

Telescope

Telemetry for Applications with Gargantuan Memory Footprints

Alan Nair (presenting)

Sandeep Kumar, Aravinda Prasad, Ying Huang
Andy Rudoff, Srinivas Subramoney

Intel Labs



Background

Terabyte-scale Applications



Machine Learning
LLMs



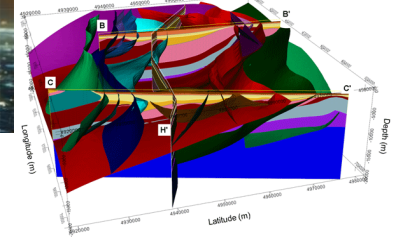
Amazon
ElastiCache

In-Memory Databases



Weather
Prediction

Seismic
Modeling

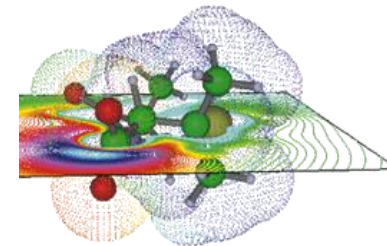


High-Performance Computing



PaLM2

Google's Next-Generation
Large Language Model

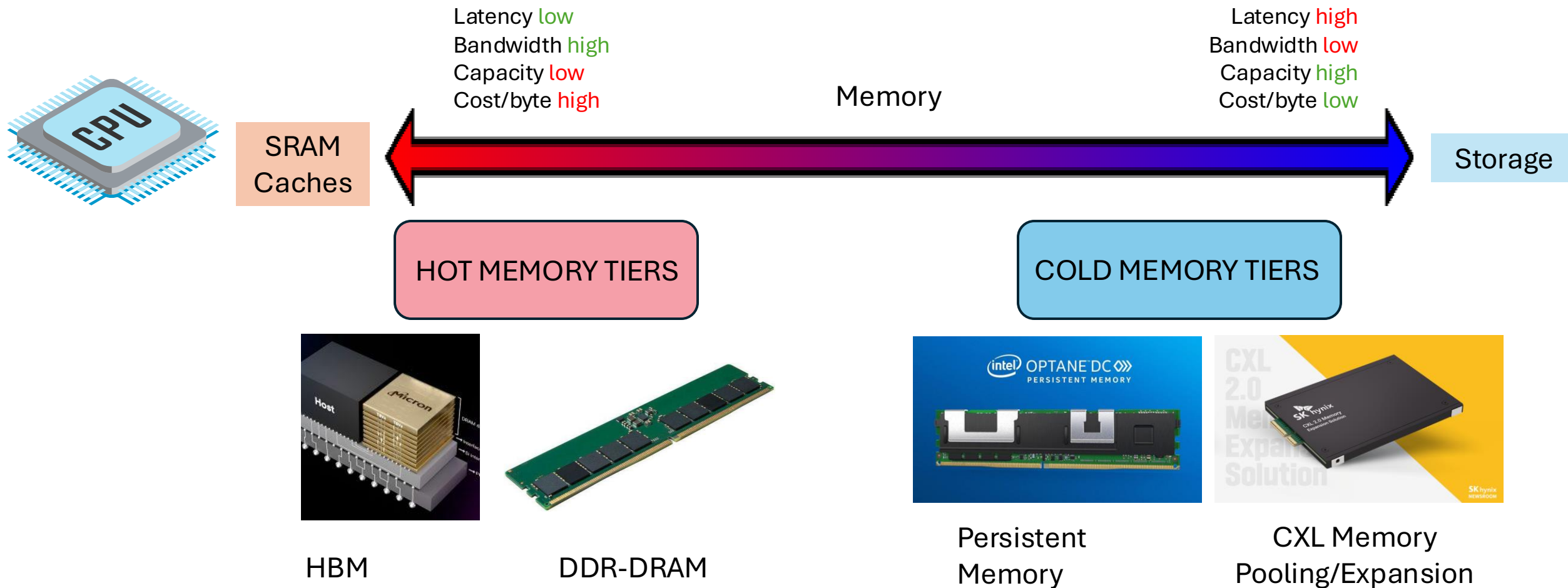


Manufacturing
Materials
Simulations

Drug Development
Cancer Analysis

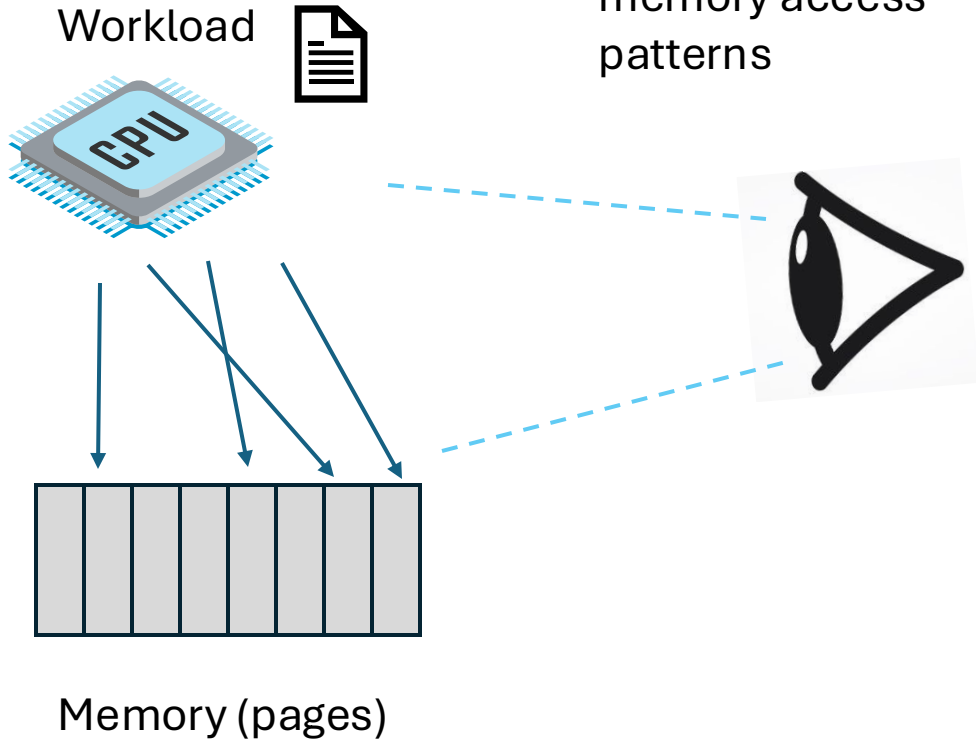


Terabyte-scale Memory Systems

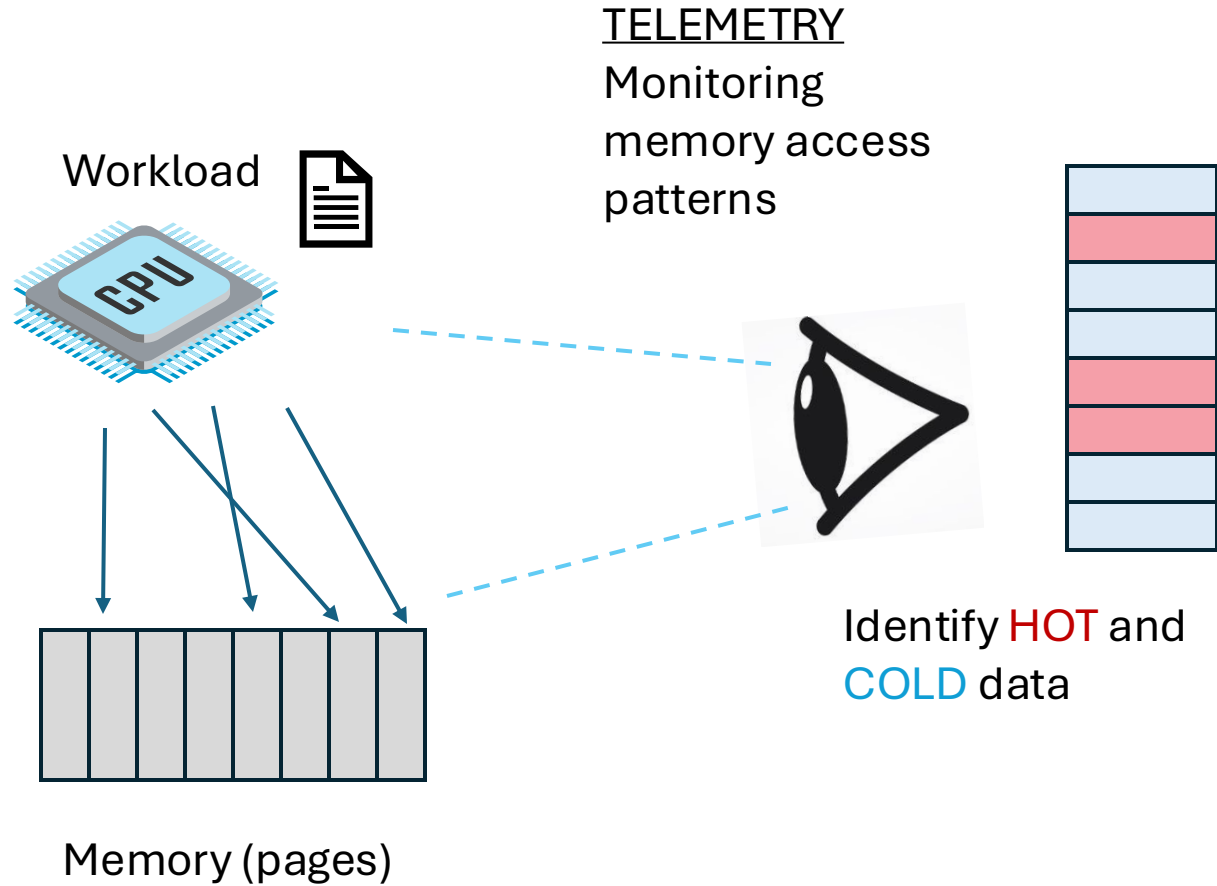


Tiered Memory needs Telemetry

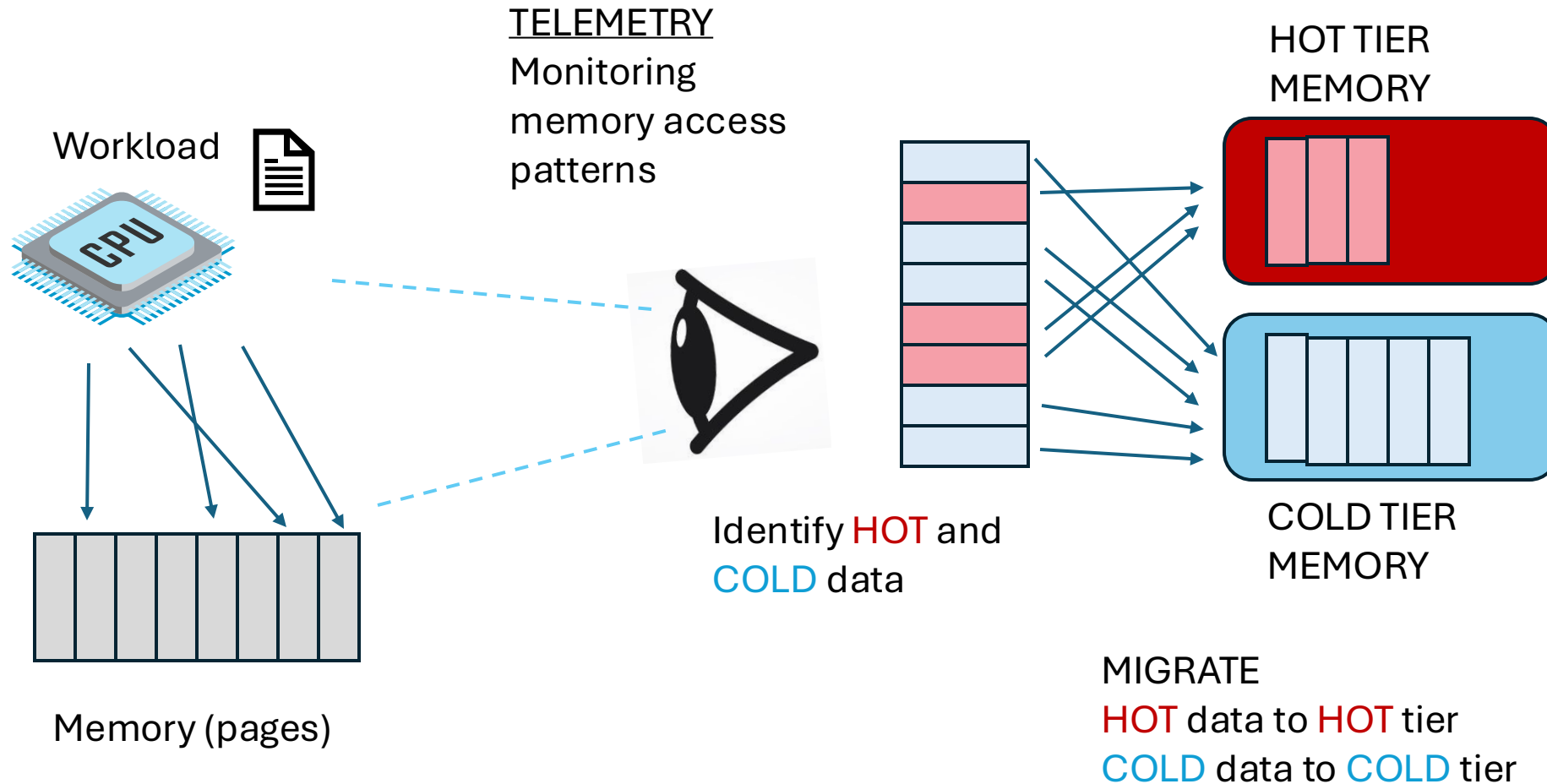
TELEMETRY
Monitoring
memory access
patterns



