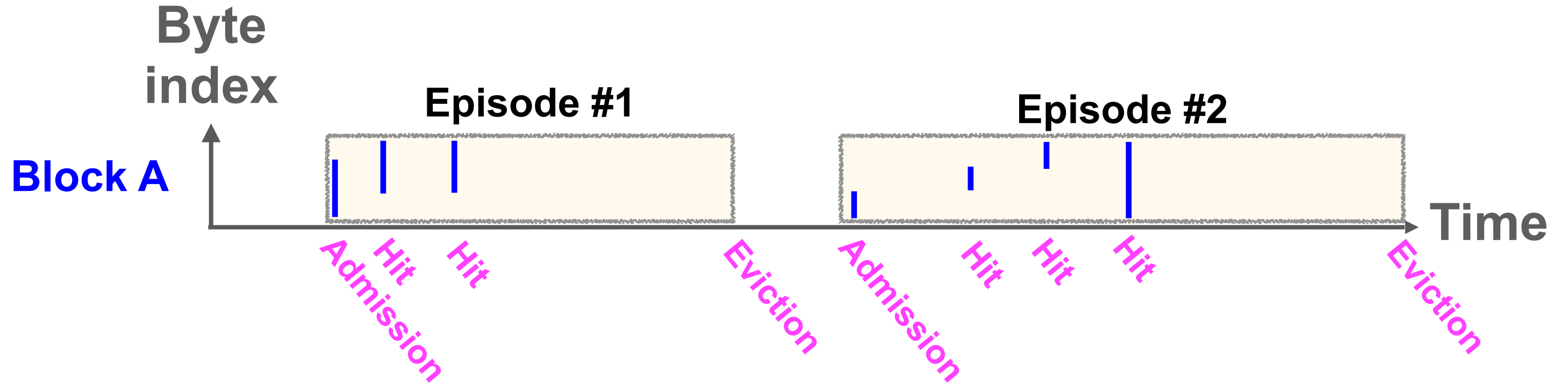if you admitted it on the 1st access
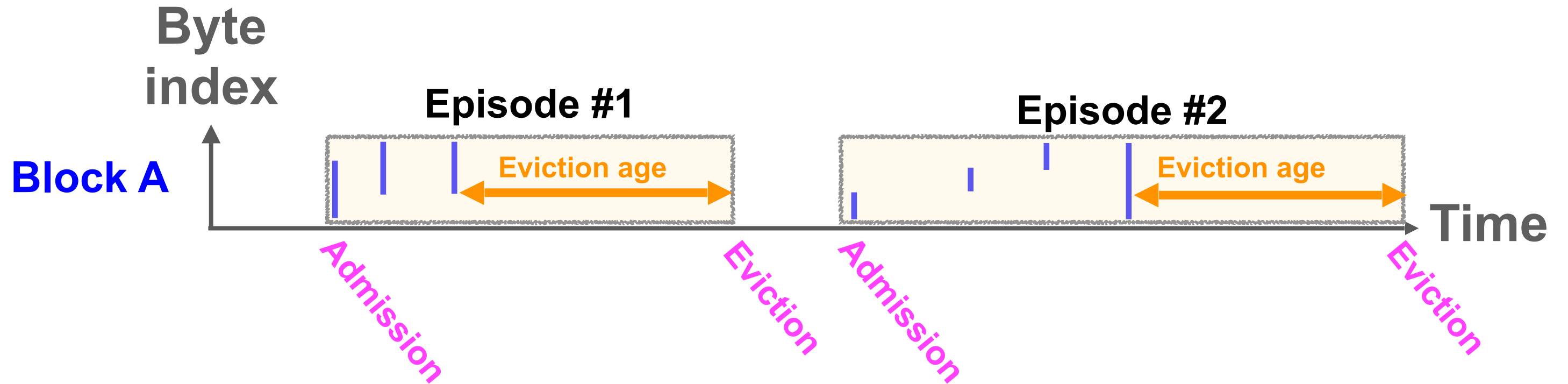
# Why use **episodes** to train ML?

