



Self-Driving DAMON/S: Controlled and Automated Access-aware Efficient Systems

SeongJae Park (SJ)

<sj@kernel.org> <sjpark@meta.com>



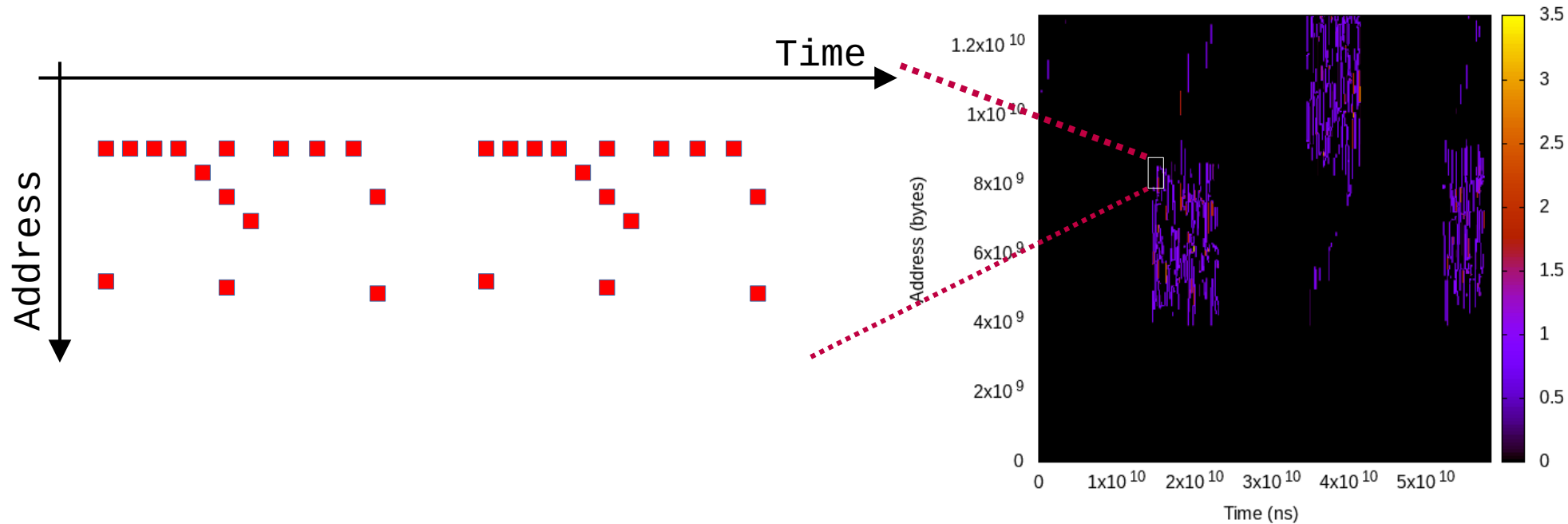
TL; DR

- Please try DAMON (again) with monitoring intervals auto-tuning



DAMON: Kernel Subsystem for Data Access Monitoring and Access-aware System Operations

Data Accesses: Events on Space/Time of Memory



Data Access Monitoring: Hope, Real and DAMON

- Hope: Precise (every bit), Complete (every moment), Light (prod online)
- Real: Expensive, YAGNI
- DAMON: Controlled and auto-tuned tradeoff for reasonable and practical monitoring

Space-Controlled Monitoring

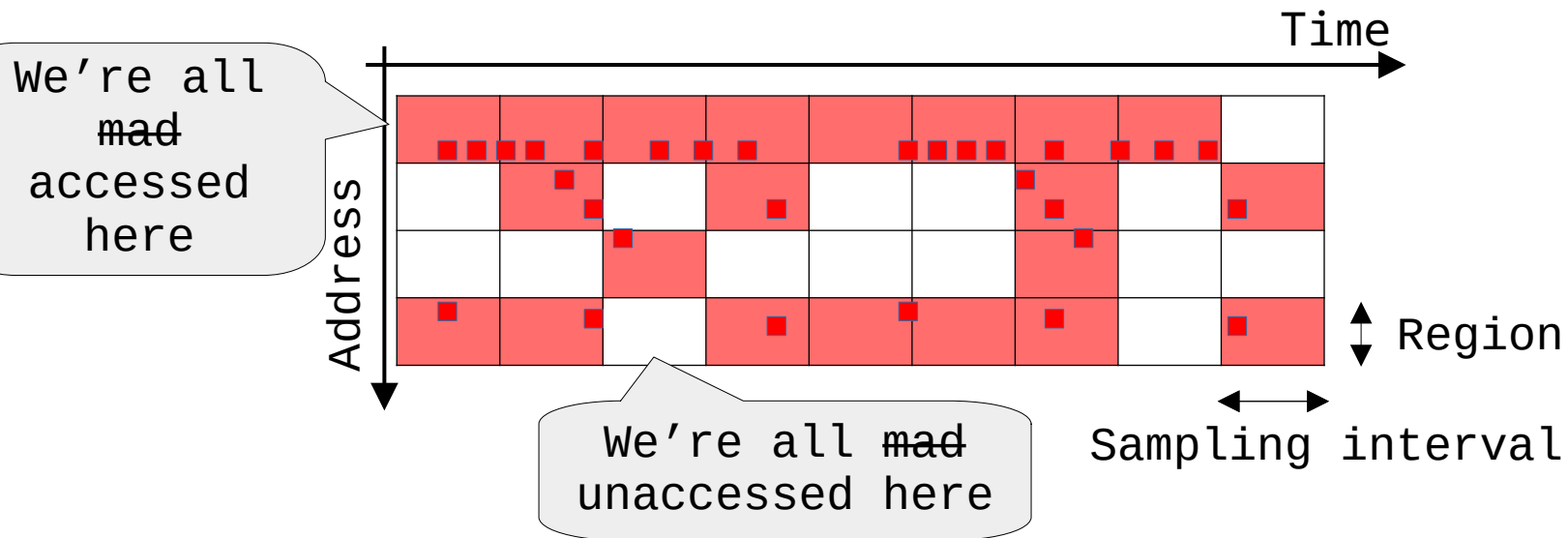
Region: Access Monitoring Unit for DAMON

- Defined as
 - A sub-area of the memory's space-time
 - A collection of adjacent elements that having similar access pattern
- Access check of one element per region is enough
- e.g., “This page is accessed within last 1 second; a cacheline is checked”

```
$ cat wonder_region_1  
We're all mad [un]accessed here
```

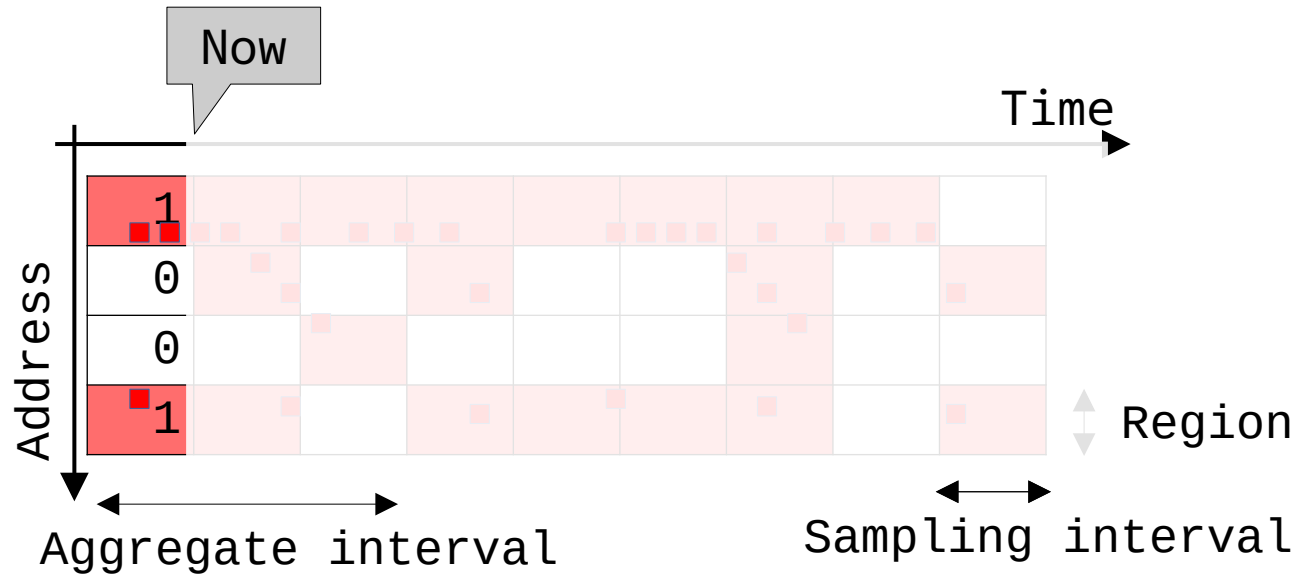
Fixed Space/Time Granularity, ≤ 1 Access Frequency

- Sort of periodic fixed granularity idleness monitoring
- Time overhead: “memory size / *space granularity*”
- Space overhead: “time overhead * monitoring time / *time granularity*”
- Reduced and controllable, but still ruled by memory size and total monitoring time



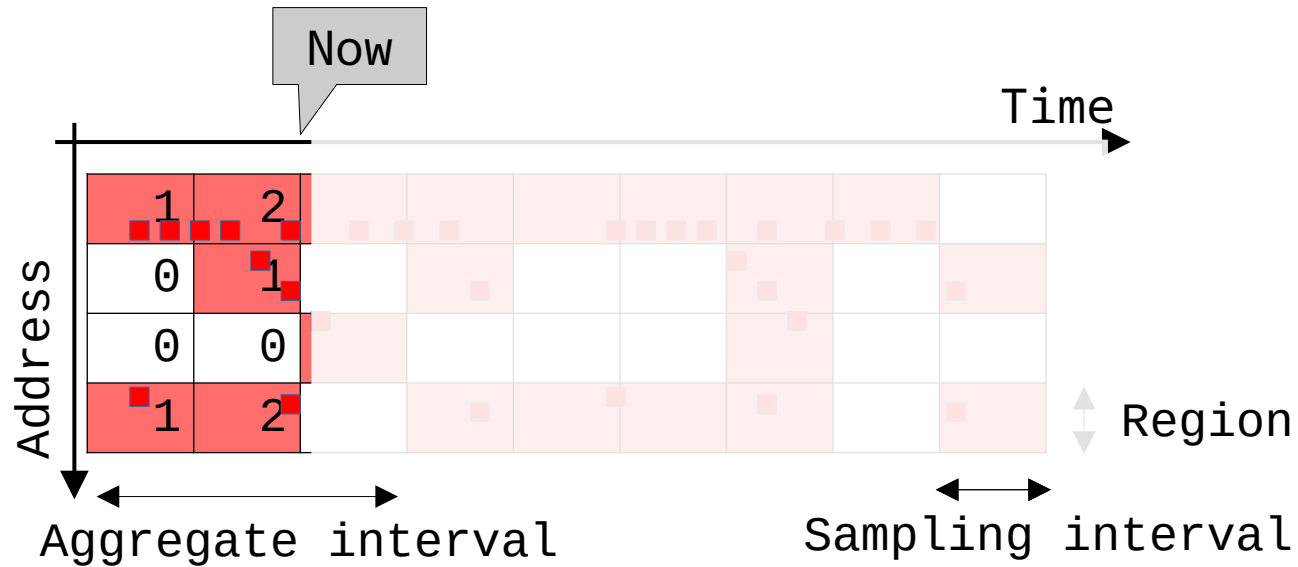
Fixed Space/Time Granularity, $\leq N$ Access Frequency