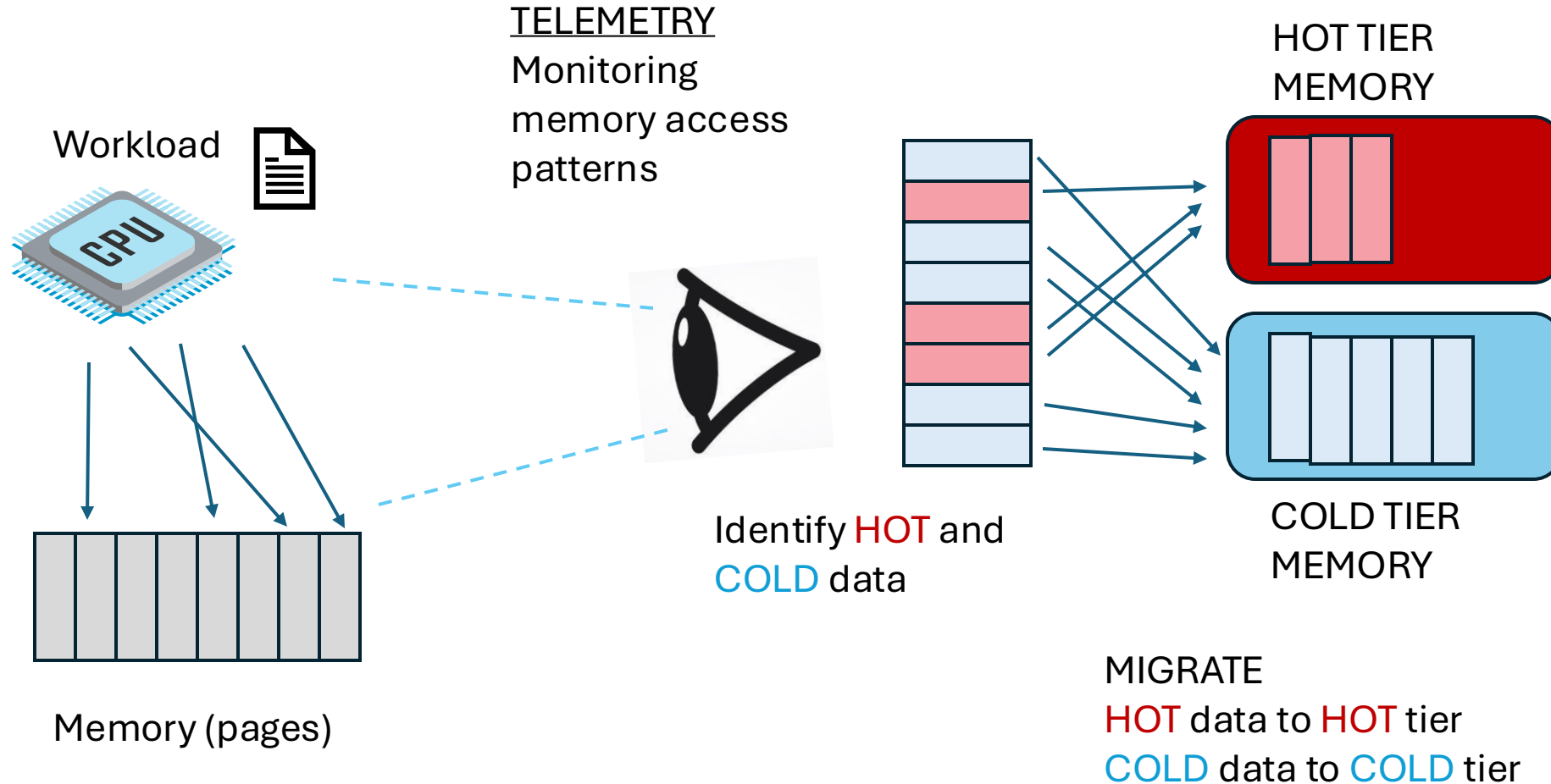
Tiered Memory needs Telemetry



Tiered Memory needs Telemetry



Tiered Memory needs Telemetry



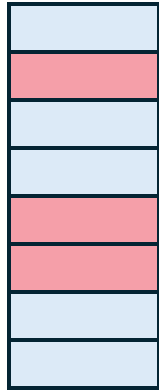
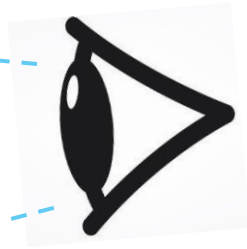
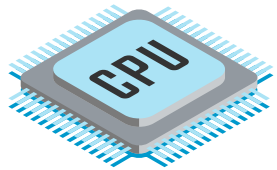
MAXIMUM PERFORMANCE !!! 

Data Placement needs Telemetry

TELEMETRY

Monitoring
memory access
patterns

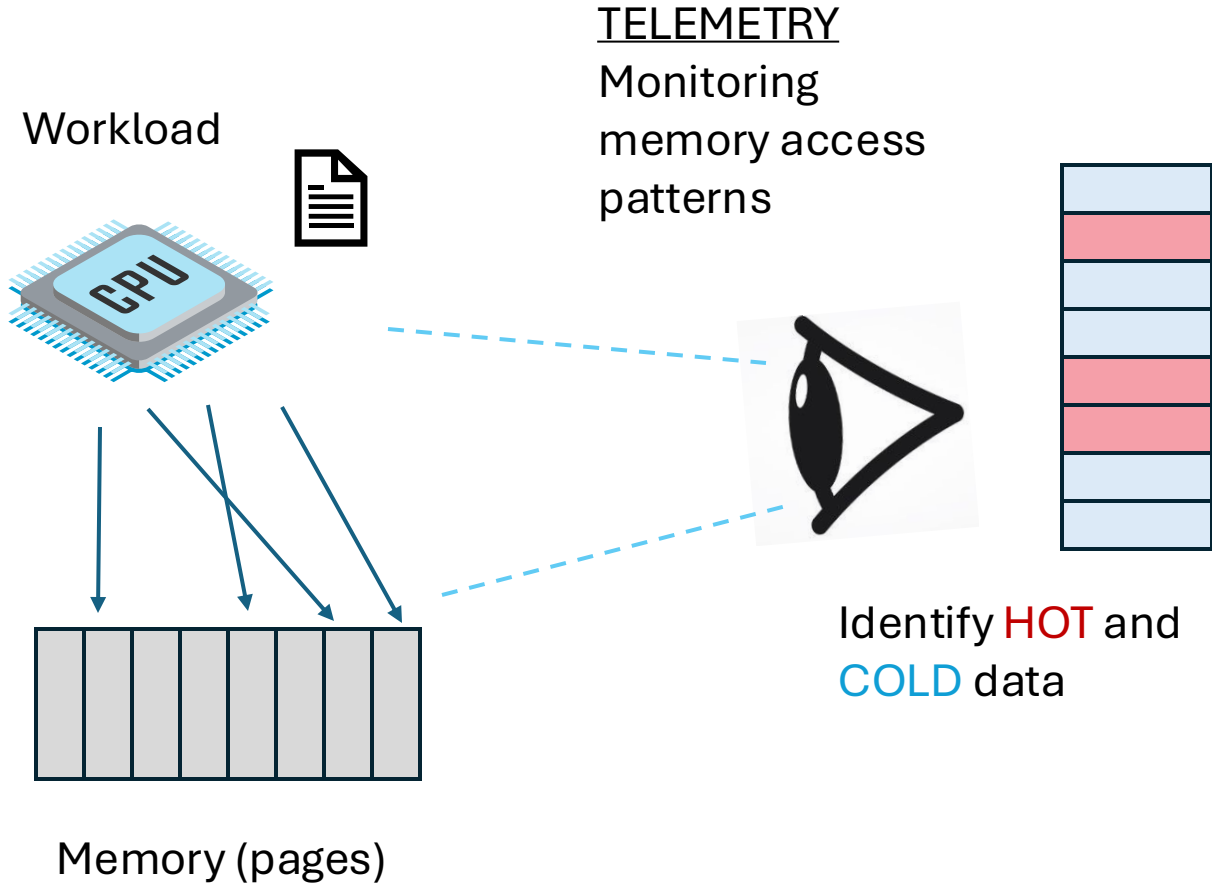
Workload



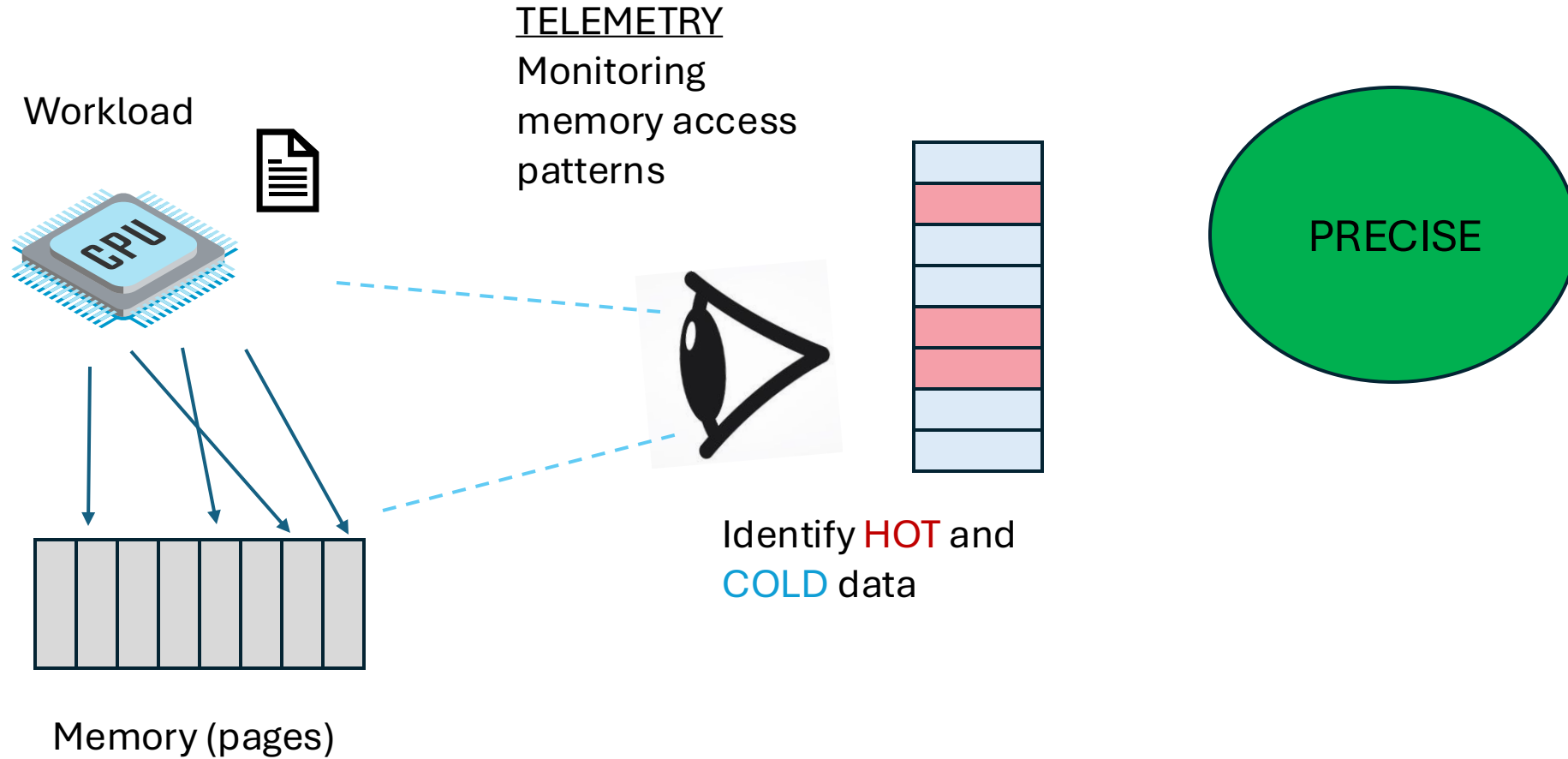
Identify **HOT** and
COLD data



Memory (pages)