- **Right granularity**

  - Focus on <u>first</u> access instead of all accesses

  - Policies see misses, not accesses

- **Right examples**

  - Avoid overfitting on popular blocks with many accesses but only 1 miss

- **Right labels**

  - Costs & benefits defined on admission to eviction

**Carnegie Mellon**
**Parallel Data Laboratory**
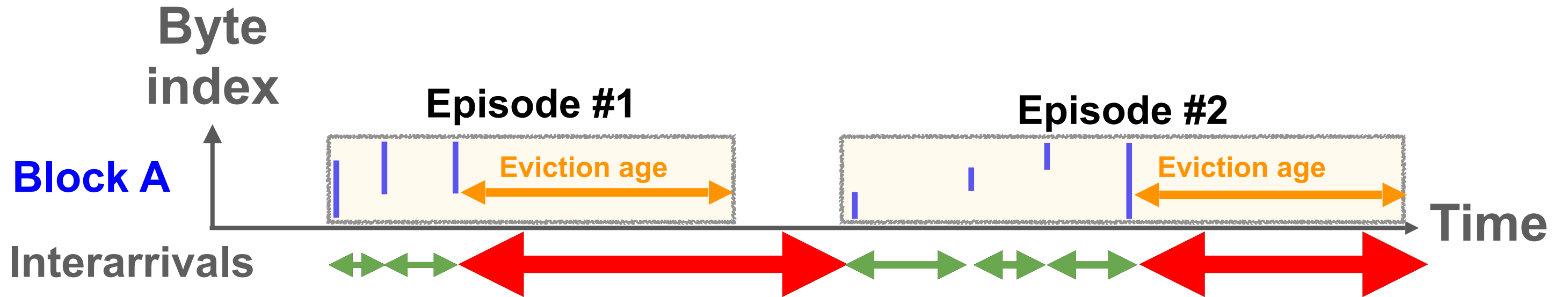
# Episodes: from admission to eviction