- Accumulate (sampled) access check results via per-region counter



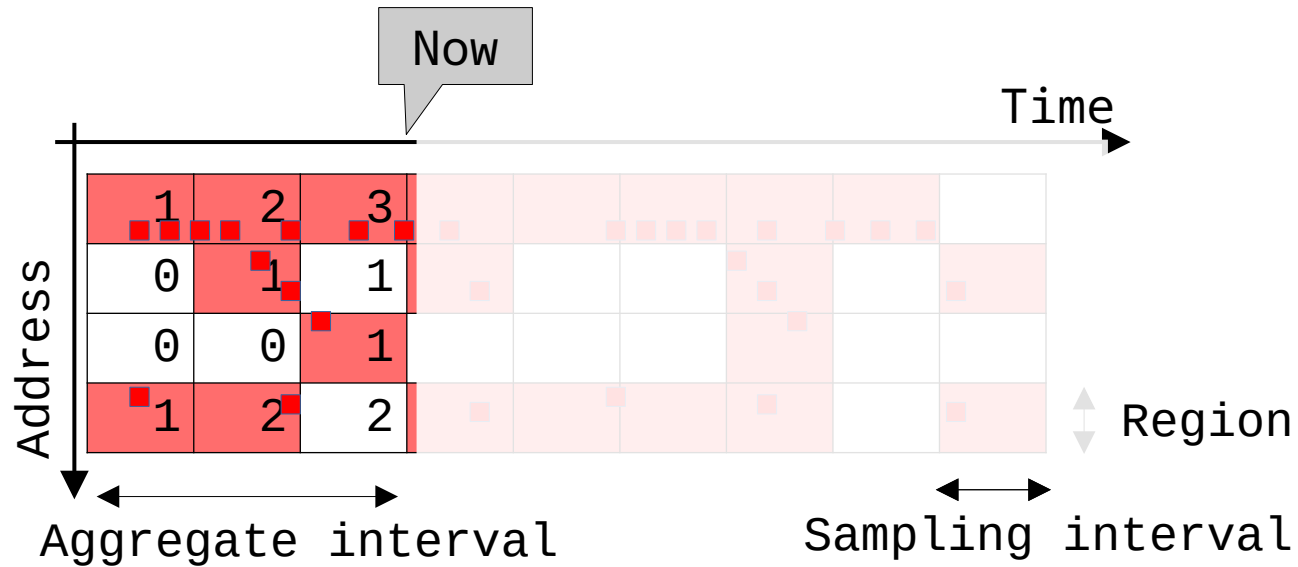
Fixed Space/Time Granularity, $\leq N$ Access Frequency

- Accumulate (sampled) access check results via per-region counter



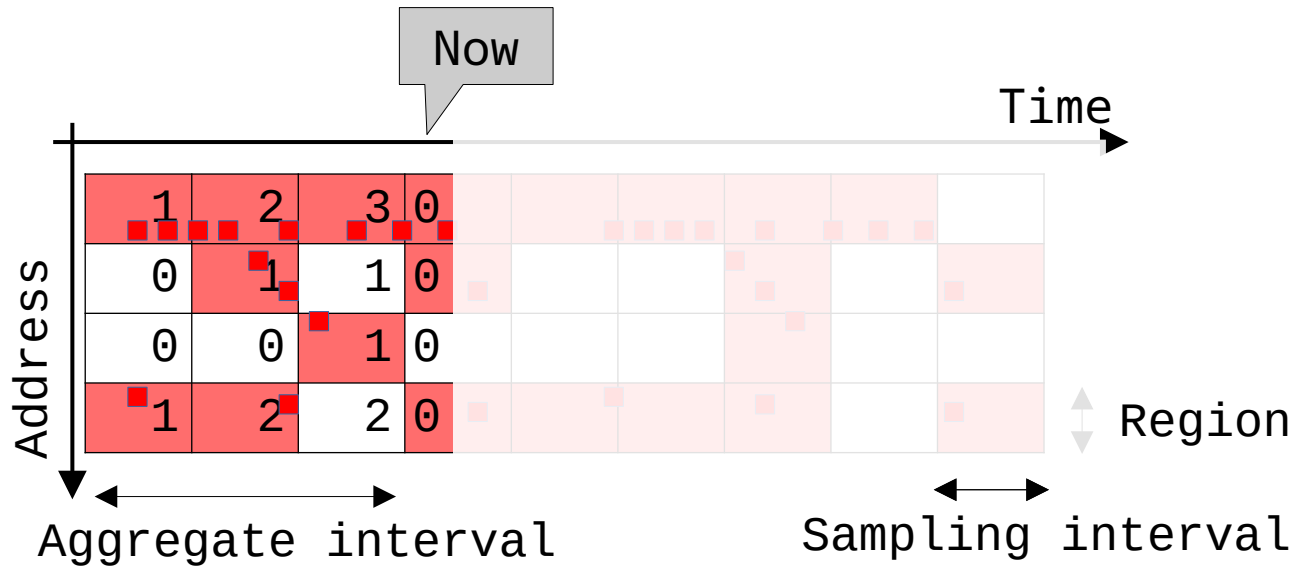
Fixed Space/Time Granularity, $\leq N$ Access Frequency

- Accumulate (sampled) access check results via per-region counter



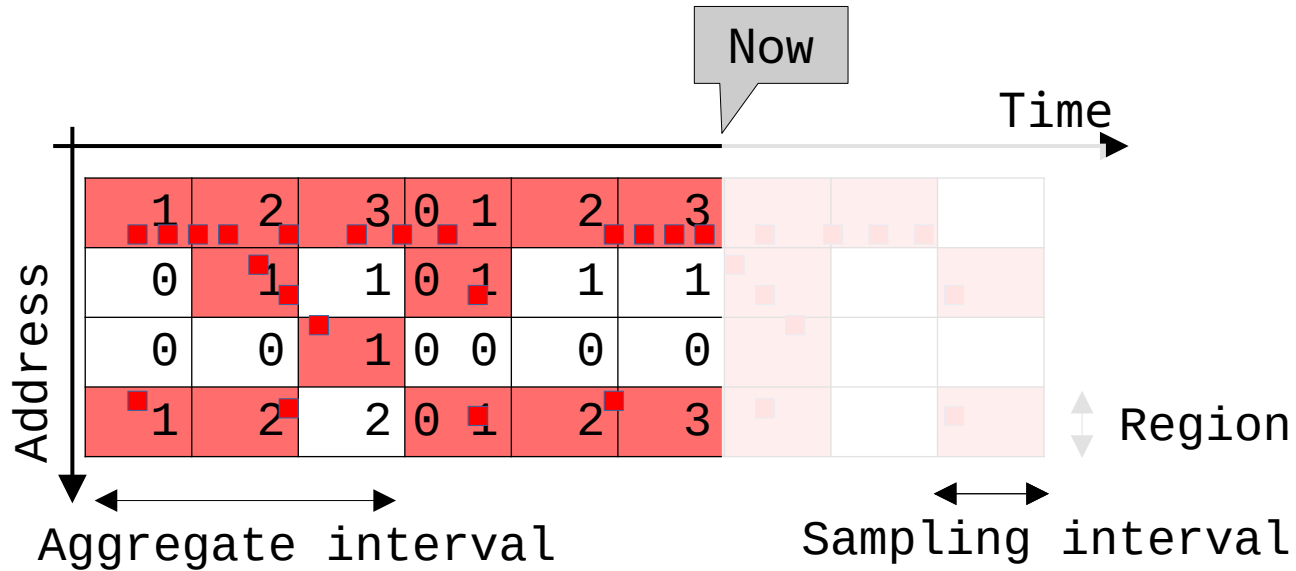
Fixed Space/Time Granularity, $\leq N$ Access Frequency

- Accumulate (sampled) access check results via per-region counter



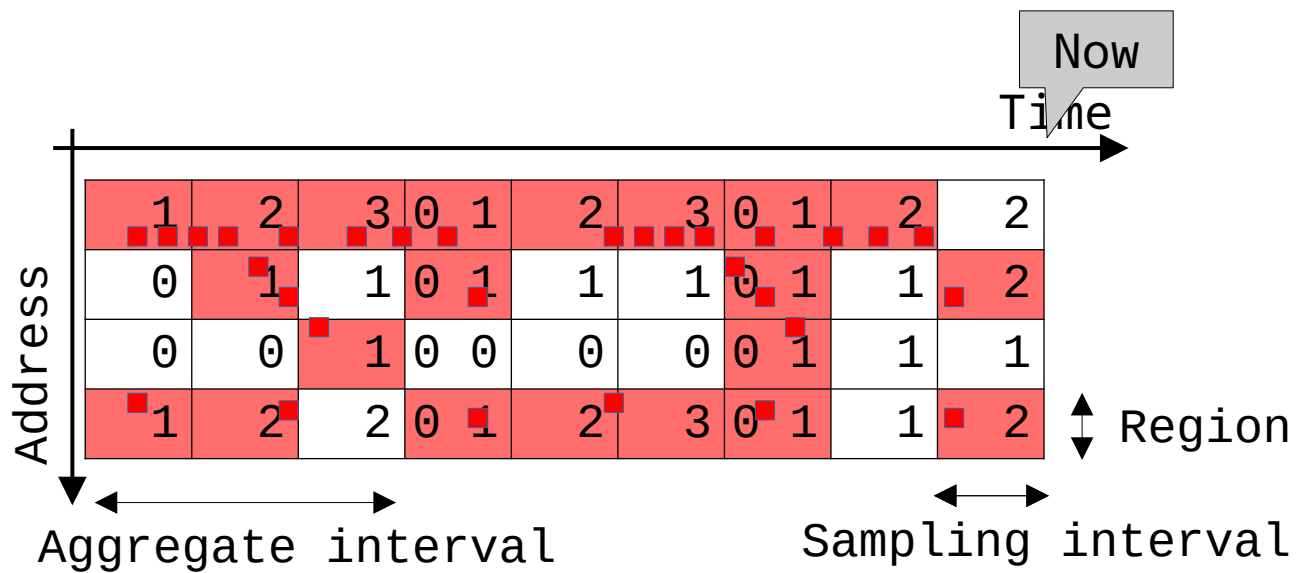
Fixed Space/Time Granularity, $\leq N$ Access Frequency

- Accumulate (sampled) access check results via per-region counter



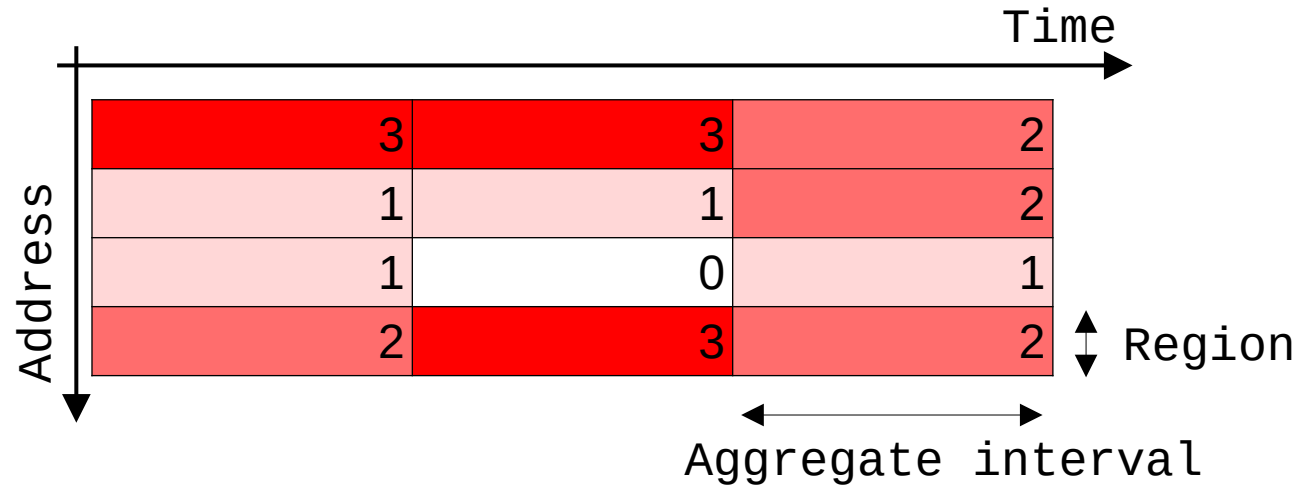
Fixed Space/Time Granularity, $\leq N$ Access Frequency

- Accumulate (sampled) access check results via per-region counter



Fixed Space/Time Granularity, $\leq N$ Access Frequency

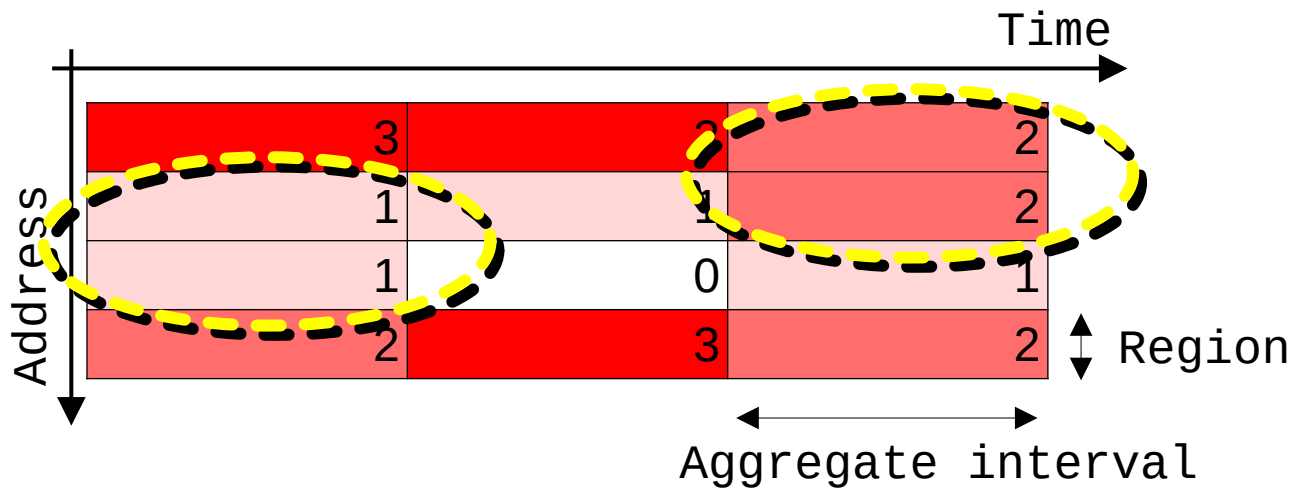
- Accumulate access checks via per-region counter
- Reduce space overhead to “ $1/N$ ”
- Still, $O(\text{memory size} * \text{total monitoring time})$