Data Placement needs Telemetry

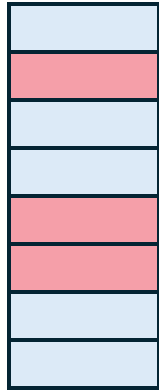
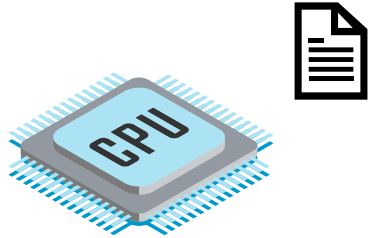


Data Placement needs Telemetry

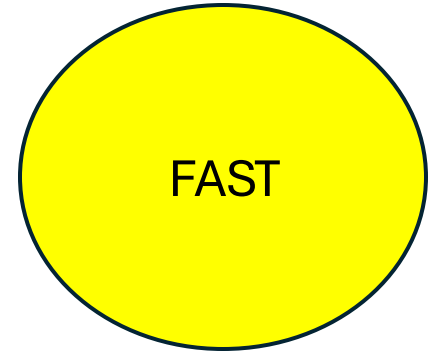
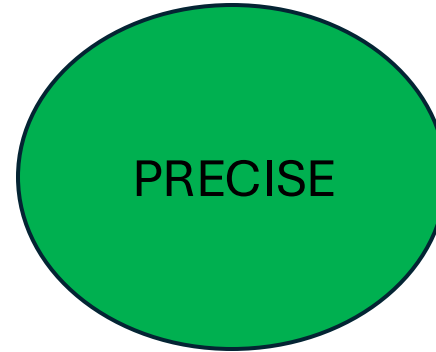
TELEMETRY

Monitoring
memory access
patterns

Workload

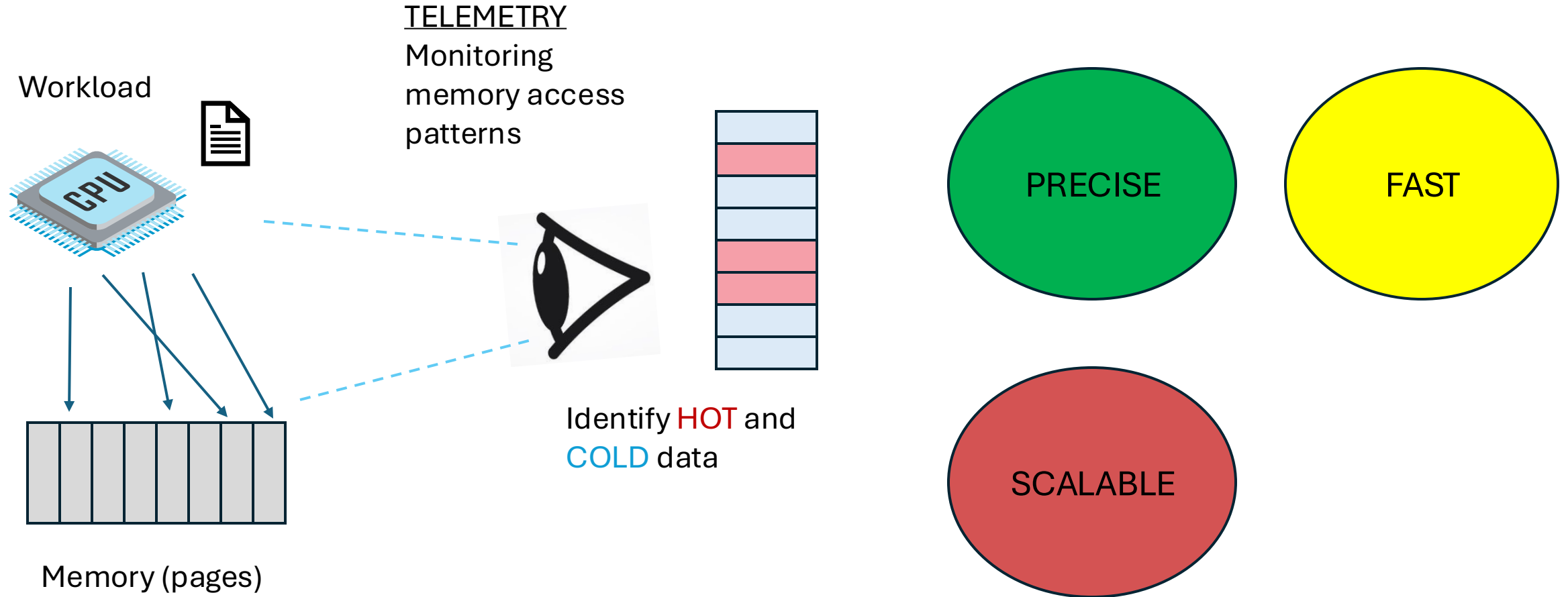


Identify **HOT** and
COLD data

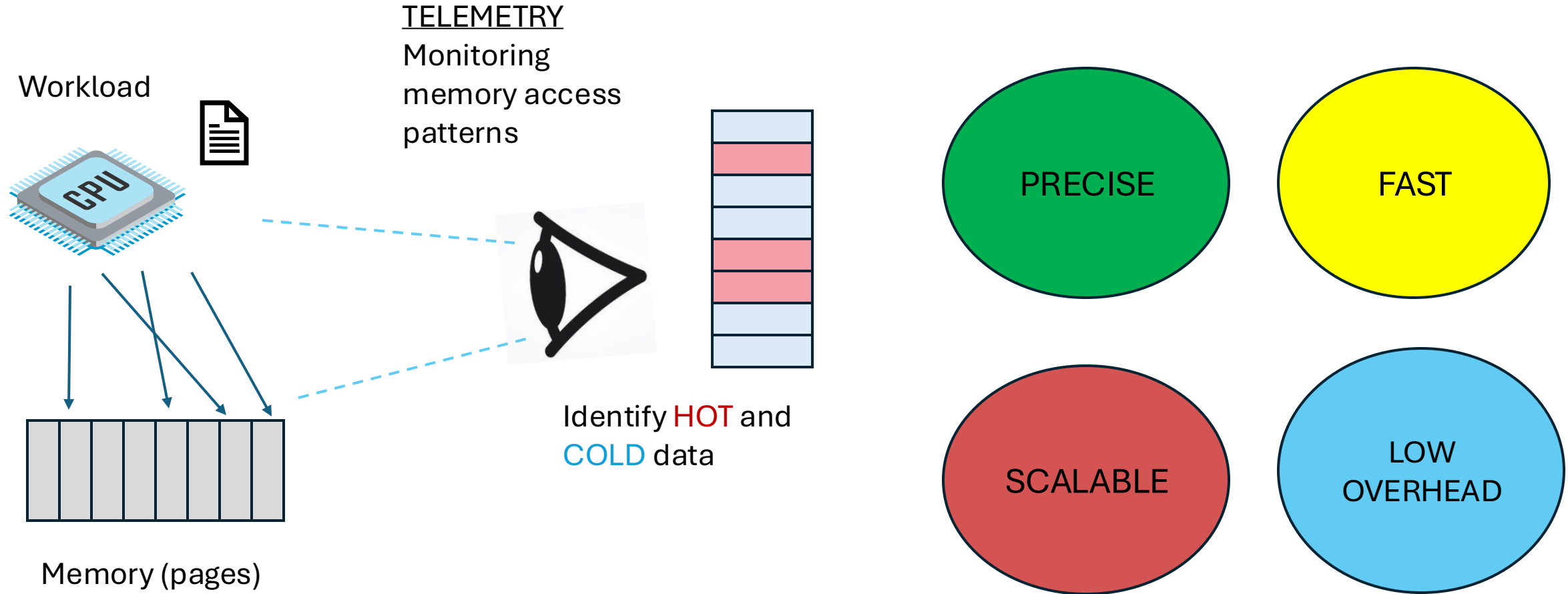


Memory (pages)