# How to know when eviction happens?



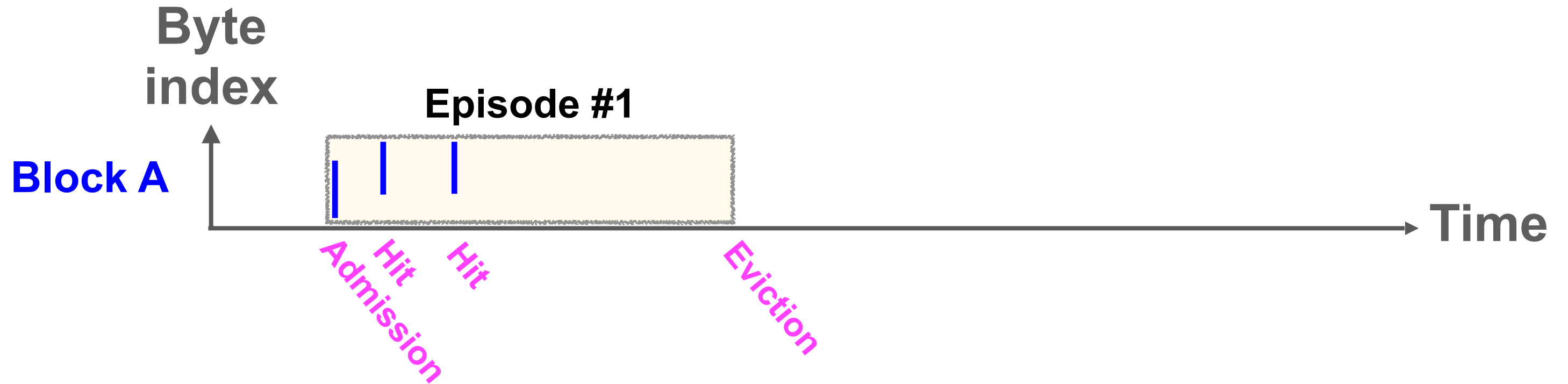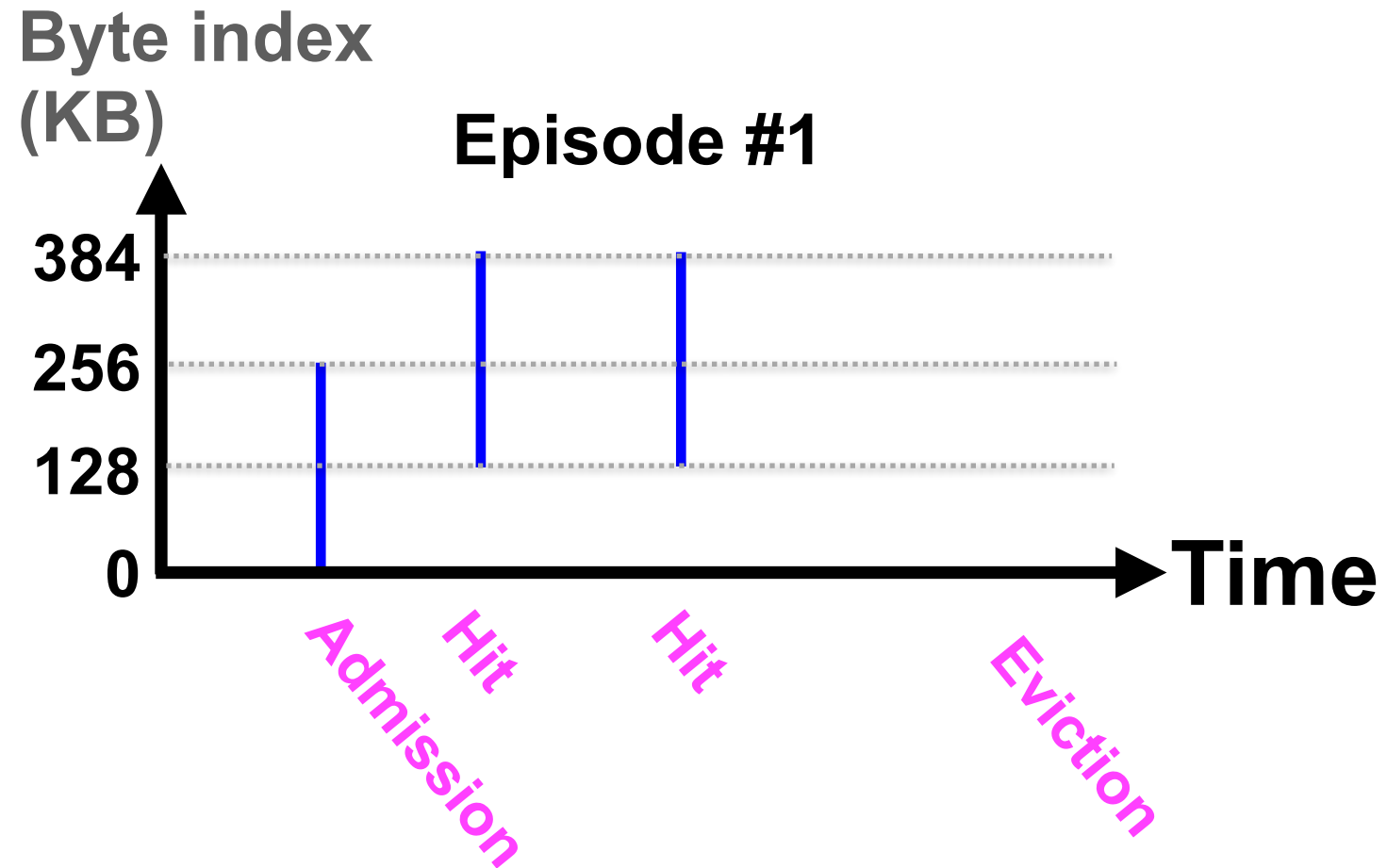**How: model LRU cache state with assumed eviction age**

# How episodes are generated



**Consider interarrival times of accesses**

**Split into episodes when interarrival > eviction age**

# Focusing on Episode 1...

Baleen: ML Admission & Prefetching for Flash Caches                    pdl.cmu.edu/CILES