Space-Auto-Tuned Monitoring

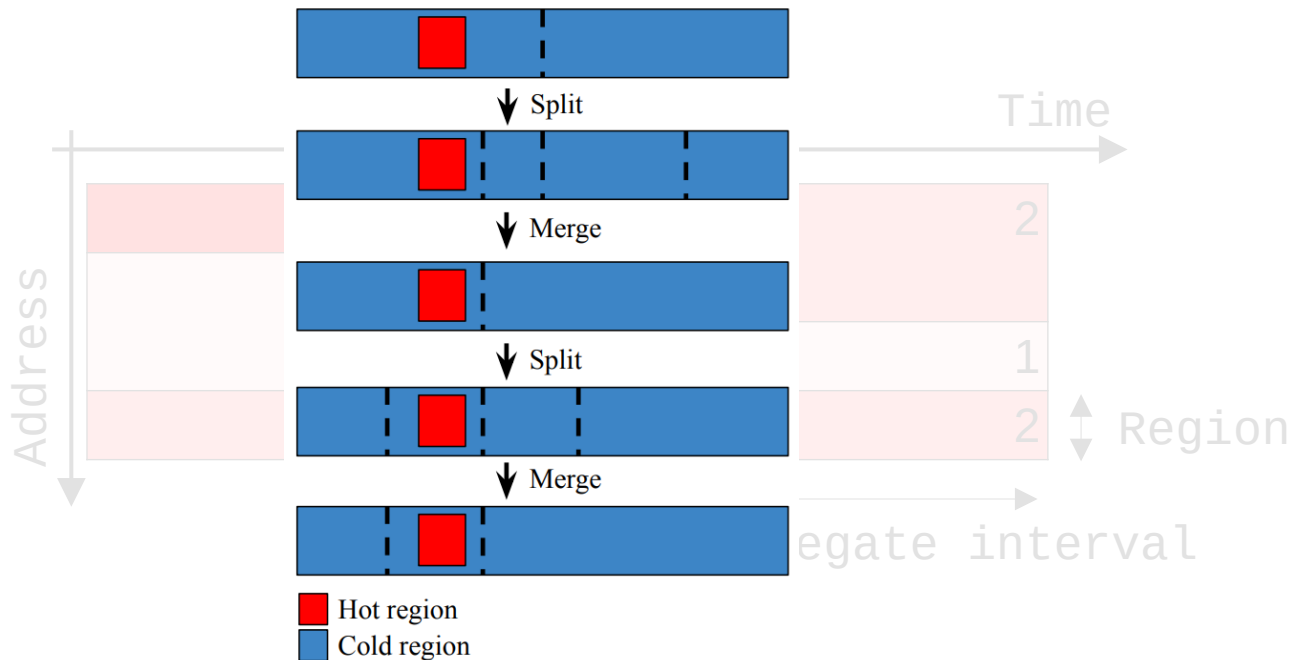
Problems of Fixed Space Granularity

- Adjacent regions of similar hotness are wastes
- Impossible to do fine-grained space monitoring



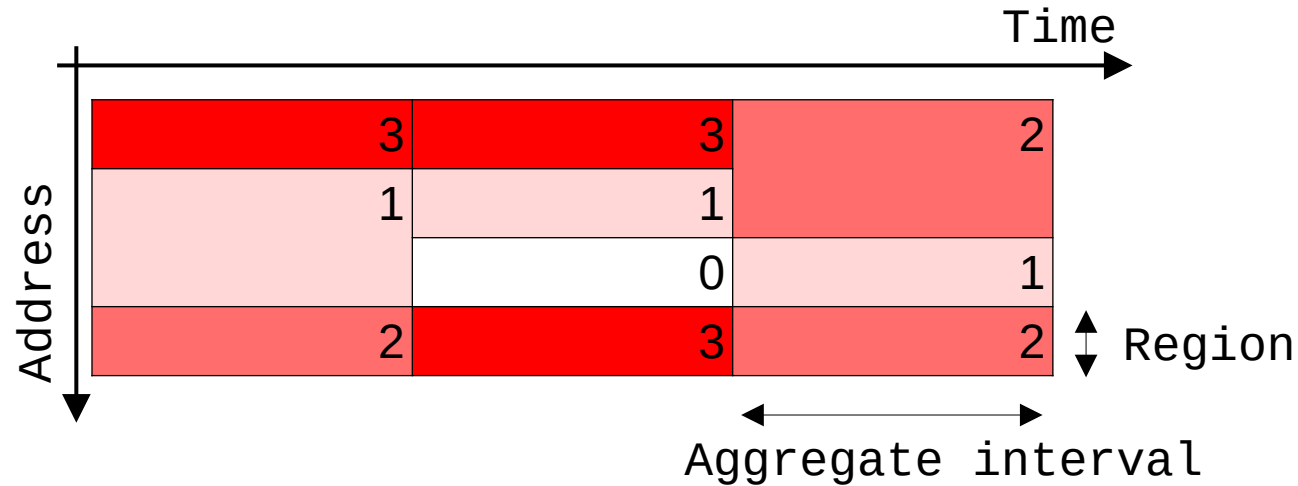
Auto-tuned Dynamic Space Granularity: Mechanisms (1/2)

- Repeat merging the wasteful regions and randomly splitting regions
 - The number of region == number of different access patterns
- Let user set min/max number of total regions



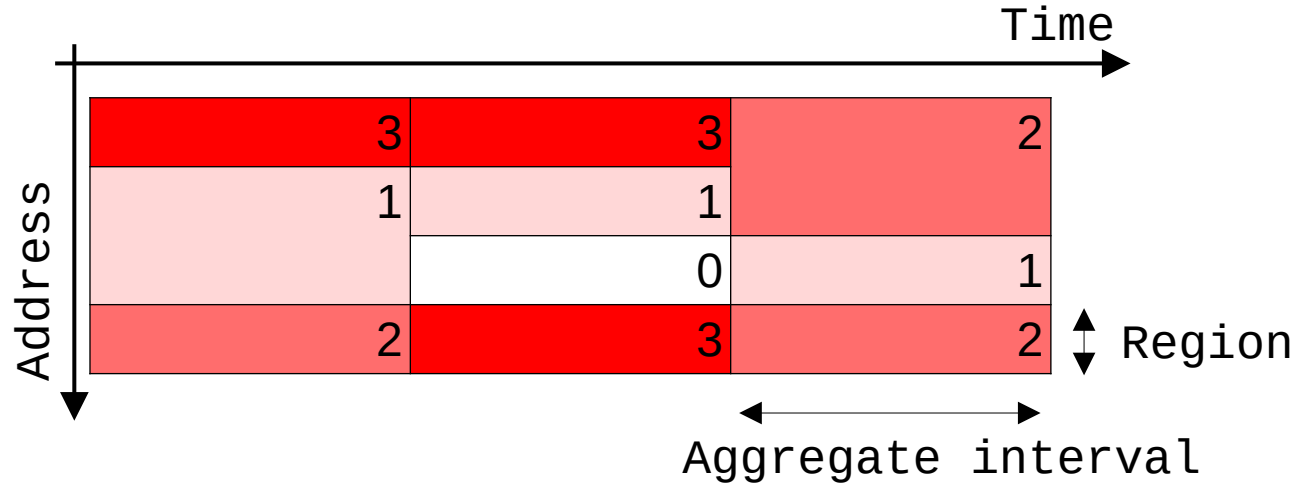
Auto-tuned Dynamic Space Granularity: Mechanisms (2/2)

- Repeat merging the wasteful regions and randomly splitting regions
 - The number of region == number of different access patterns
- Let user set min/max number of regions



Auto-tuned Dynamic Space Granularity: Overhead/Accuracy

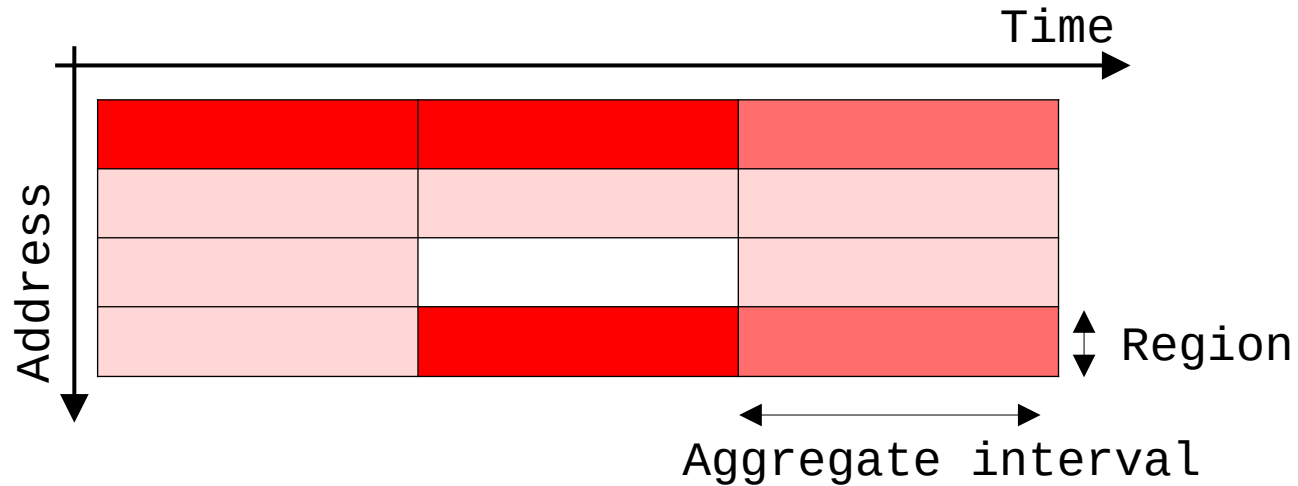
- Time overhead: $\min(\text{different access patterns}, \text{max number of regions})$
 - No more ruled by memory size, fully controlled and auto-tuned
- Accuracy: best-effort high
 - Auto-tuned dynamic granularity can find even bit level accesses



Time-Controlled Monitoring

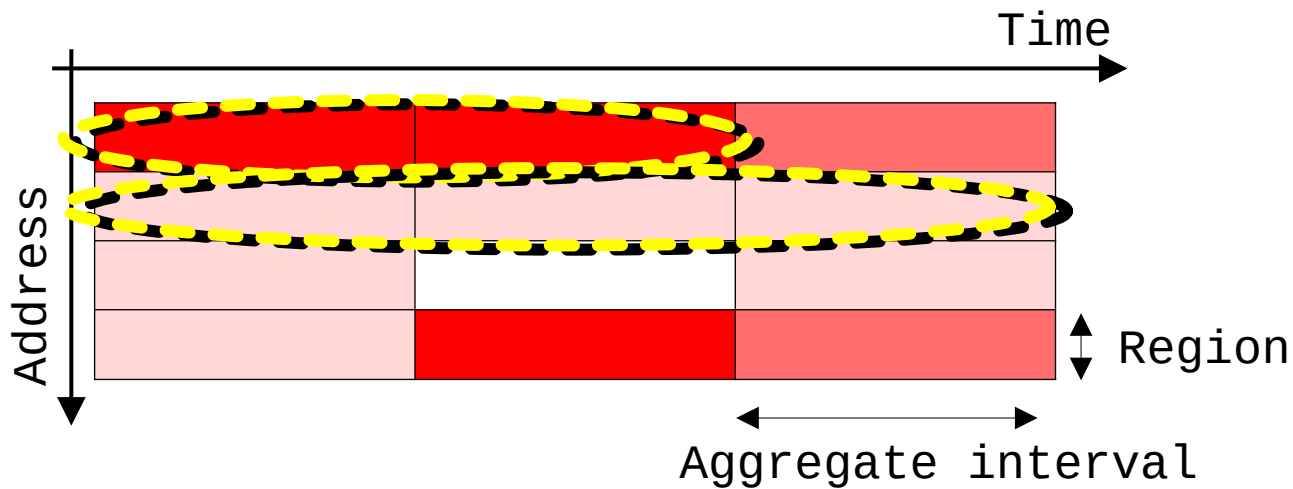
Problems of Fixed Time Granularity Regions (1/2)