Data Placement needs Telemetry



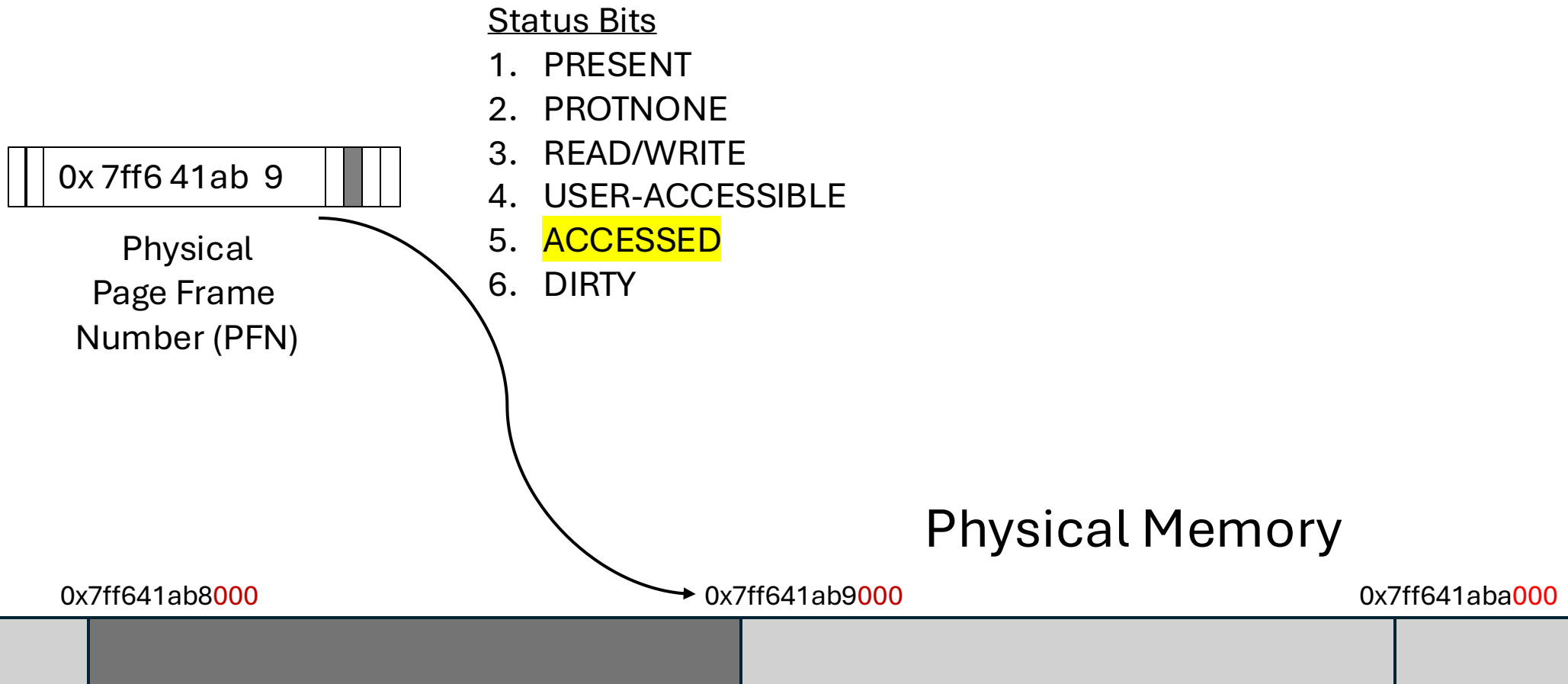
Data Placement needs Telemetry



Prior Approaches

1. LINEAR SCANNING
OF PAGE TABLE ENTRIES
(PTEs)
2. REGION-BASED SAMPLING
3. PERFORMANCE
COUNTERS

Linear Scanning of PTEs



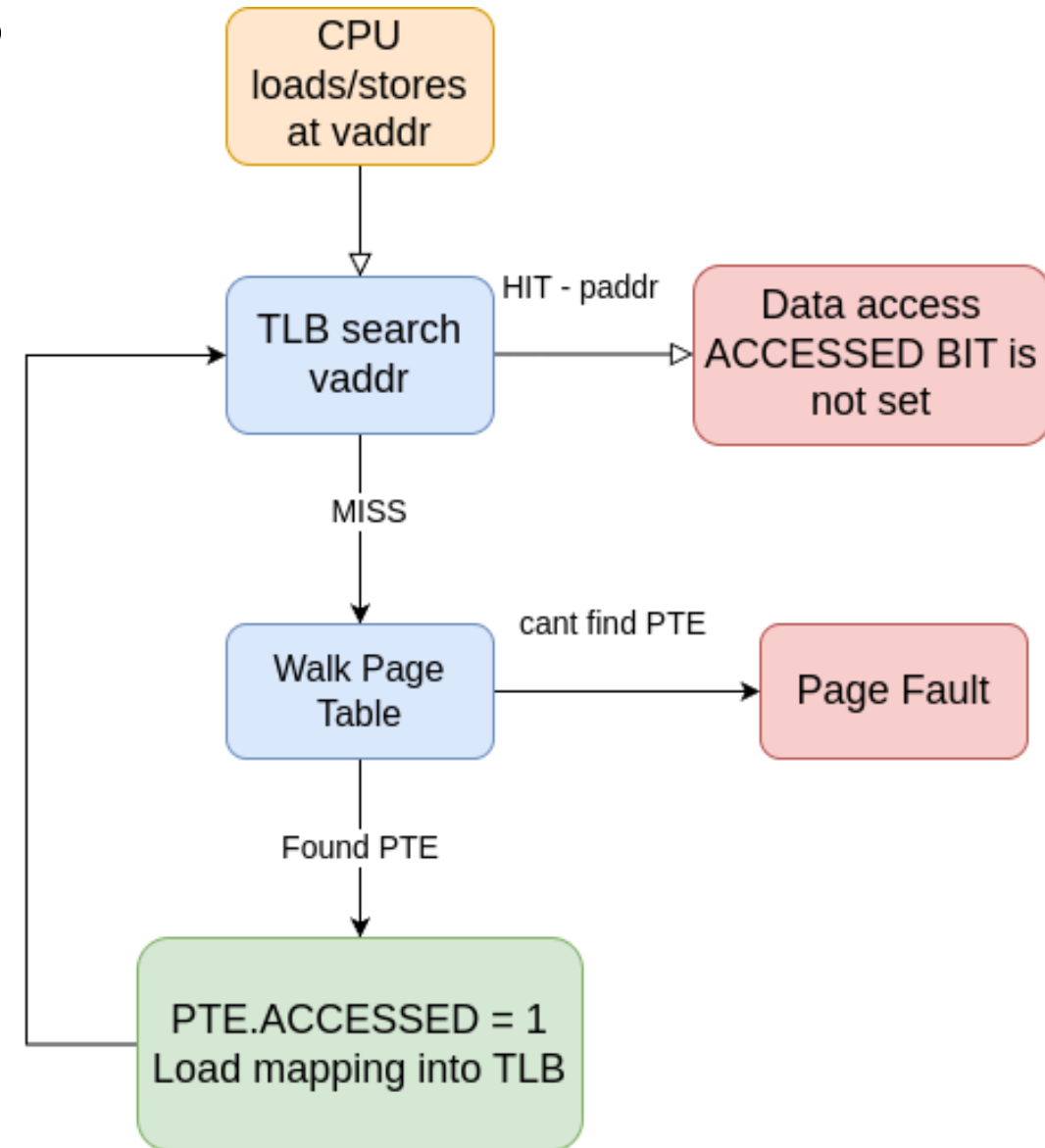
#1 Linear Scanning of PTEs

0x 7ff6 41ab 9

Physical
Page Frame
Number (PFN)

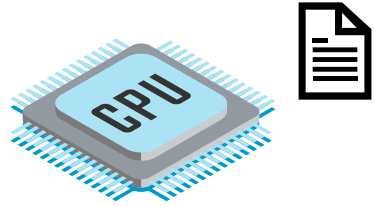
Status Bits

1. PRESENT
2. PROTNONE
3. READ/WRITE
4. USER-ACCESSIBLE
5. **ACCESSED**
6. DIRTY

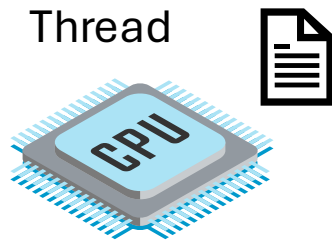


Linear Scanning of PTEs

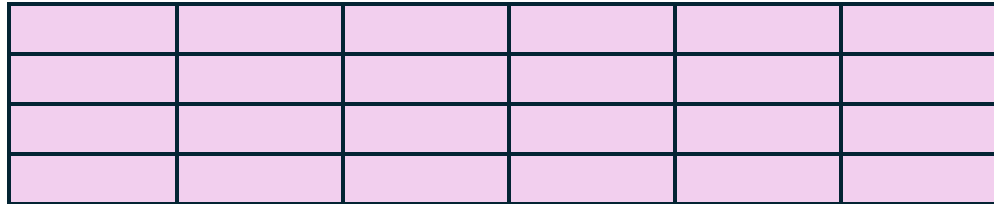
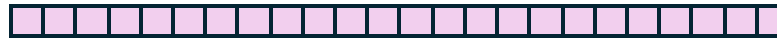
Workload



PTE Scanning
Thread



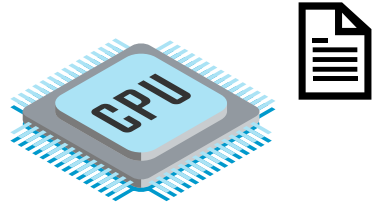
Workload's PTEs



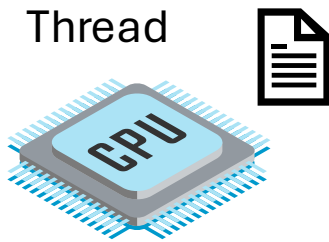
Workload's
pages

Linear Scanning of PTEs

Workload



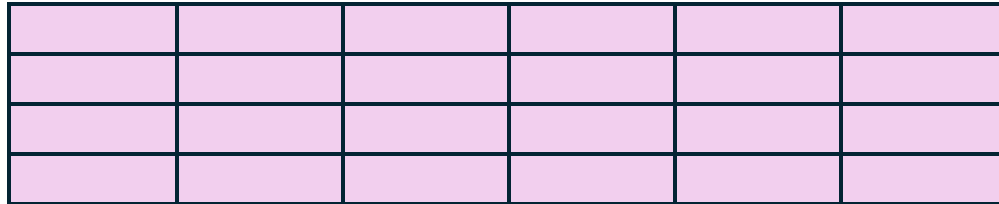
PTE Scanning
Thread



Workload's PTEs



Workload's
pages



Sleep 1 ms, while application continues
Wake up and check all ACCESSED bits



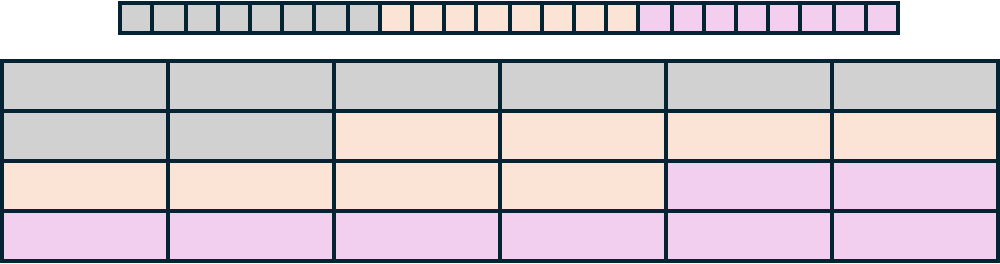
Linear Scanning of PTEs - Limitation

1 TB memory = 256M pages (4K each)
One full scan takes > 1 minute

Does not scale well
with increasing
number of pages !

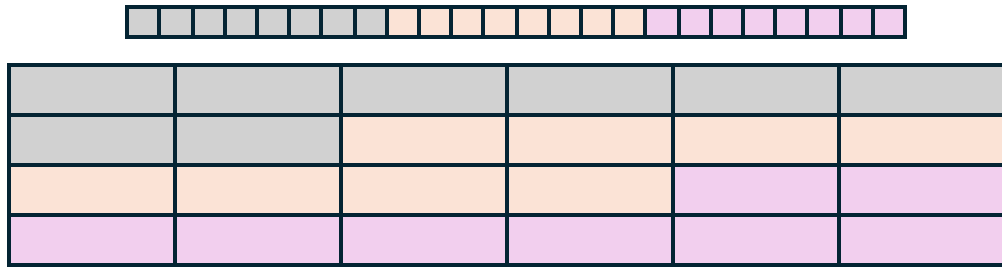
Need Telemetry that
converges to access
pattern in SECONDS

#2 Region-Based Sampling



Workload's Mapped Address Space is divided into Regions.

#2 Region-Based Sampling

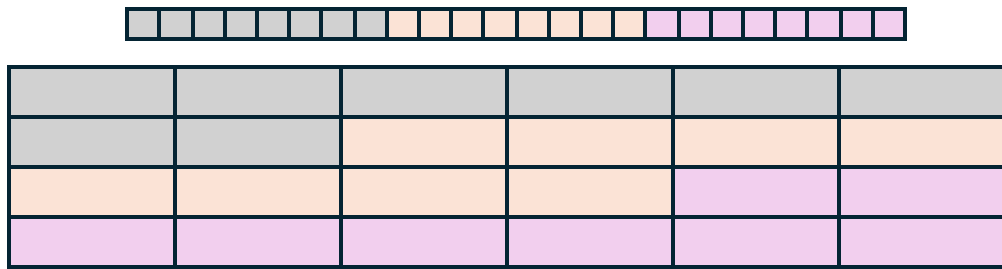


Workload's Mapped Address
Space is divided into Regions.

Randomly pick one PTE per region
RESET these PTEs' ACCESSED bits



#2 Region-Based Sampling



Workload's Mapped Address
Space is divided into Regions.

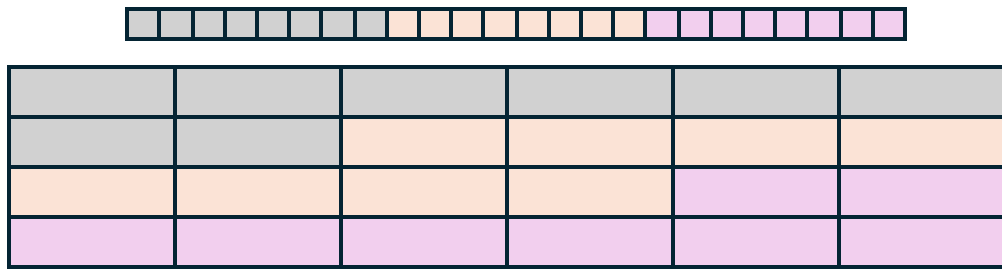
Randomly pick one PTE per region
RESET these PTEs' ACCESSED bits



Sleep 1ms while app executes
Wake up & check ACCESSED bits



#2 Region-Based Sampling

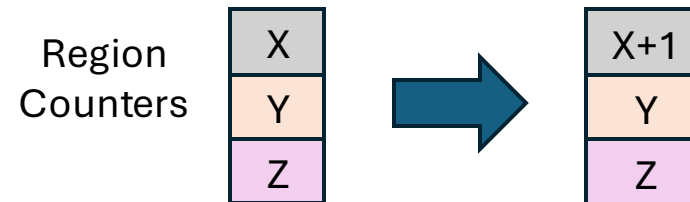


Workload's Mapped Address Space is divided into Regions.

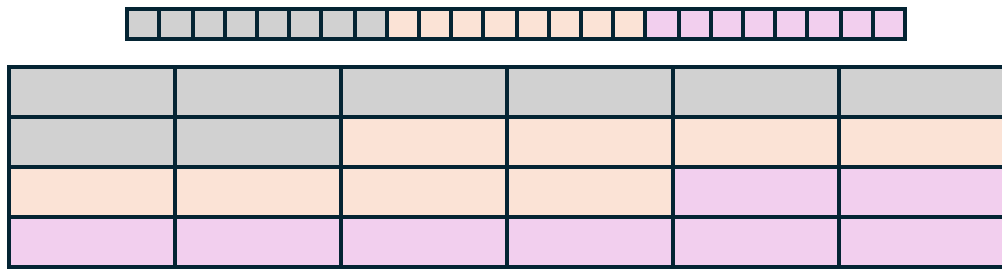
Randomly pick one PTE per region
RESET these PTEs' ACCESSED bits



Sleep 1ms while app executes
Wake up & check ACCESSED bits

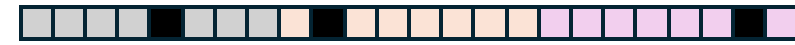


#2 Region-Based Sampling

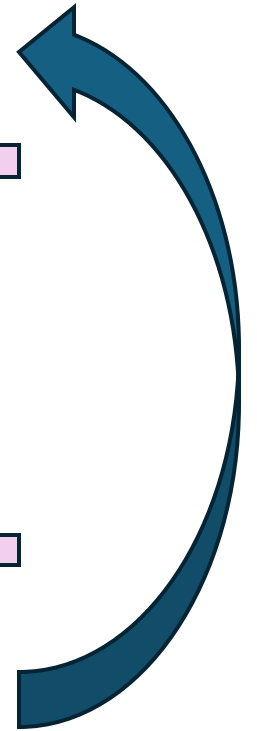
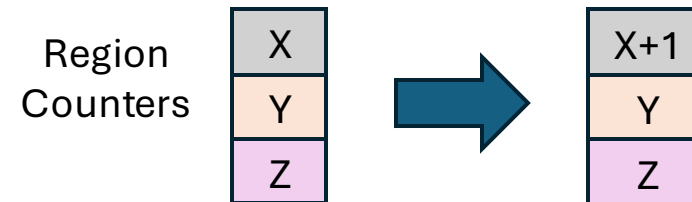


Workload's Mapped Address Space is divided into Regions.