# Reason about episodes instead of accesses

# Benefits & costs defined on episodes

**Episode #1**

**Misses, Prefetches, Hits**



- **Benefit:** 27ms of DT saved

- **Cost:** 3 flash writes needed

**Carnegie Mellon**
**Parallel Data Laboratory**

# Design
Using episode-based policies to answer "What does good look like?"

# Admission: Baleen learns from episode-based OPT

**OPT** (approx. optimal) admits highest scoring episodes

$$\text{Score}(Ep) = \frac{\text{DTSaved}(Ep)}{\text{FlashWrites}(Ep)} = \frac{27 \text{ ms}}{3 \text{ flash writes}}$$ *Episode #1*
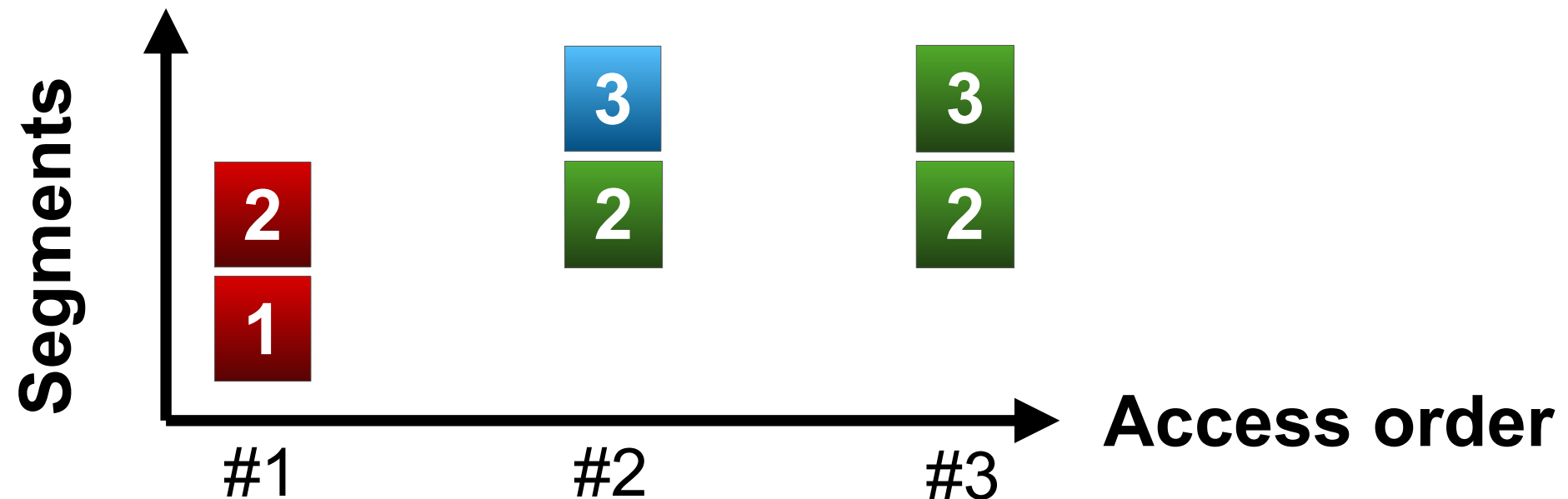
**OPT** emits binary labels based on flash write budget

Yes if $\text{Score}(Ep) > \text{Cutoff}_{\text{TargetFlashWriteRate}}$

**Baleen imitates OPT admission**

**Carnegie Mellon**
**Parallel Data Laboratory**
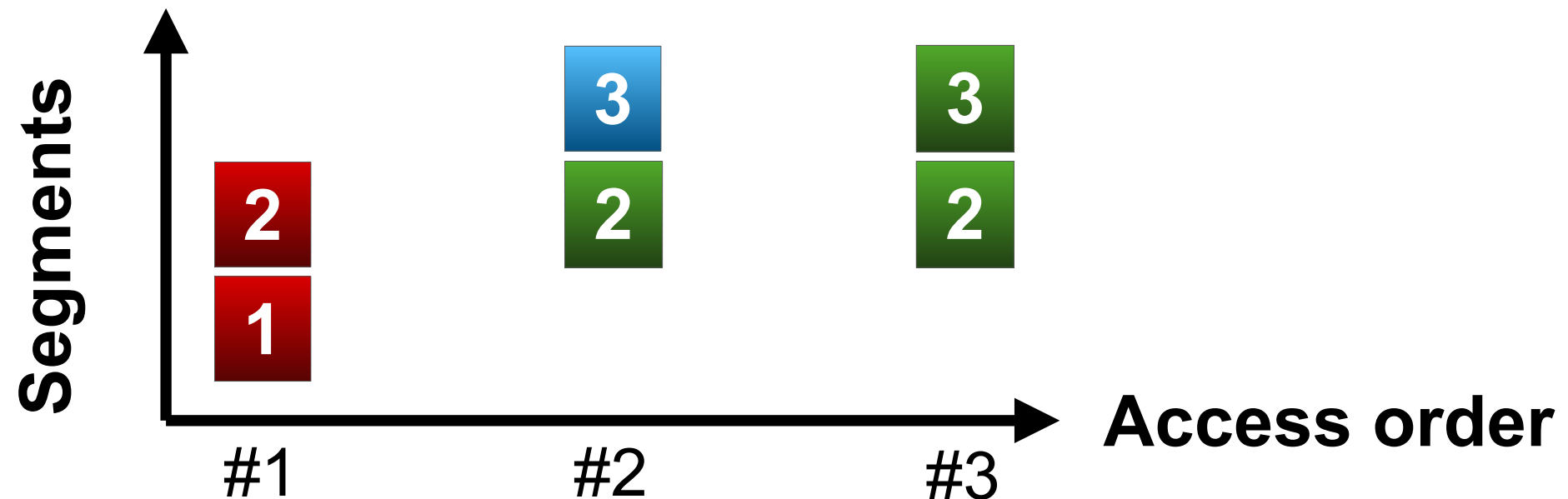
# Baleen's ML-Range learns **what** to prefetch