- The definition of regions is not only about space, but also about time



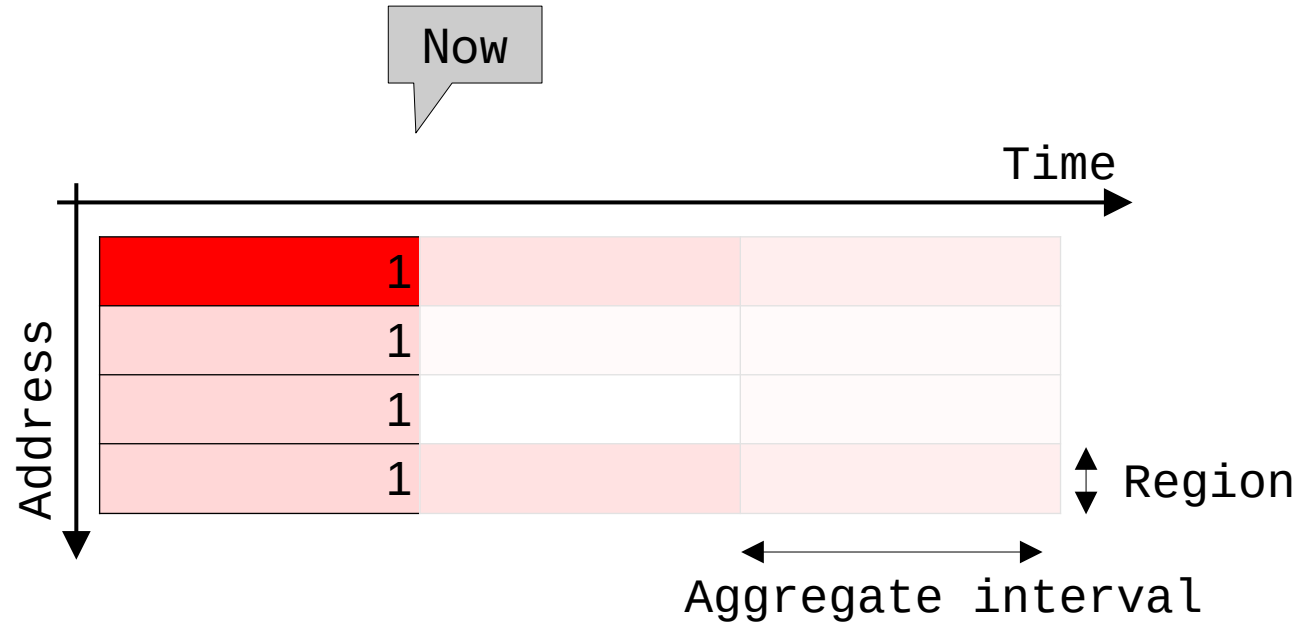
Inefficiency of Fixed Time Granularity Regions (2/2)

- The definition of regions is not only about space, but also about time
- Multiple time-adjacent regions of similar hotness: only waste



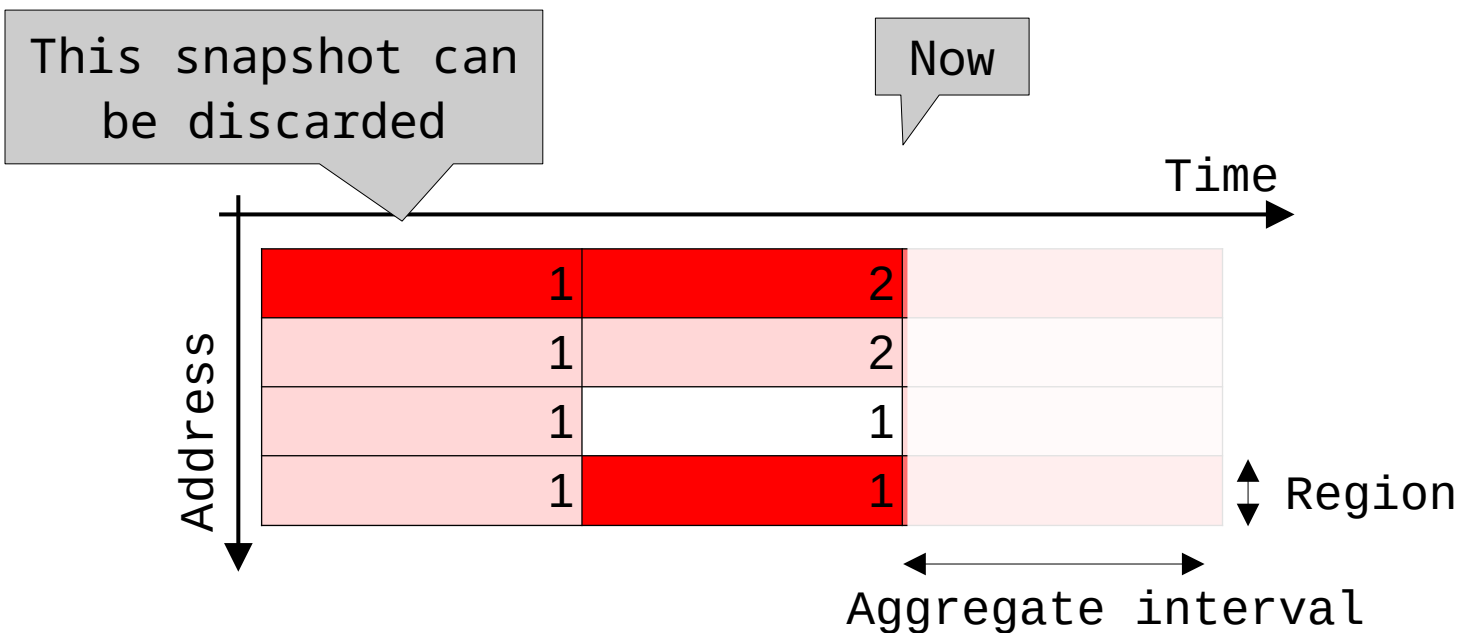
Dynamic Time Granularity (1/3)

- Count how long the hotness has kept
- Snapshot contains history of useful length



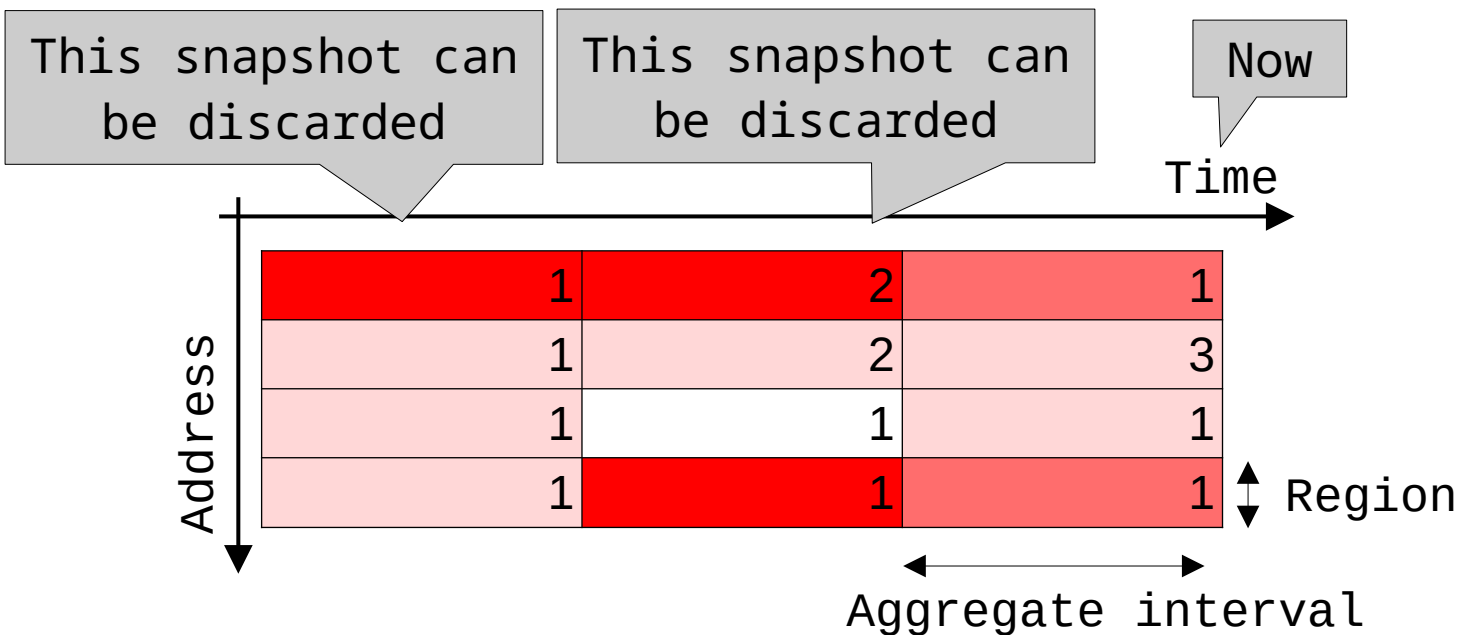
Dynamic Time Granularity (2/3)

- Count how long the hotness has kept
- Snapshot contains history of useful length



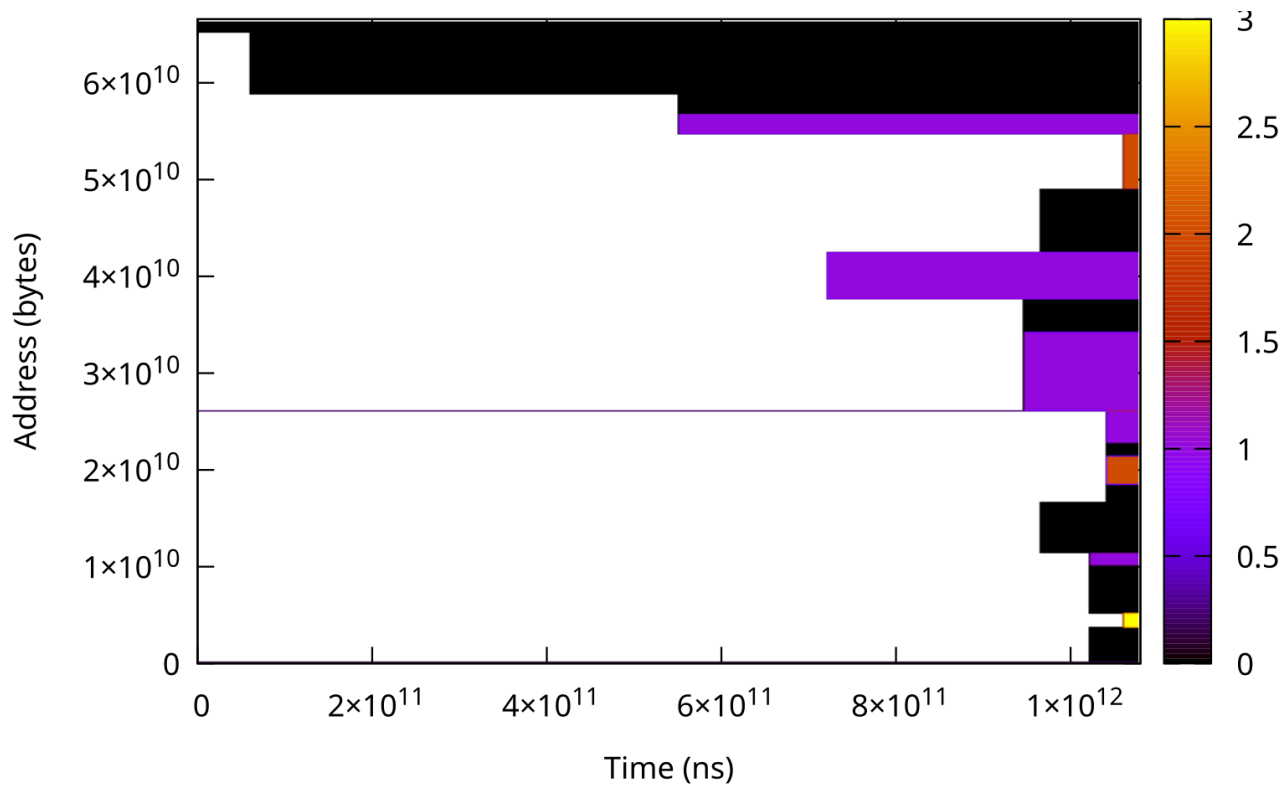
Dynamic Time Granularity (3/3)