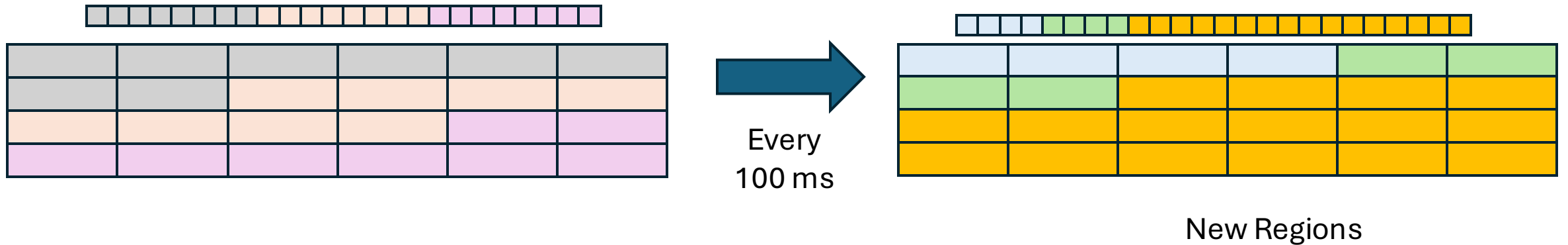
Randomly pick one PTE per region
RESET these PTEs' ACCESSED bits



Sleep 1ms while app executes
Wake up & check ACCESSED bits

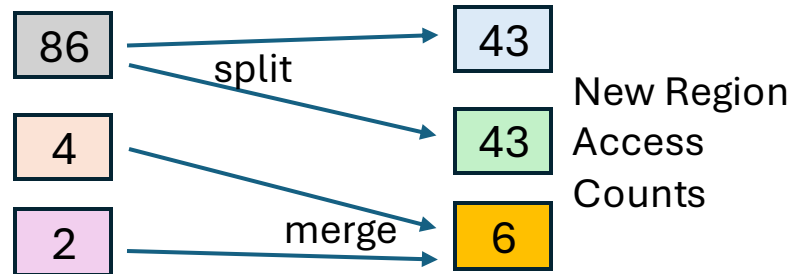


Region-Based Sampling



Adjacent Regions with identical access patterns (access counts) are MERGED.

Regions are split for precision.



Example: DAMON

Region-Based Sampling - Limitation

Does not scale well
with increasing
number of pages per
region !

PREMISE

"Sampled page represents
its whole region"

Premise holds only if
 $(\# \text{hot pages}) / (\# \text{pages in region})$
is significant

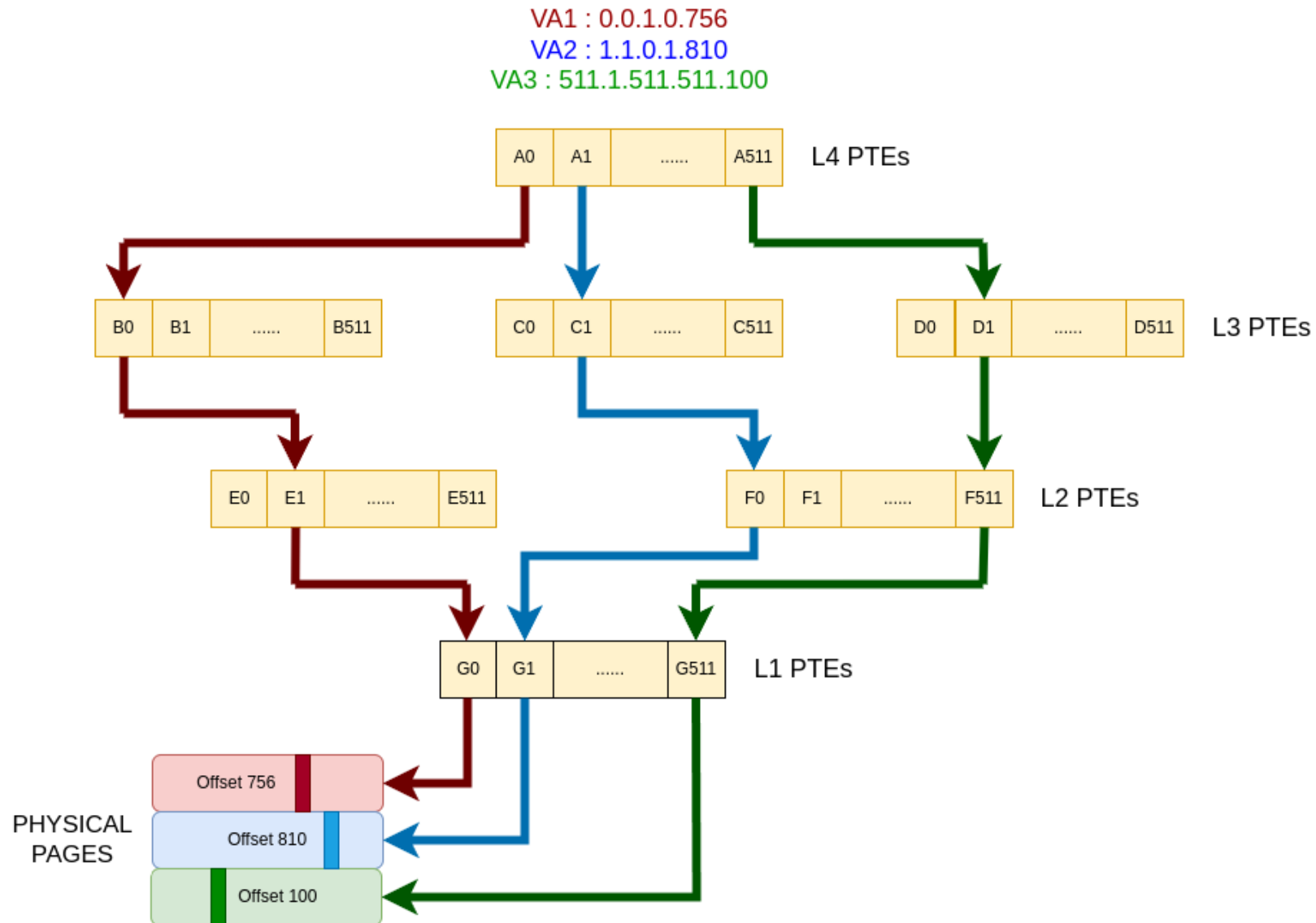
10 GB hot data in a 1 TB address space
DAMON fails to detect ANYTHING in 1 minute

#3 Performance Counters

- Sample hardware events - LLC Miss / TLB Miss / ...
 - Sample contains data address
 - Example: Intel PEBS (Processor Event-Based Sampling) on x86_64
-
- Limitation: HIGH OVERHEADS
 - Huge memory workload => more samples needed => high sampling frequency => Workload Slowdown => UNACCEPTABLE !

Telescope

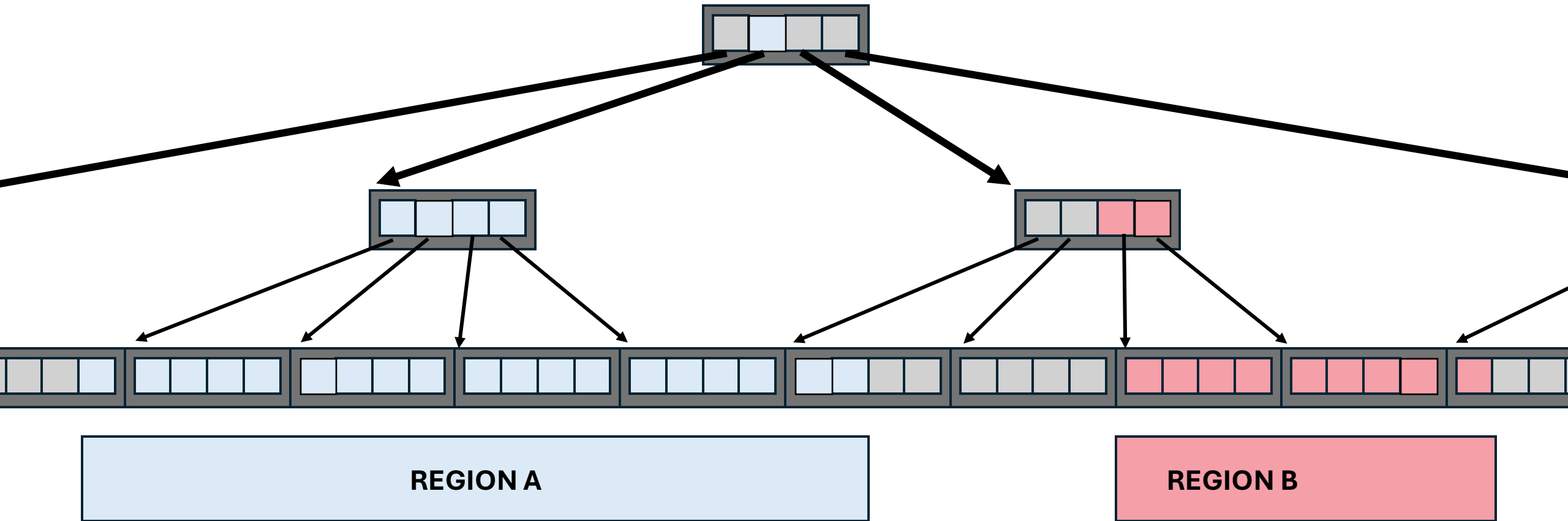
Radix Tree Structure of the Page Table



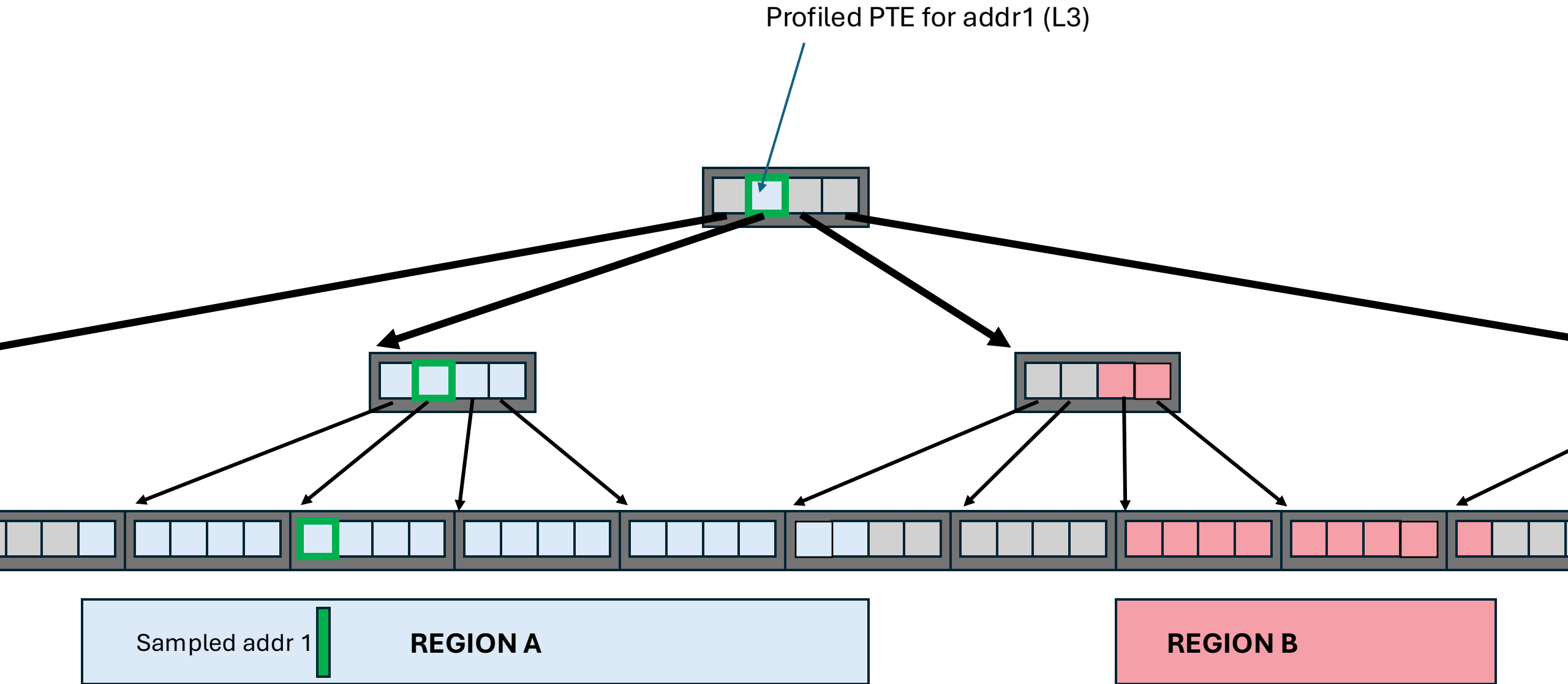
Key Insight

- On Page Walk, Page Table Walker sets the ACCESSED bit at every level of the Page Table tree.
- If at a higher level PTE, ACCESSED=0, then all its lower PTEs will have ACCESSED=0.
- Check the ACCESSED bit at higher level for fast but coarse-grained profiling of access patterns.