- **What** range to prefetch
  - OPT-Range Start: lowest segment
  - OPT-Range End: highest segment
- **ML-Range** is trained on OPT-Range

$$\text{OPT-Range}(Ep) = [1, 3]$$

**Carnegie Mellon**
**Parallel Data Laboratory**

# Baleen's ML-When learns **when** to prefetch

- **When** to prefetch
  - Bad prefetching hurts: wasted DT & cache space
  - Prefetch only when confident of benefits
  - **ML-When**: Yes <span style="color:red">if</span> $\mathrm{PrefetchBenefit(Ep)} > \epsilon$

# Baleen-TCO balances HDD savings against SSD cost

- Q: How to balance #HDD against #SSDs?

**HDD cost** + **SSD cost**

*Measure*　　　　　　　　　　　　　*Vary*

$$\mathbf{TCO_1} \propto \frac{\boxed{\text{PeakDT}_\mathbf{1}}}{\text{PeakDT}_0} \cdot \#\text{HDDs}_0 + \frac{\text{Cost}_\text{SSD}}{\text{Cost}_\text{HDD}} \cdot \frac{\boxed{\text{FlashWR}_\mathbf{1}}}{\text{FlashWR}_0} \cdot \#\text{SSDs}_0$$

- Baleen-TCO picks optimal flash write rate
  - for each workload

*TCO function based on Google's CacheSack [Yang23]

# Evaluation

- Production workloads from Meta's Tectonic

  - 7 clusters from 3 years (2019, 2021, 2023)