- Count how long the hotness has kept
- Snapshot contains history of useful length



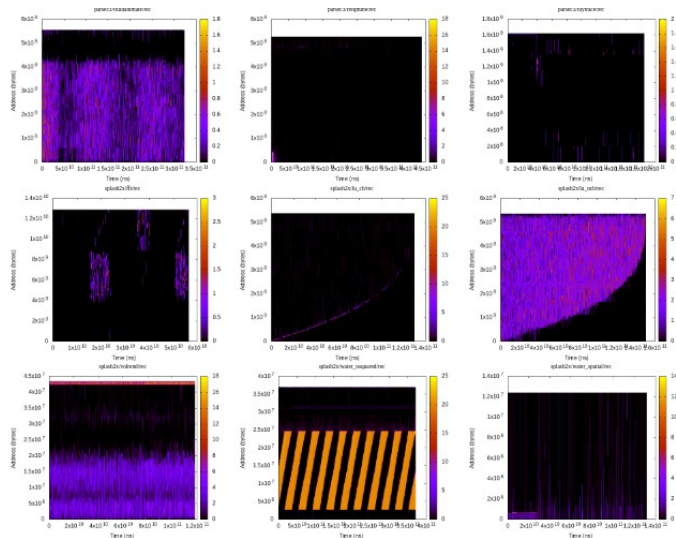
Snapshot: The Output of DAMON

- $O(\text{max_nr_regions})$ time/space overhead
- Both time/space overheads are not ruled by memory size/monitoring time



Potential, or Aimed Usages

- Profiling (e.g., GIF demo [link](#))
 - Help better understand and find rooms for improvements
- Profiling-guided Optimizations
 - Could be done on both offline and online
- Why not let kernel just (transparently) works?

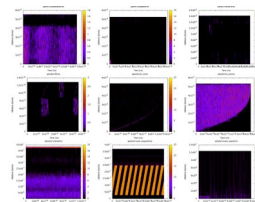


# Memory Footprints Distribution					
percentile	0	25	50	75	100
wss	0 B	9.520 MiB	9.543 MiB	9.785 MiB	107.039 MiB
rss	104.820 MiB	104.820 MiB	104.820 MiB	104.820 MiB	104.820 MiB
vsz	108.352 MiB	108.352 MiB	108.352 MiB	108.352 MiB	108.352 MiB
sys_used	2.348 GiB	2.417 GiB	2.424 GiB	2.436 GiB	2.453 GiB

# Hotspot functions			
# Samples: 589K of event 'cpu-clock:ppp'			
# Event count (approx.): 147266750000			
#			
# Overhead	Command	Shared Object	Symbol
#			
#			
57.73%	swapper	[kernel.kallsyms]	[k] pv_native_safe_halt
40.26%	masim	masim	[.] do_seq_wq
0.11%	python3	python3.11	[.] _PyEval_EvalFrameDefault
0.09%	ps	[kernel.kallsyms]	[k] do_syscall_64
0.05%	ps	[kernel.kallsyms]	[k] memset_orig
0.04%	ps	libc.so.6	[.] open64

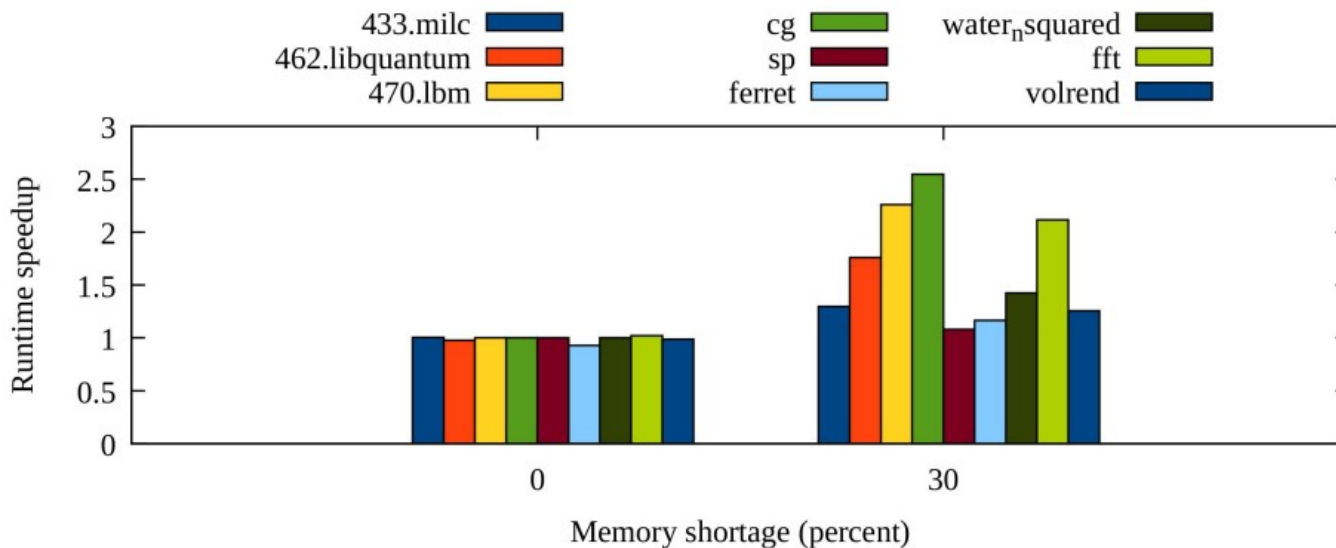
Potential, or Aimed Usages

- Profiling (e.g., GIF demo [link](#))
 - Help better understand and find rooms for improvements
- Profiling-guided Optimizations
 - Could be done on both offline and online
- Why not let kernel just (transparently) works?



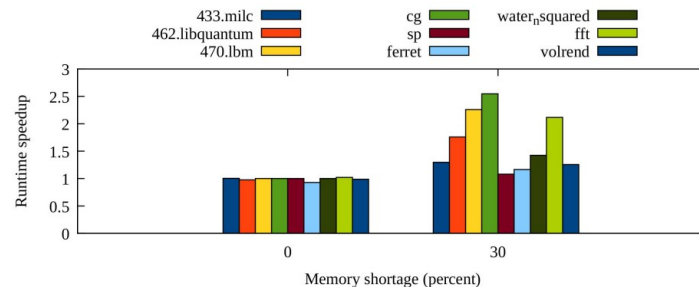
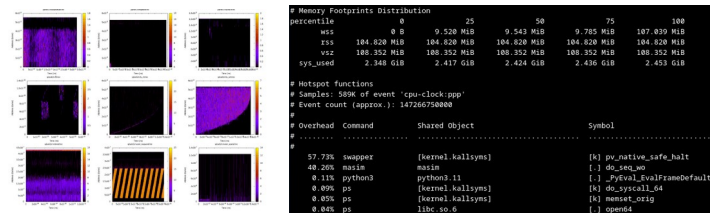
```
# Memory Footprints Distribution
percentile 0 25 50 75 100
rss 0 0 9.520 MiB 9.543 MiB 9.785 MiB 107.839 MiB
vss 104.820 MiB 104.820 MiB 104.820 MiB 104.820 MiB 104.820 MiB
vms 188.352 MiB 188.352 MiB 188.352 MiB 188.352 MiB 188.352 MiB
vsl_used 2.348 GiB 2.437 GiB 2.424 GiB 2.436 GiB 2.433 GiB

# Hotspot functions
# Samples: 586 of event 'cpu-clock:ppp'
# Event count (approx.): 14726750000
#
# Overhead Command Shared Object Symbol
# .....
#
# 57.73% swapper [kernel.kallsyms] [k] pv_native_safe_halt
# 48.20% masin [.] do_seq_no
# 0.11% python3 [python3] _pyeval_evalframeDefault
# 0.09% ps [kernel.kallsyms] [k] do_syscall_64
# 0.05% ps [kernel.kallsyms] [k] memset_orig
# 0.04% ps [libc.so.6] [.] open64
```



Potential, or Aimed Usages

- Profiling (e.g., GIF demo [link](#))
 - Help better understand and find rooms for improvements
- Profiling-guided Optimizations
 - Could be done on both offline and online
- Why not let kernel just (transparently) works?



DAMON: Kernel Subsystem for Data Access Monitoring and Access-aware System Operations

DAMOS: Data Access Monitoring-based Operation Schemes

- Another side of DAMON
- Let users define schemes
 - Memory operation actions to apply to regions of specific access pattern
- Once per user-defined time interval
 - find the regions of the condition from the snapshot and apply the action
- Finding optimum “access pattern” on dynamic environments is challenging
 - Uncontrolled DAMOS could byte you!

```
# # pageout memory regions that not accessed for >=5 seconds
# damo start --damos_action pageout --damos_access_rate 0% 0% --damos_age 5s max
```


DAMOS Filter: Fine-Control Access-aware System Operation Targets

- Define target memory with non-access-pattern information
 - Page level filters: anon, owned cgroup, hugepage, LRU-activeness
 - Non-page level filters: address
 - “pageout cold pages *of NUMA node 1 that associated with cgroup A and file-backed*”
 - Can be useful for fine-grained monitoring, too
 - (“stat”, instead of “pageout”)

DAMOS Quota: Control Access-aware System Operation Aggressiveness

- Six fixed thresholds (min/max size, access frequency, age) are unnecessary in many cases
- Setting thresholds flexibly and controlling aggressiveness works in many cases
 - Single control knob
- Quota set the aggressiveness limit as amount of memory to apply action per a time interval
- Access pattern based prioritization is applied under the quota
- “pageout cold pages *up to 100 MiB per second using <2% CPU time, coldest ones first*”

Quota Auto-tuning: Auto-tuned Access-aware System Operations

- Quota tuning is manual and repetitive
- Change the question for user: How to do (mechanism) → What to achieve (final goal)
- Let users specify goal of the quota as a value of a metrics
 - Metrics: PSI level, NUMA node memory utilization, workload's latency, bandwidth, TPS, ...
 - e.g., “reclaim cold pages aiming 0.5% memory PSI”
- DAMOS adjusts quota using feedback loop, for current value of the metric
 - e.g., If memory PSI is 0.1% increase quota for reclaiming cold pages (reclaim more warm pages)

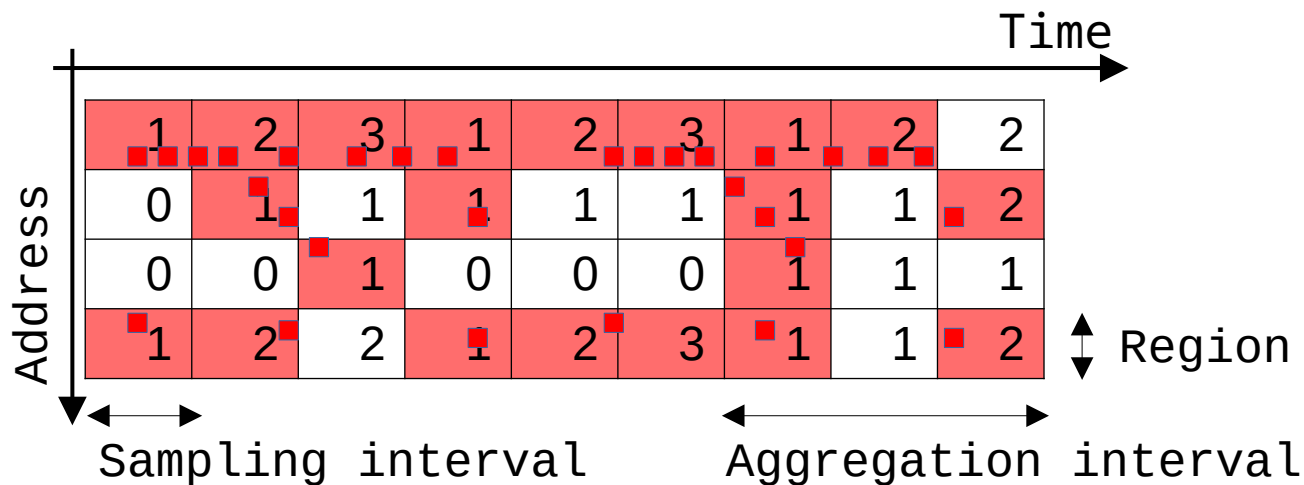
Auto-Tuning Time (Monitoring Intervals)