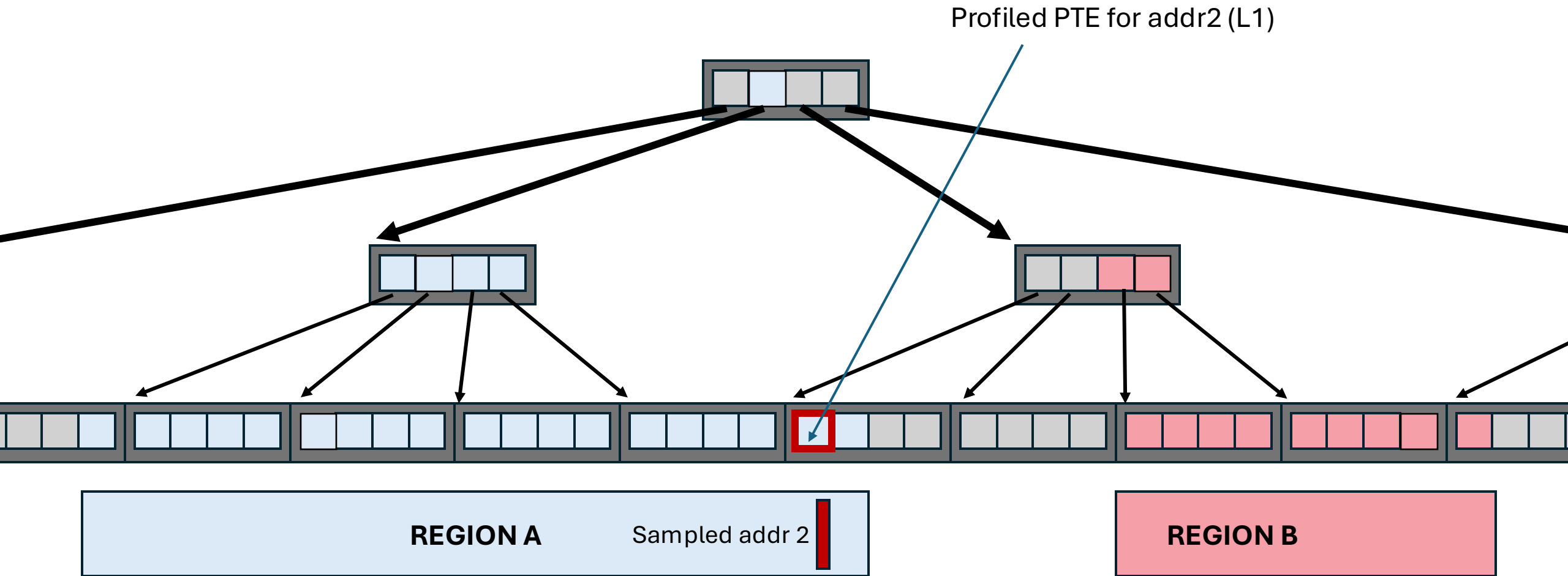
Telescope - Profiling



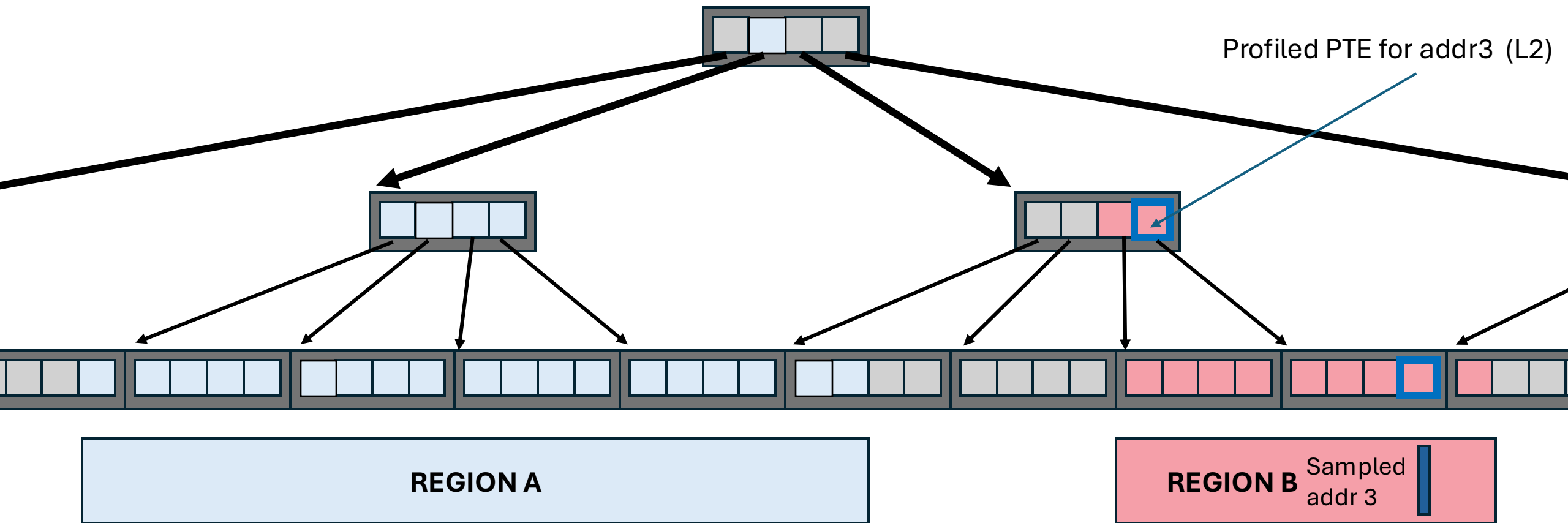
Telescope - Profiling



Telescope - Profiling



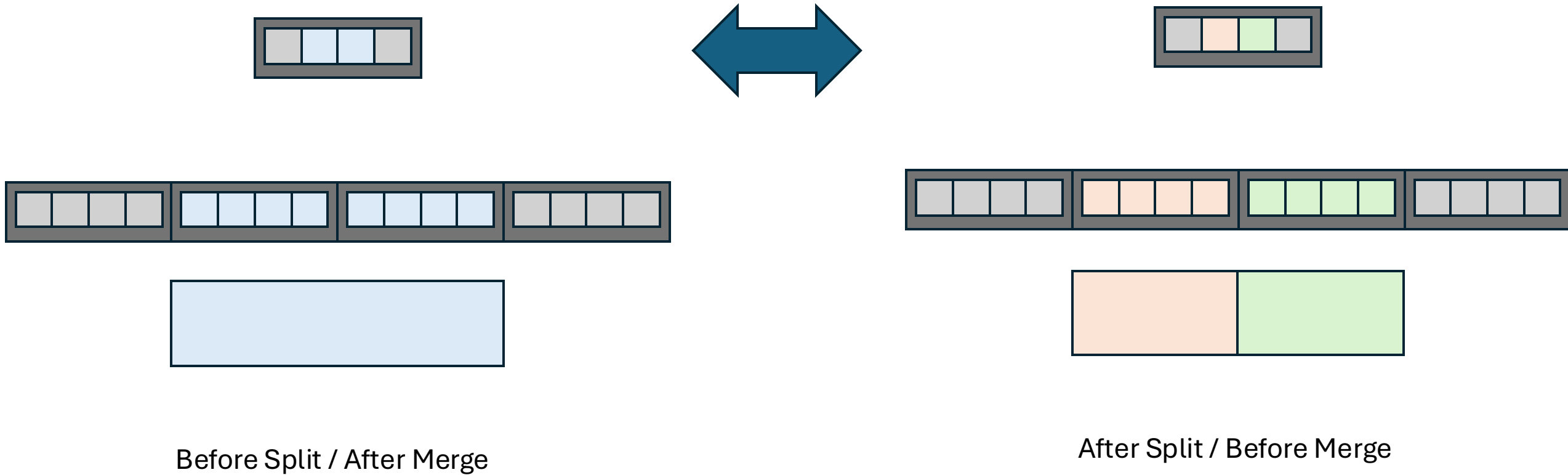
Telescope - Profiling



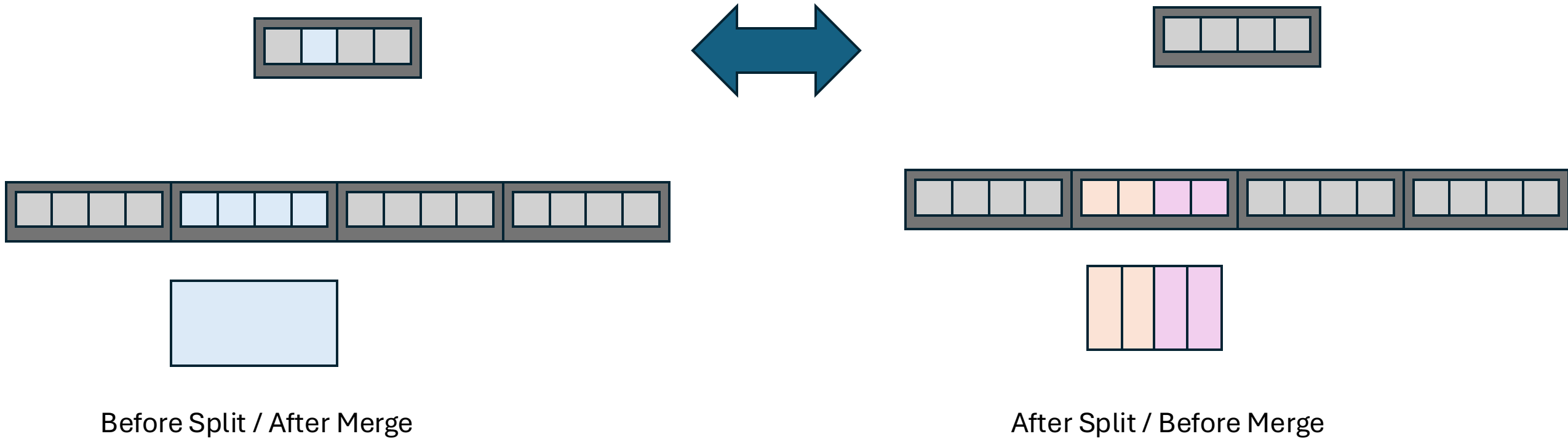
Telescope - FLEX variant

- Can pick a higher PTE that overshoots the region boundary.
- The overshooted region must be within an error threshold.
- Called FLEXIBLE Telescope, as opposed to BOUNDED Telescope.

Telescope - Region Split/Merge



Telescope - Region Split/Merge



Evaluation

Telemetry Techniques Evaluated

DAMON

Region-based Sampling

MODERATE Config

5 ms sampling window

AGGRESSIVE Config

1 ms sampling window

200 ms profiling window
1 sec vma-scan interval

PEBS

Performance Counters

MODERATE Config

5 kHz sampling freq

AGGRESSIVE Config

10 kHz sampling freq

Events sampled:

- MEM_INST_RETIRED.ALL_LOADS
- MEM_INST_RETIRED.ALL_STORES

TELESCOPE

Our Work

BOUNDED Config

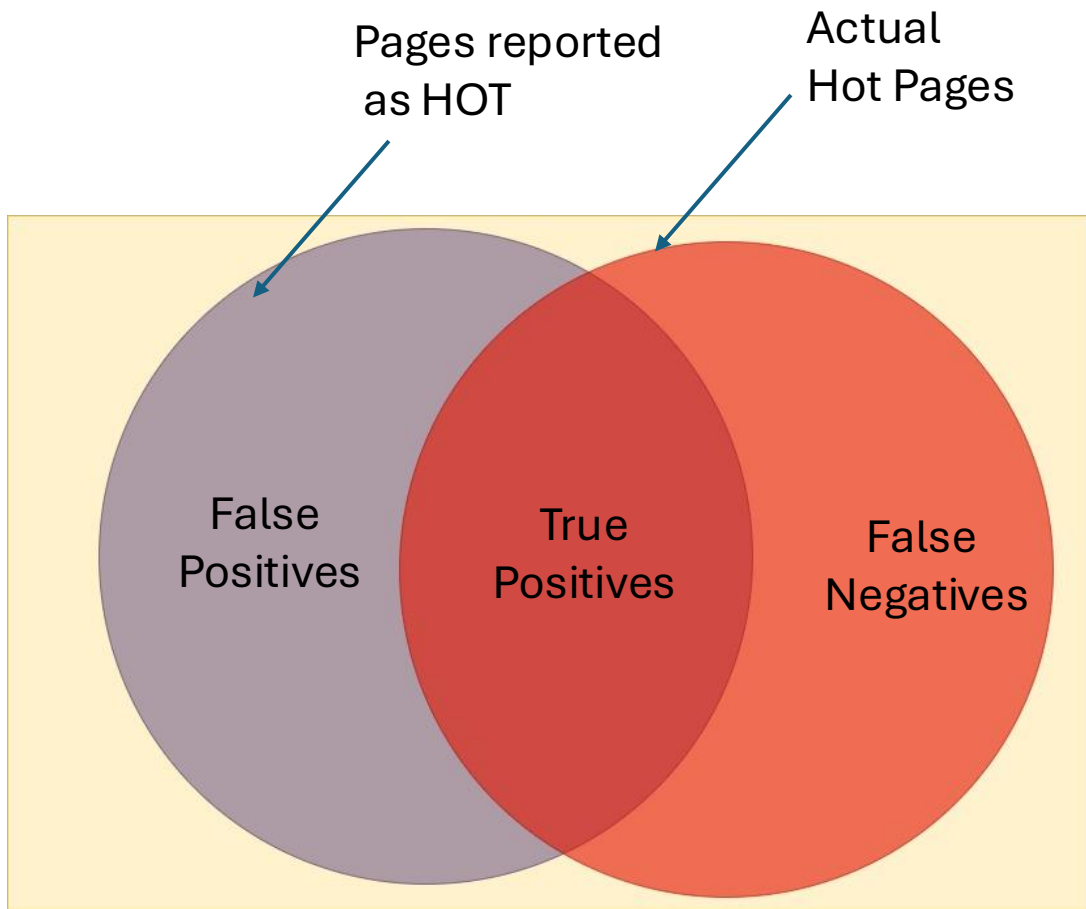
FLEX Config

L2: n = 25%

L3: n = 15%

5 ms sampling window
200 ms profiling window
1 sec vma-scan interval

Key Metrics



$$\text{PRECISION} = \frac{(\text{Actual HOT pages reported as HOT})}{(\text{Total pages reported as HOT})}$$

$$\text{OR } (\text{True Positives}) / (\text{True Positives} + \text{False Positives})$$