  - Each serves 1-10 tenants, e.g., data warehouse

  - Each tenant serves 100s of applications

- *More details on Tectonic in Pan et al (FAST 2021)*
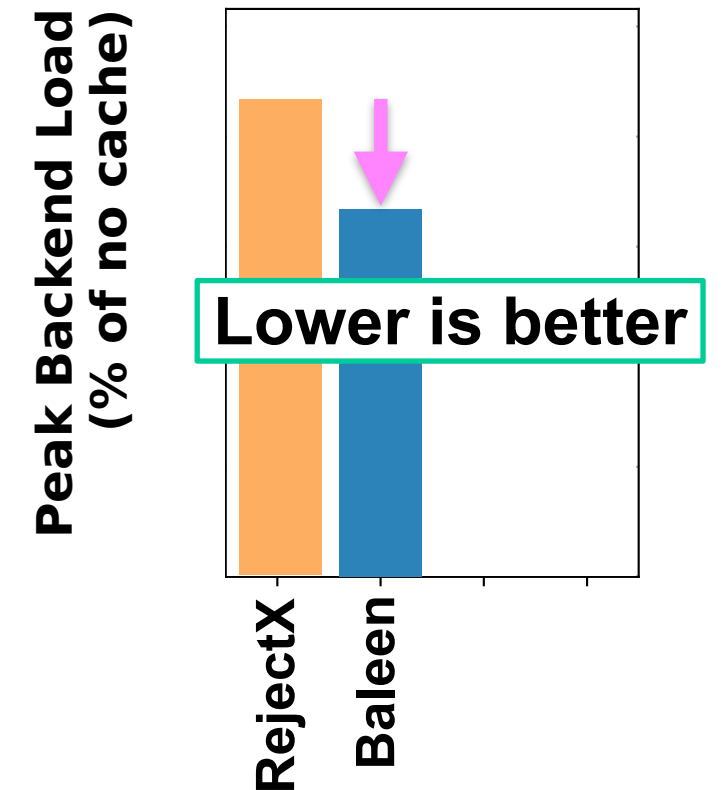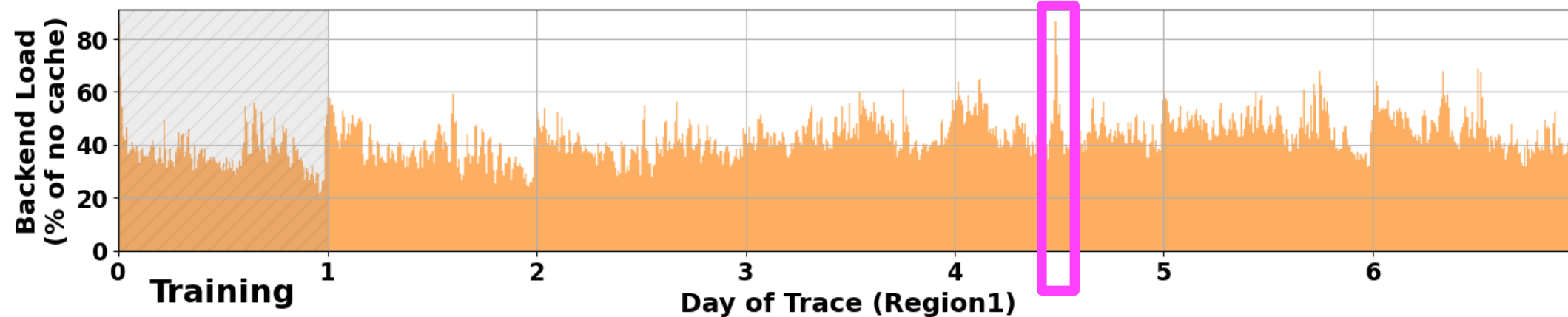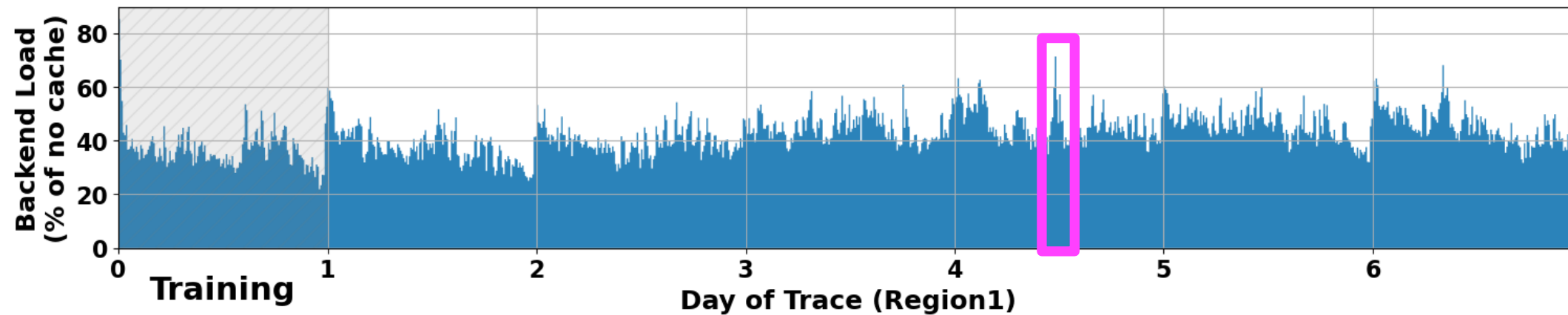
- *Traces & simulator code released*

# Baseline admission policies

- ## CoinFlip: flip a coin for each IO
  - ### Simplest, requires no state

- ## RejectX (e.g., X=1: accept segment after 1 reject)
  - ### Used by Meta, Google as baseline
  - ### 2nd access is always a miss

- ## CacheLib-ML
  - ### Used by Meta in production for 3 years
  - ### Trained on accesses, not episodes