FAQ: DAMON output looks only cold, only hot, or just random

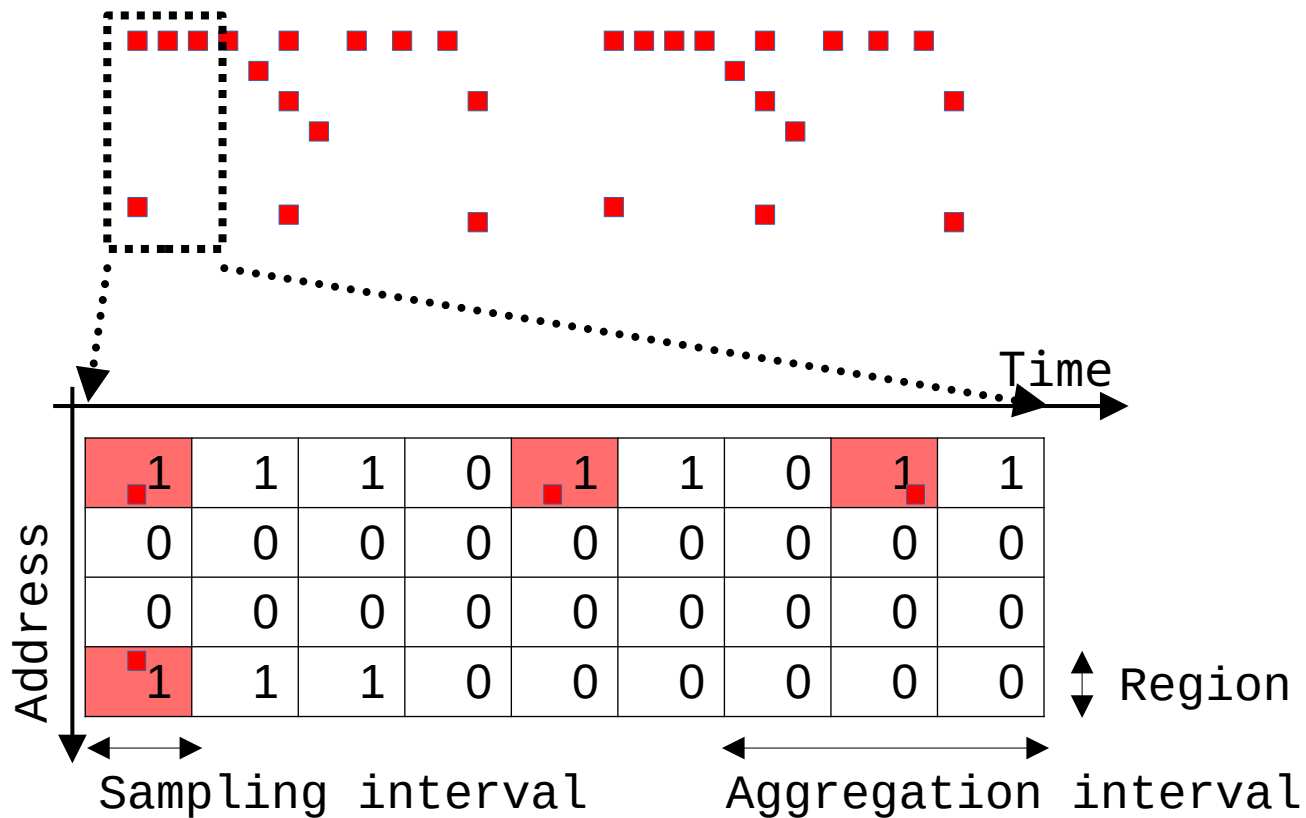
- Frequent Answer: Have you ~~turned DAMON off and on~~ tuned the monitoring intervals?
 - IOW, the default intervals (5ms sampling, 100ms aggregation) are not really suggested ones

If Intervals Are Appropriate: Meaningful Hot/Cold Regions

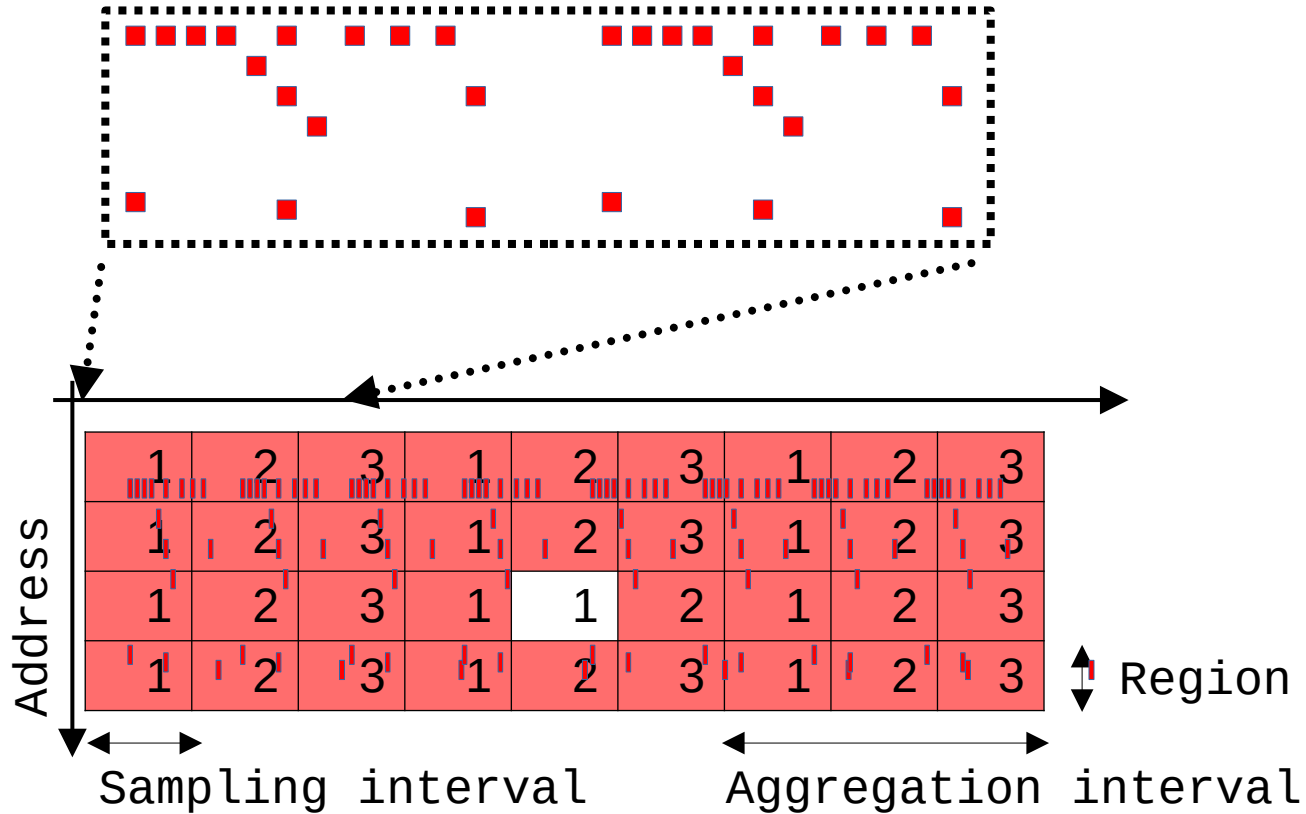
- Meaningful enough to make some memory management decisions



If Intervals Are Too Short: Everything Looks Cold

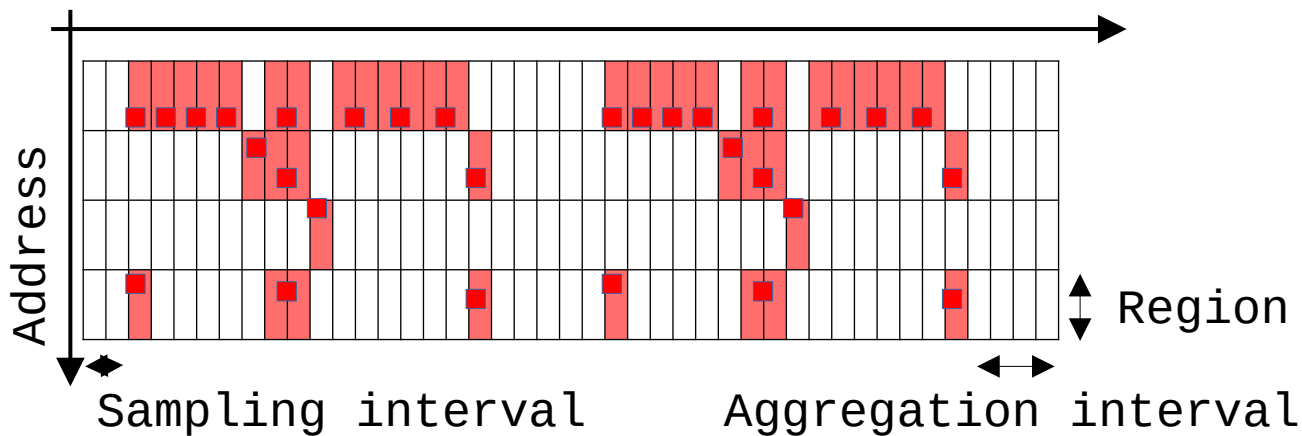


If Intervals Are Too Long: Everything Looks Hot



If Sampling:Aggregation Interval ratio is Too Low: Meaningless Samples

- Most sampling returns “negative”: unnecessary CPU cycle waste
- On large systems, sampling quality can also degrade
 - Not enough time for workloads to leave footprints



Intervals Tuning **Guideline**

- If it looks too cold, increase intervals
- If it looks too hot, decrease intervals
- Change both sampling and aggregation intervals in same ratio
- Repeat until meaningful snapshot is made
- Time consuming and repetitive for different workloads
- Guideline was added more than 3.5 years after DAMON be merged upstream!
- Could be automated?

Monitoring Intervals Auto-tuning (Motivated by Quotas Auto-tuning)

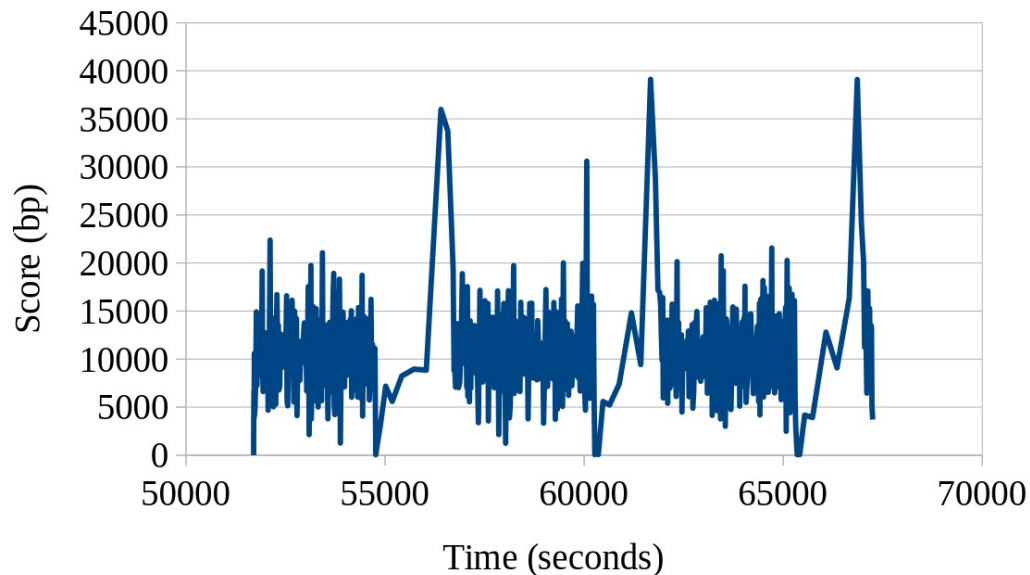
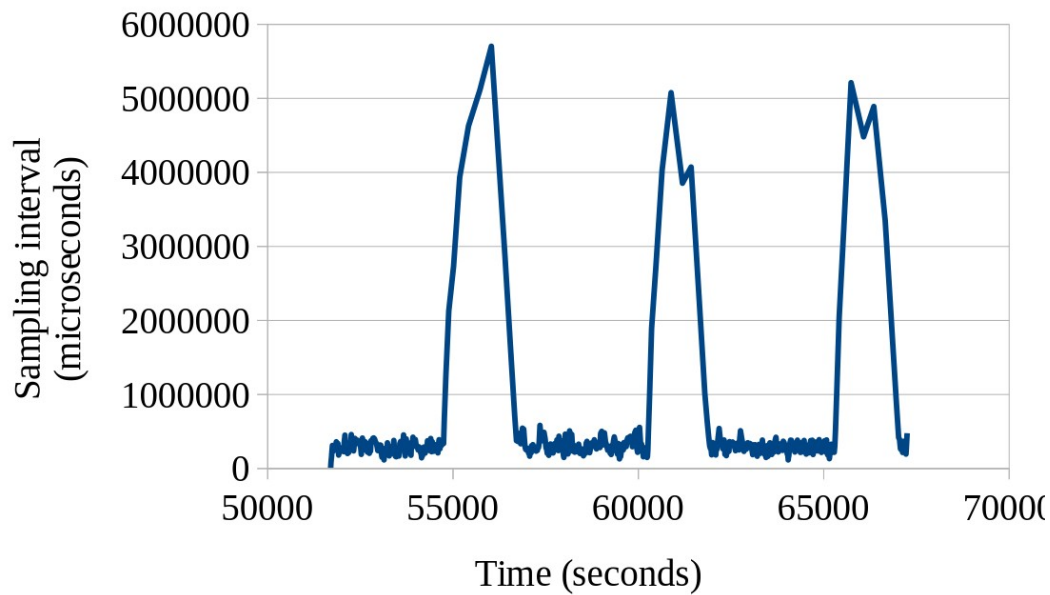
- Change Question: How to do? (mechanism) → What to achieve? (final goal, policy)
- Let users specify
 - Desired amount of access events to capture in each snapshot
 - Minimum and maximum sampling intervals
- Find proper sampling/aggregation intervals for the desire using feedback loop
 - If less than desired events are captured in current snapshot, increase intervals
 - If more than desired events are captured in current snapshot, decrease intervals
 - Min/max sampling intervals ensure auto-tune goes no too long
- Available on mainline from v6.15