A measure of Accuracy

$$\text{RECALL} = \frac{(\text{Actual HOT pages reported as HOT})}{(\text{Actual HOT pages})}$$

$$\text{OR } (\text{True Positives}) / (\text{True Positives} + \text{False Negatives})$$

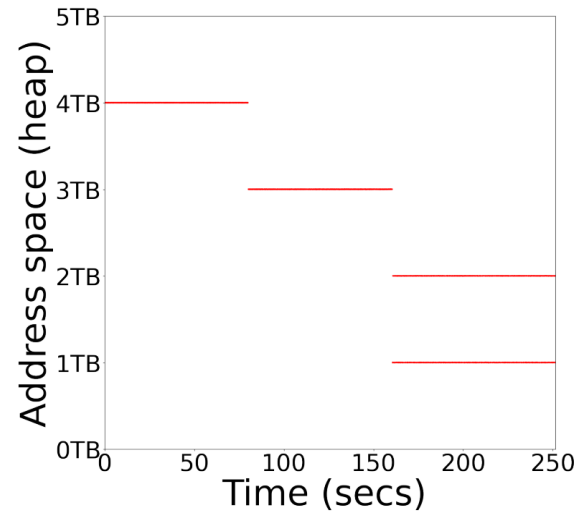
A measure of Coverage

Microbenchmark

- MASIM (Memory Access Simulator)



sjp38/masim

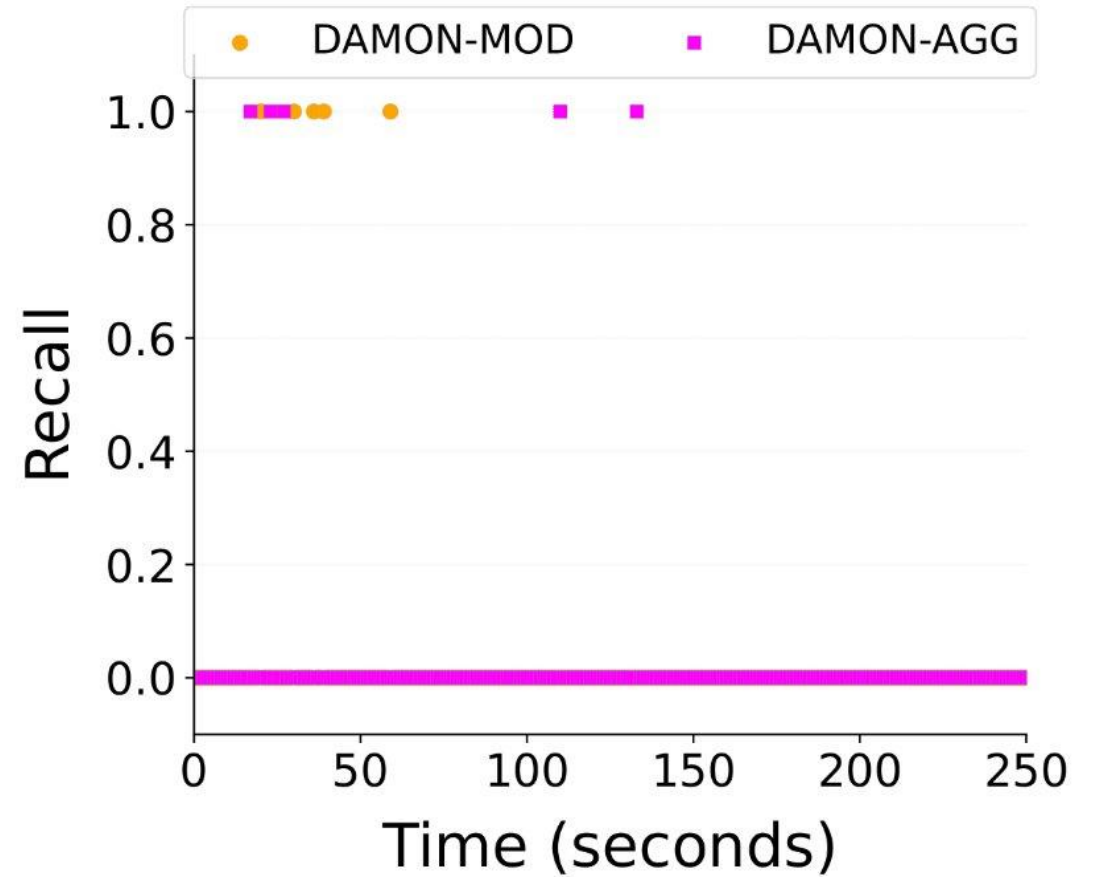
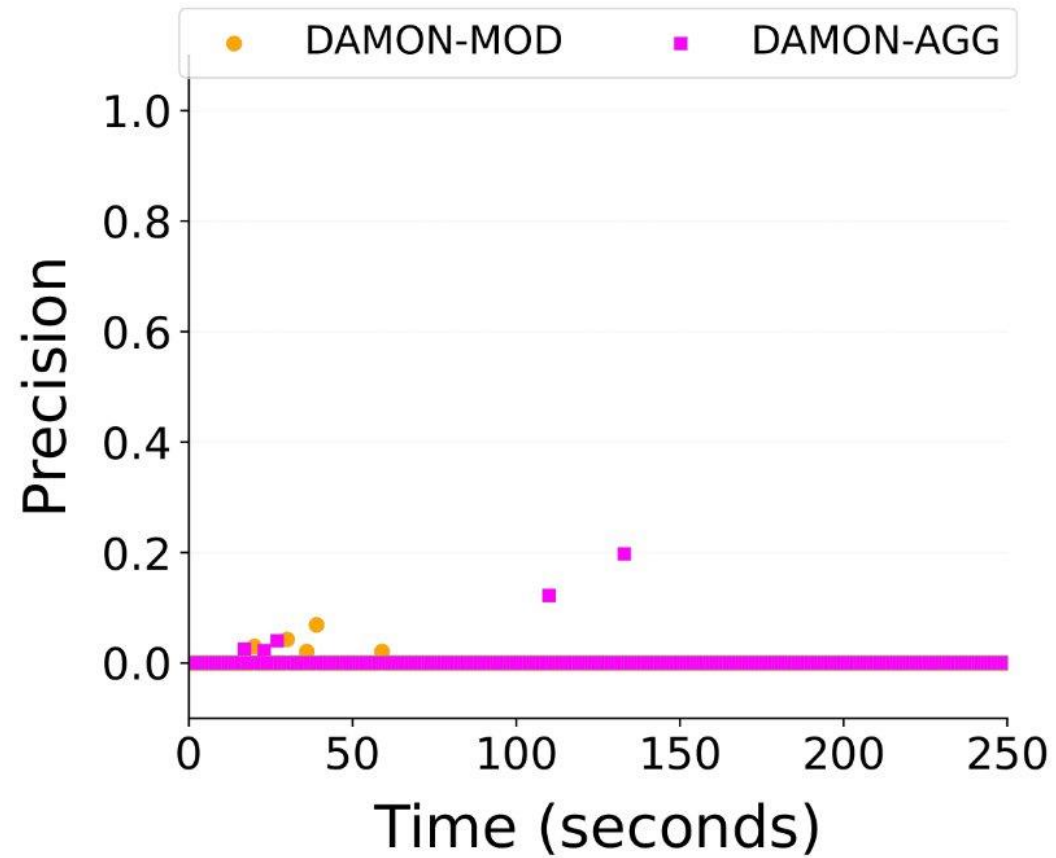


Multi-Phase

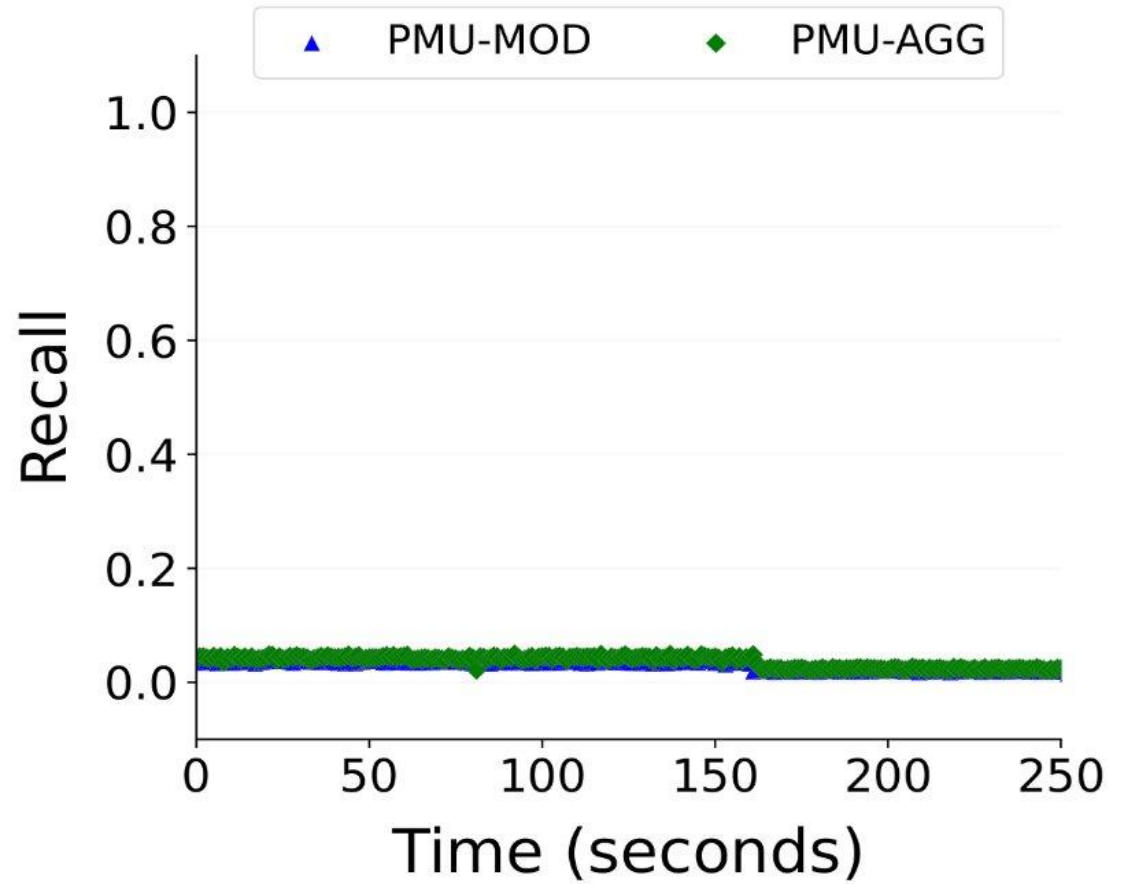
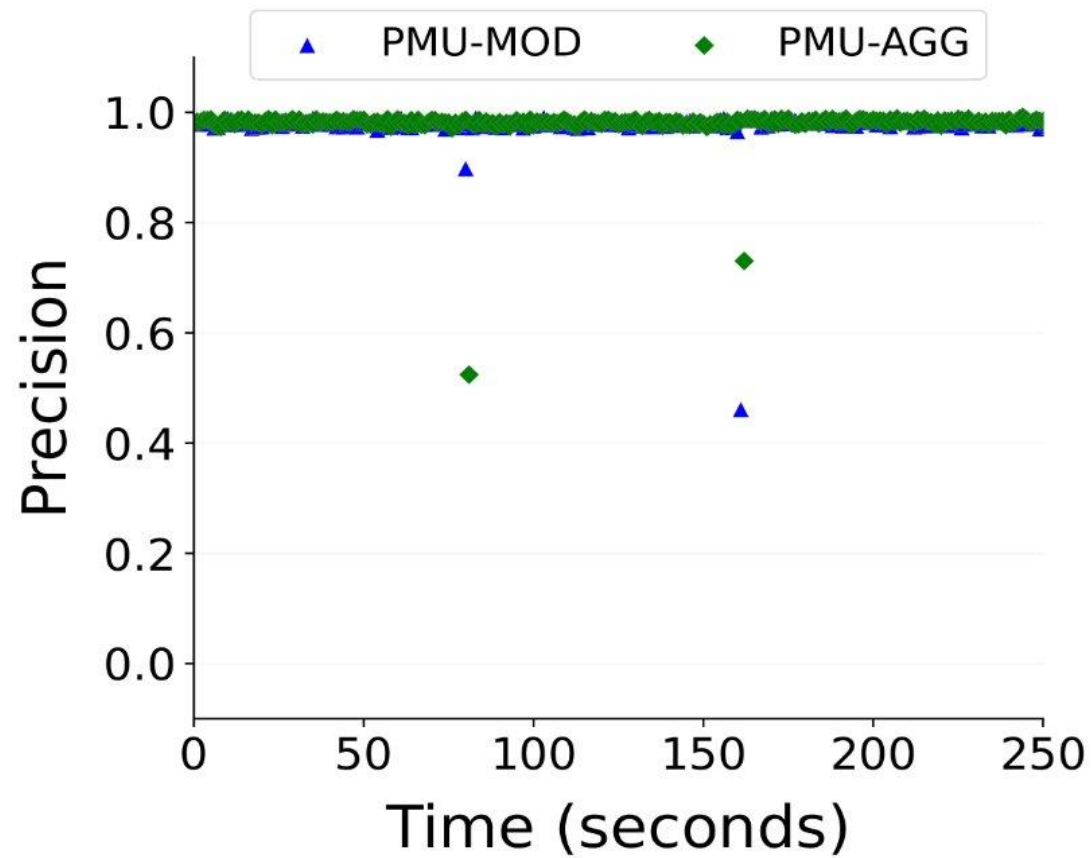
10 GB hot *needles* in a 5 TB heap