**Admitted**    **Hit**

**Segments**

**Access order**

**Misses, Hits**

**Carnegie Mellon**
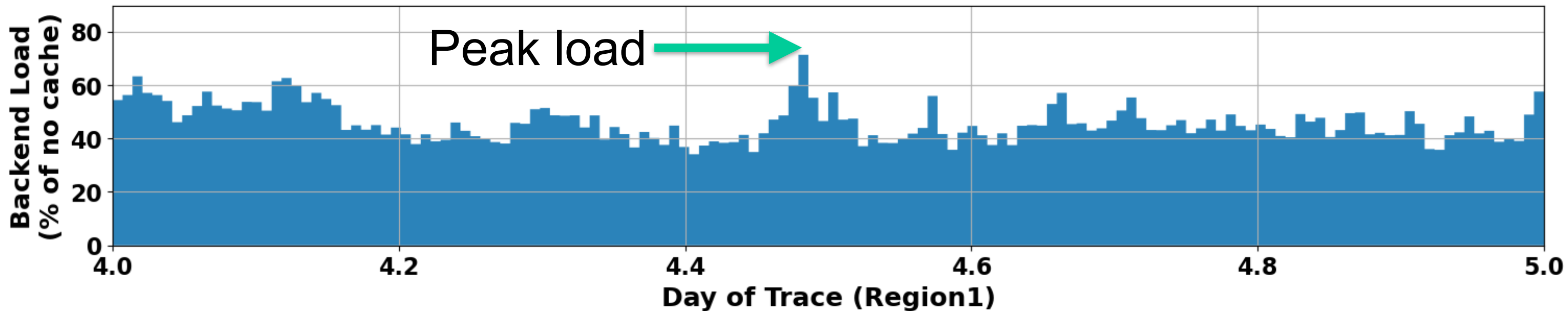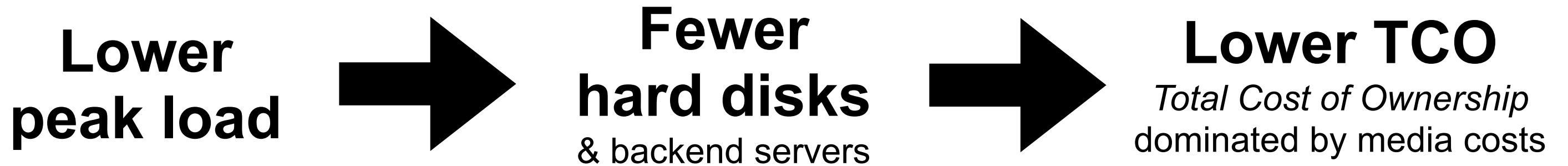**Parallel Data Laboratory**

# Minimize peak backend load to minimize cost

- We train (offline) on Day 1 and evaluate on Day 2-7
- We compare policies' **Peak DT** (as a % of no caching)

**Carnegie Mellon**
**Parallel Data Laboratory**

# Reduce peak load to lower total cost

**Lower peak load** → **Fewer hard disks**
& backend servers → **Lower TCO**
*Total Cost of Ownership*
dominated by media costs



Peak load →

**Carnegie Mellon**
**Parallel Data Laboratory**

# Baleen saves most cost



Estimated TCO (% of RejectX)

Carnegie Mellon
**Parallel Data Laboratory**

# Baleen saves most cost



**17% savings in TCO**

*In paper:* **14% less IO misses**

Legend:
- CoinFlip
- RejectX
- CacheLib-ML
- Baleen-TCO

Y-axis: Estimated TCO (% of RejectX)

**Carnegie Mellon**
**Parallel Data Laboratory**

# Prefetching accounts for most benefit



**Prefetching saves most**

Estimated TCO (% of RejectX)

Legend:
- CoinFlip
- RejectX
- CacheLib-ML
- Baleen (No Prefetch)
- Baleen
- Baleen-TCO

# Prefetching depends on good admission decisions

- Choice of admission policy matters

  - ML prefetching makes admission baselines worse

- Even with ML admission, 2 models required

  - ML-Range to know what to prefetch

  - ML-When to select when to prefetch

# Conclusion

- Baleen reduces cost by 17%

- Episodes guide ML training

- Optimize for Disk-head Time metric

- Smart admission & prefetching

  - ML-Range predicts what to prefetch

  - ML-When estimates confidence in ML-Range

- Ongoing work: workload drift mitigation

  - **Seeking longer traces with features! (>1 week)**

The Boy and the Big Blue Whale
Dr Rose Wadenya, Maria Andrieieva

Questions / data / jobs? wonglkd@cmu.edu / www.cs.cmu.edu/~dlwong

**Carnegie Mellon**
**Parallel Data Laboratory**