Monitoring Intervals Auto-tuning Parameters

- Parameters for parameters auto-tuning, but easy to set
- Suggestion
 - Desired access events per snapshot: 4% of maximum events that can be captured in snapshot
 - Applies Pareto principle (80:20 rule) twice, assume to capture $80\% * 80\% = 64\%$ real access
 - Min/max sampling intervals: 5ms and 10s
 - Sampling:aggregation intervals ratio: 1:20
 - Only a few different actions are required, 20 is high enough

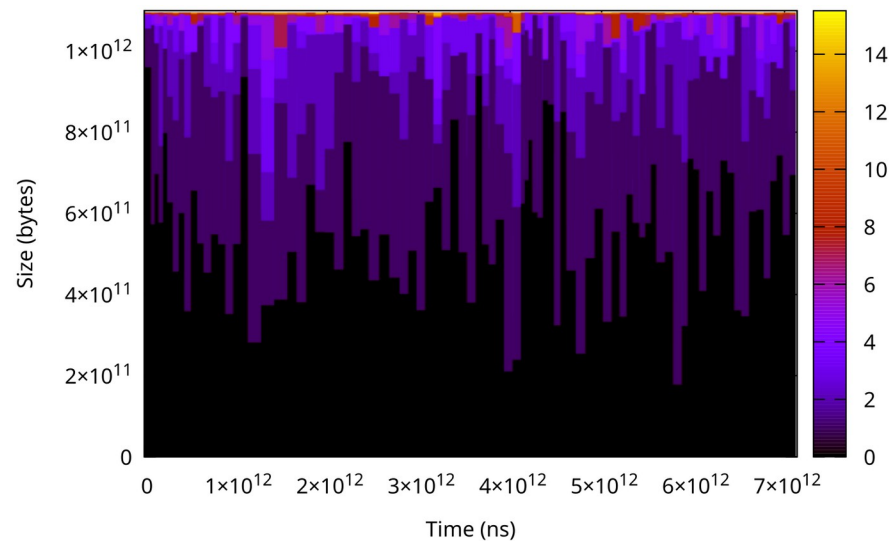
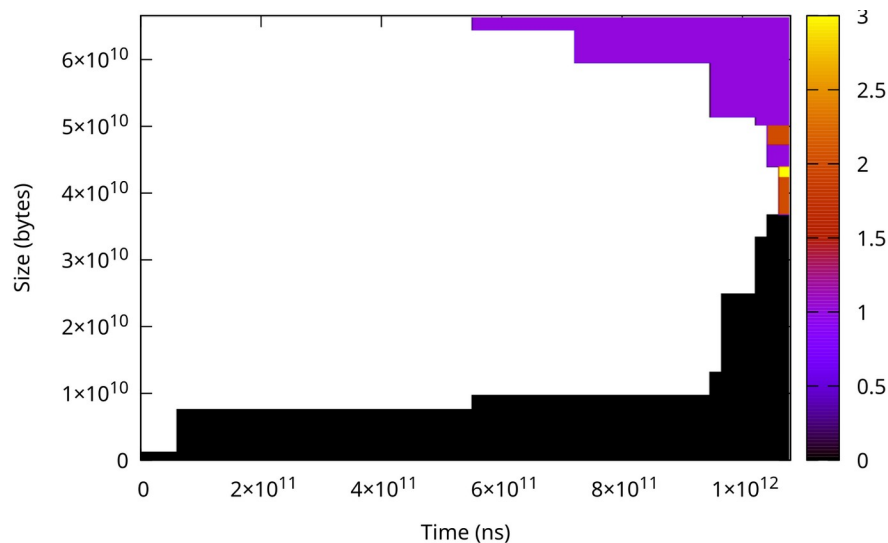
Intervals Auto-tuning on a Real-world Server Workload

- Sampling interval and tuning score continuously change, and converge for given situation
 - Sampling interval converges to 370ms under usual load, ~4-5 seconds under light load
 - Tuning score converges to the goal (10,000 bp)



Intervals Auto-Tuning on Real World Server Workloads

- Meaningful access patterns found on three different workloads including 1 TiB memory size workload
- 0.0% CPU time consumed for the monitoring

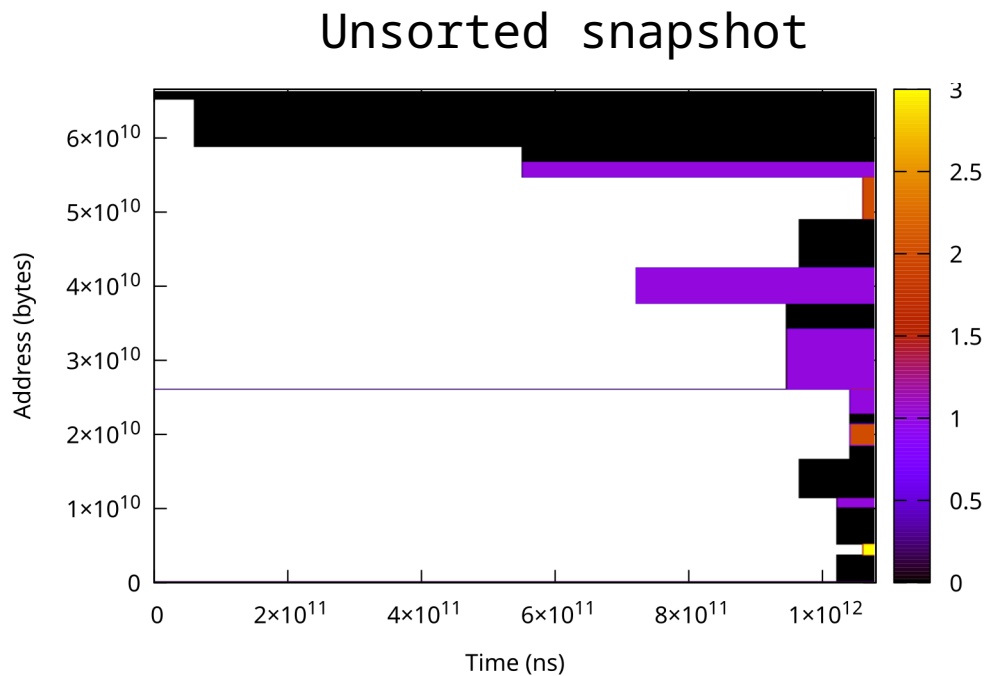


DAMON_STAT: Recommended Way For System-wide Access Monitoring

- Kernel module running DAMON for the entire physical address space
- Use intervals auto-tuning with the suggested auto-tune parameters
- Extract Idle time percentile
 - distribution of per-byte memory idle times (time the byte was not accessed)
 - P75 idle time 2minutes: 75 percent of the memory was accessed at least once in last 2 minutes; rest 25 percent of memory was not accessed at all for last 2 minutes
- Extract estimated memory bandwidth
 - Memory bandwidth estimated based on access events that captured in the last snapshot
- Recommended way for system-wide access monitoring
 - Easy to enable (`CONFIG_DAMON_STAT_DEFAULT_ENABLED=y`), aggregate, compare
 - Can be enabled/disabled at build, boot time and runtime

Idle Time Percentiles

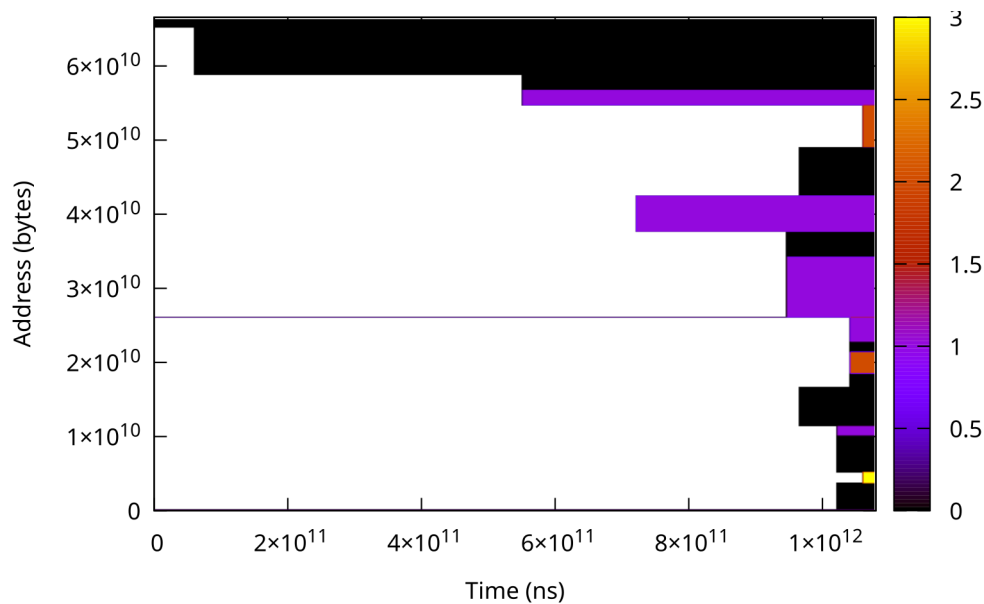
- Idle time: How long the region kept being not accessed (access frequency 0)
- Idle time percentiles: Percentiles of sorted per-byte idle times



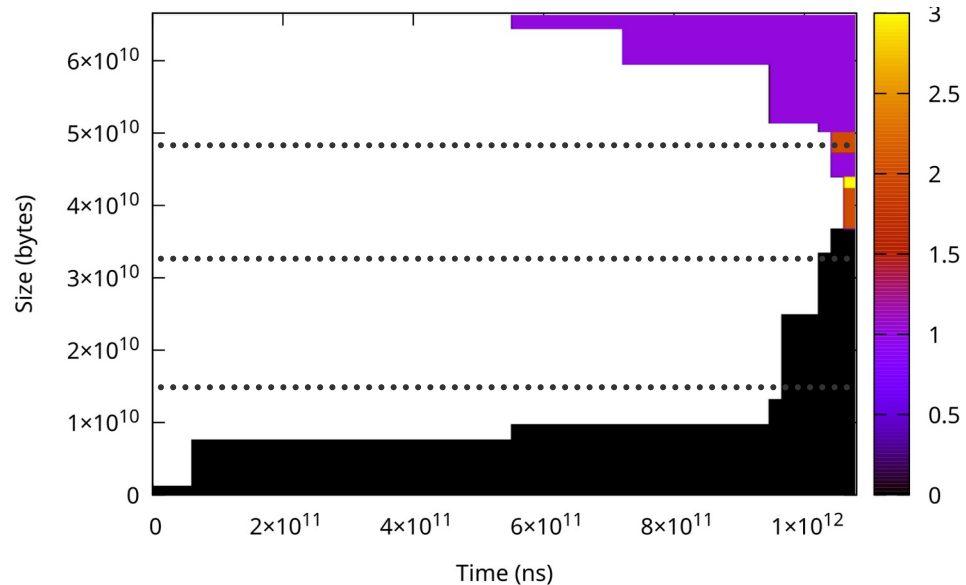
Idle Time Percentiles

- Idle time: How long the region kept being not accessed (access frequency 0)
- Idle time percentiles: Percentiles of sorted per-byte idle times