4 KB pages

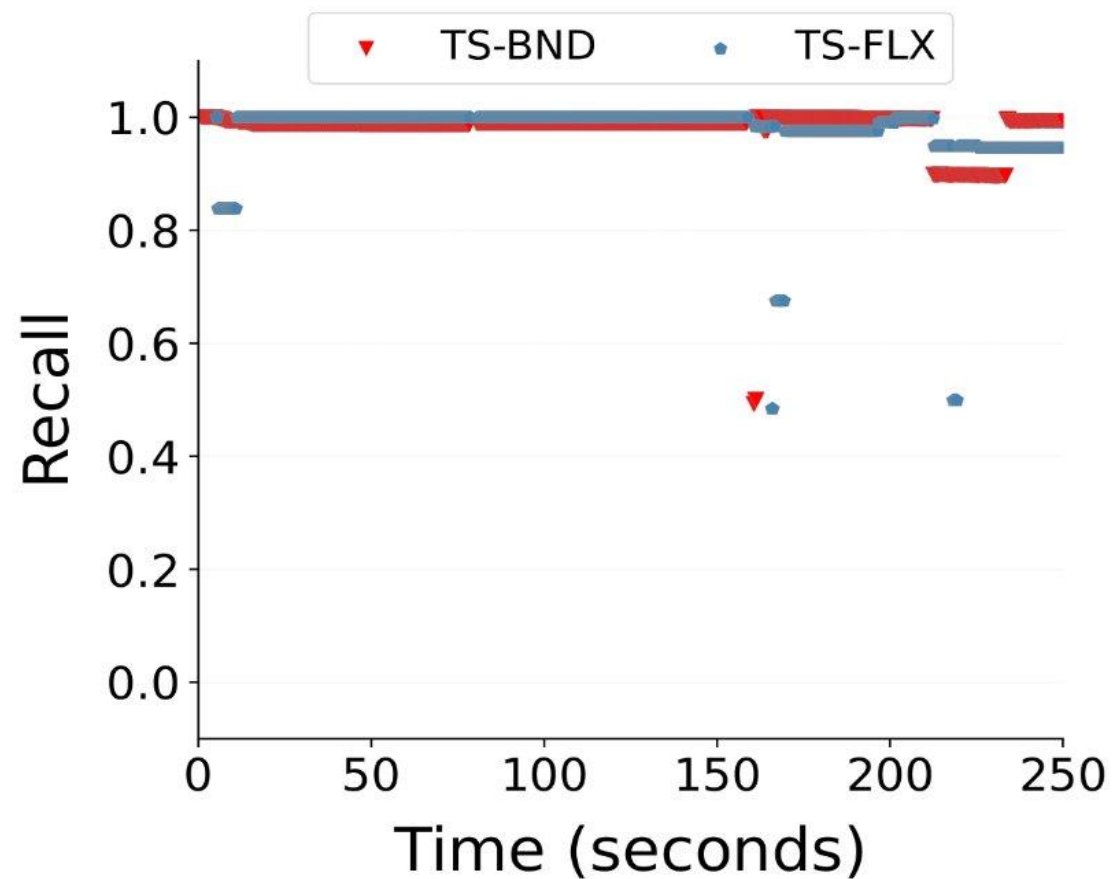
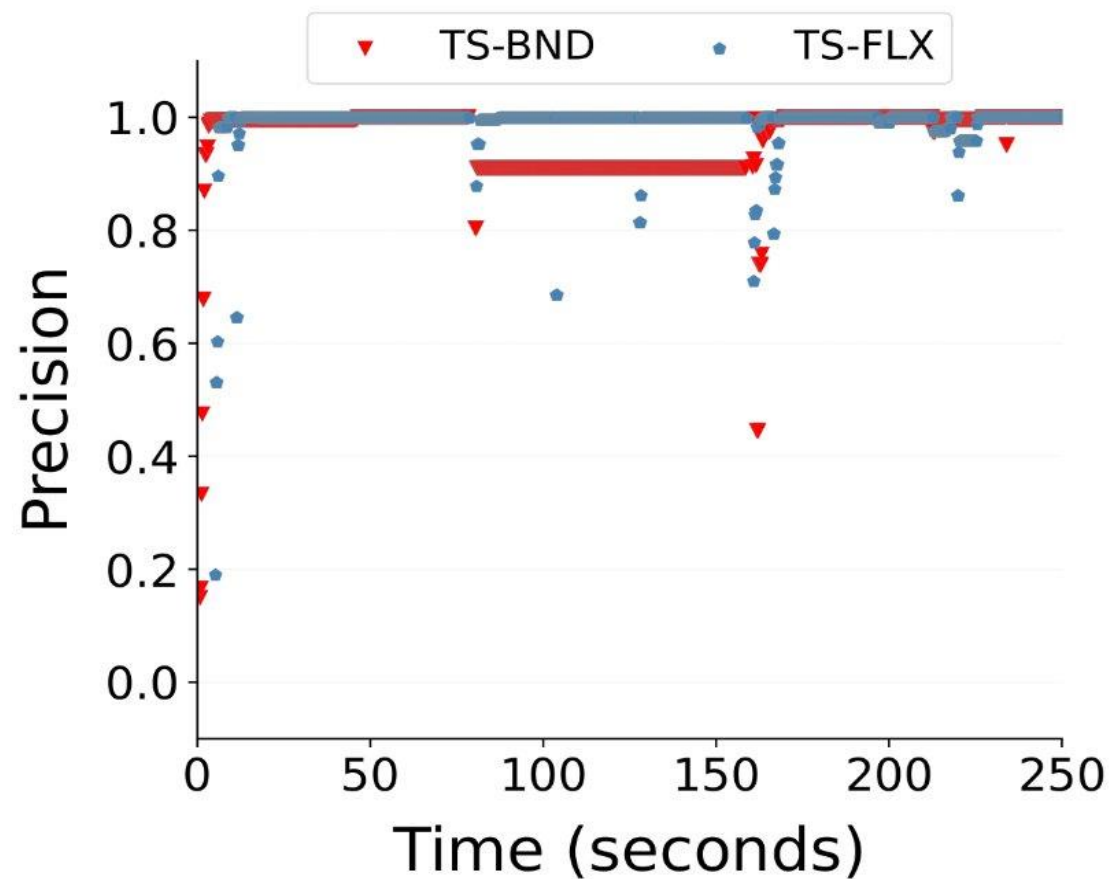
Results - DAMON



Results - PEBS



Results - Telescope

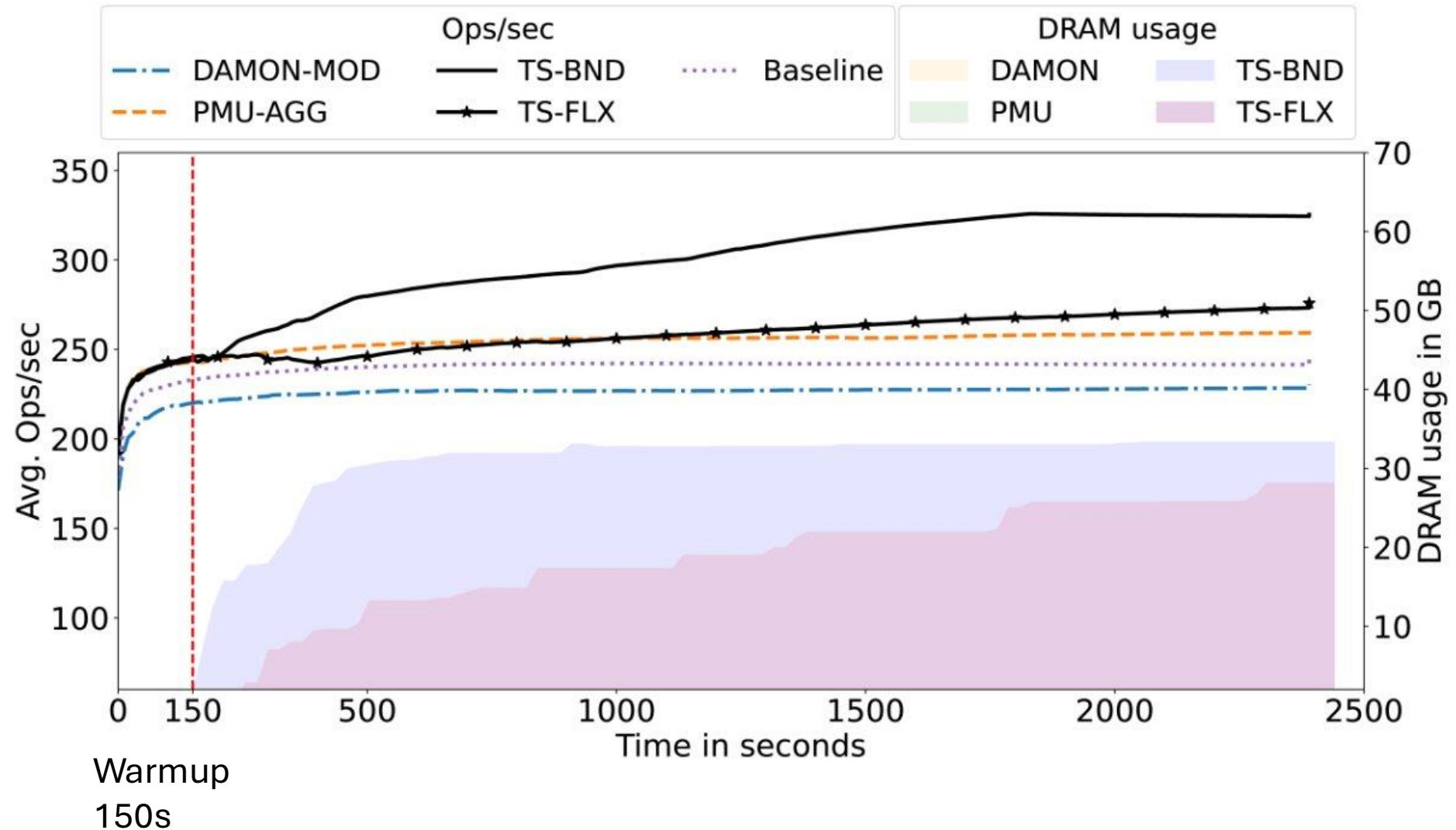


Real World Benchmark - Data Tiering

- Redis with YCSB for Load Generation
- 2TB data initialized in Optane NVM (cold tier)
- Telemetry outputs list of hot regions
- We migrate hottest regions to DRAM (hot tier)

- Metrics
 - Redis Throughput (ops/sec)
 - Tail Request Latency (95p)
 - DRAM Usage (GB)

Real World Benchmark - Results



Real World Benchmark - Application Impact

	Config.	95th percentile latency (ms)	
Redis	DAMON-MOD	850	59.13
	PMU-AGG	757	57.50
	Telescope-BND	696	54.01
	Telescope-FLX	741	55.55

Conclusion

- Effectiveness of terabyte-scale tiered memory systems depends on precise and timely identification of hot/cold data.
- Telescope introduces a novel page table profiling technique to quickly converge upon memory access patterns for workloads with huge memory footprints.
- We evaluate Telescope and compare it with other State-Of-The-Art telemetry techniques, and demonstrate its benefits for various applications with large memory footprints.
- Telescope future-proofs memory access telemetry for tiered memory systems with memory capacity up to and beyond the terabyte scale.



We are in the process of
upstreaming Telescope into the Linux Kernel

!!!

Scan QR code to see activity on
lore.kernel.org

THANK YOU!

Backup Slides

Region-Based Sampling

Successive Merge-and-Split ensures convergence to real access pattern over time.

Example: DAMON