Unsorted snapshot



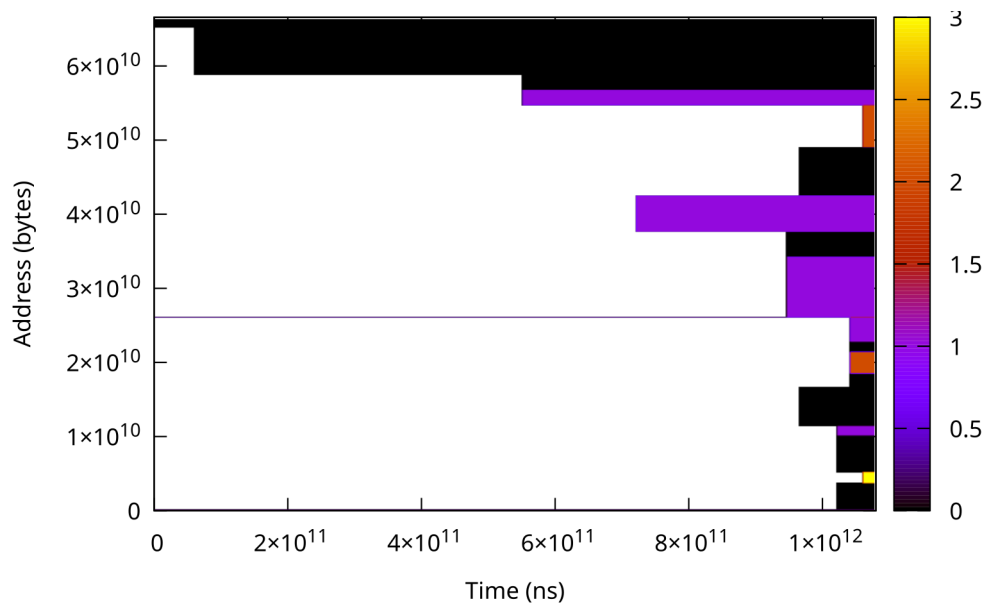
Sorted by access frequency



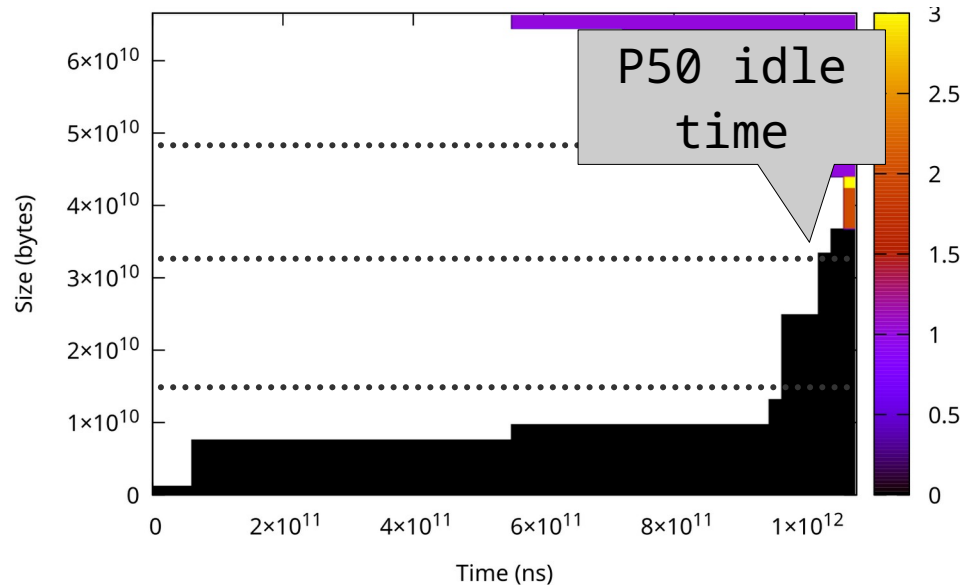
Idle Time Percentiles

- Idle time: How long the region kept being not accessed (access frequency 0)
- Idle time percentiles: Percentiles of sorted per-byte idle times

Unsorted snapshot



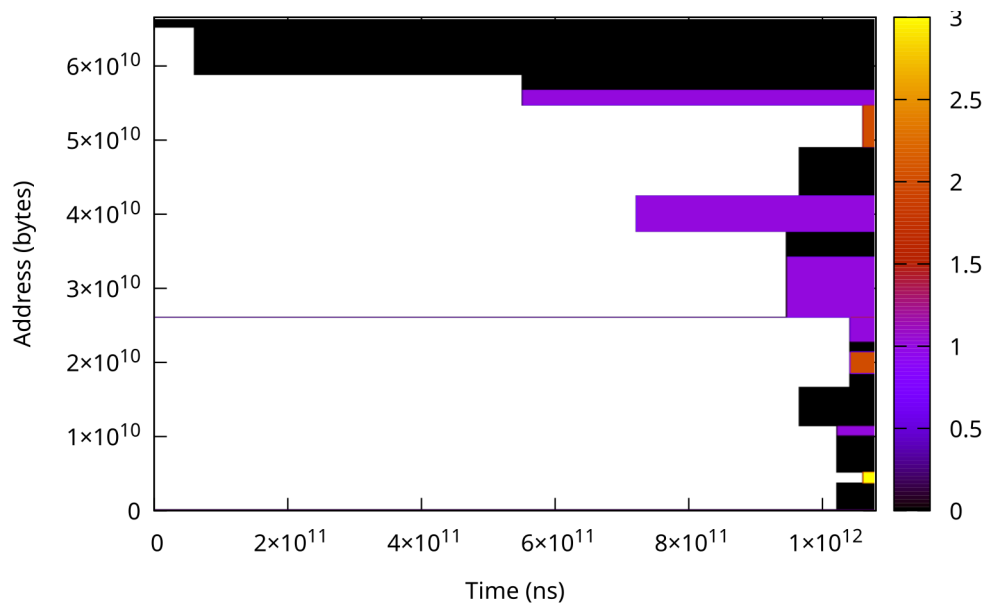
Sorted by access frequency



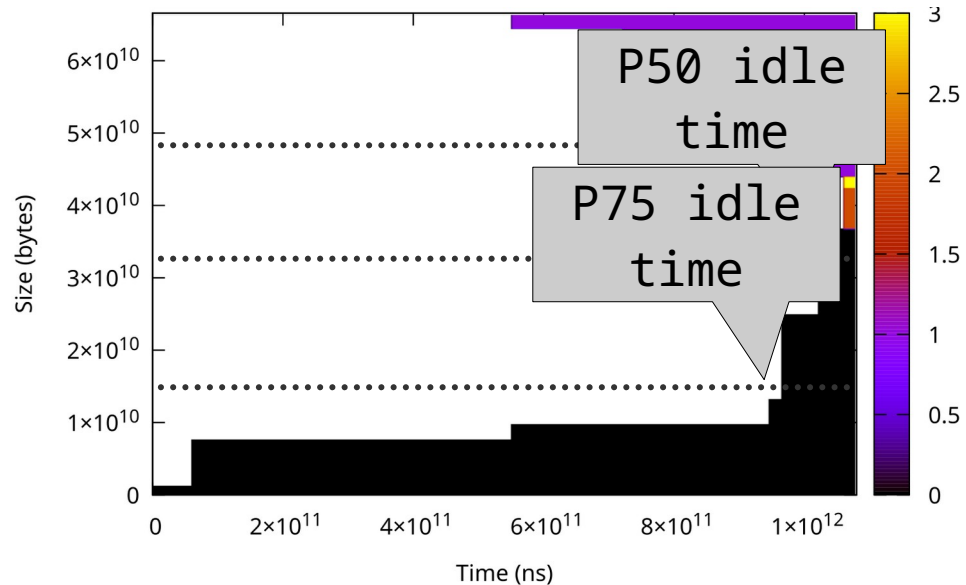
Idle Time Percentiles

- Idle time: How long the region kept being not accessed (access frequency 0)
- Idle time percentiles: Percentiles of sorted per-byte idle times

Unsorted snapshot

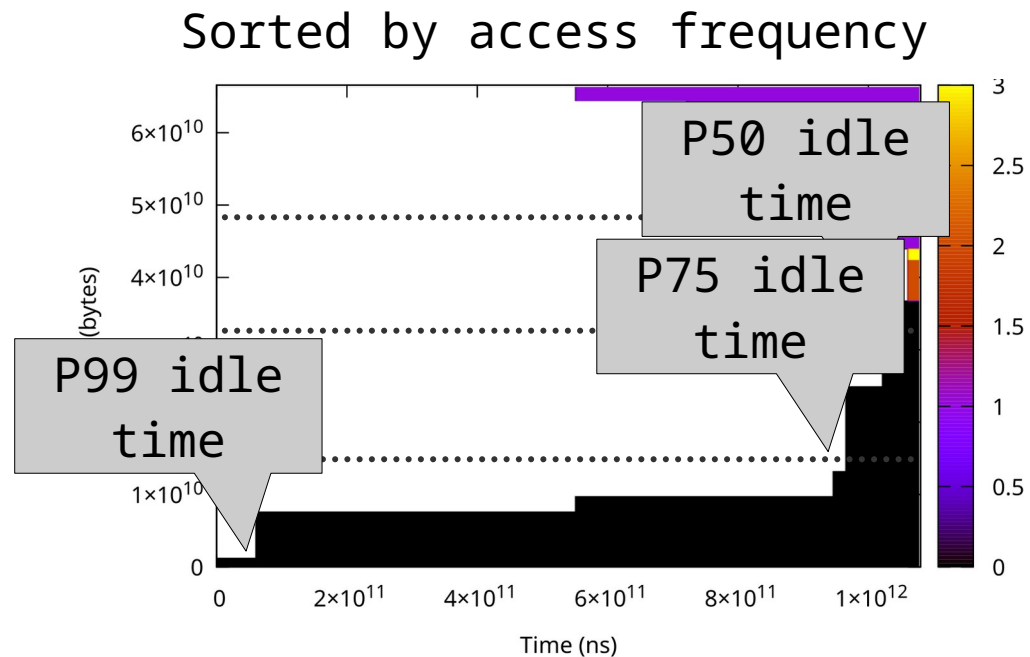
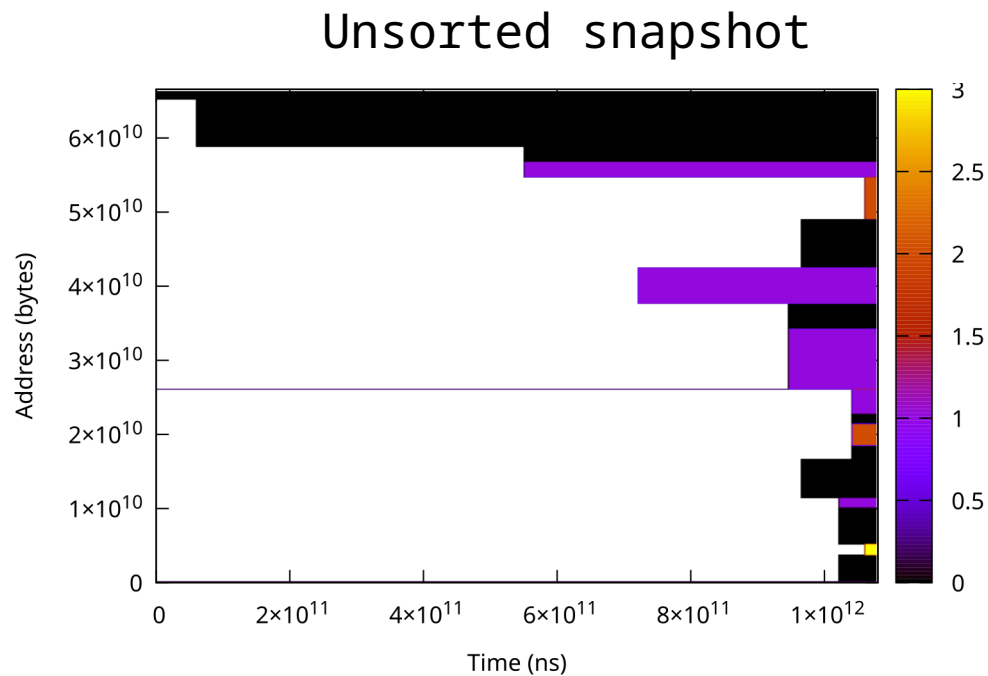


Sorted by access frequency



Idle Time Percentiles

- Idle time: How long the region kept being not accessed (access frequency 0)
- Idle time percentiles: Percentiles of sorted per-byte idle times

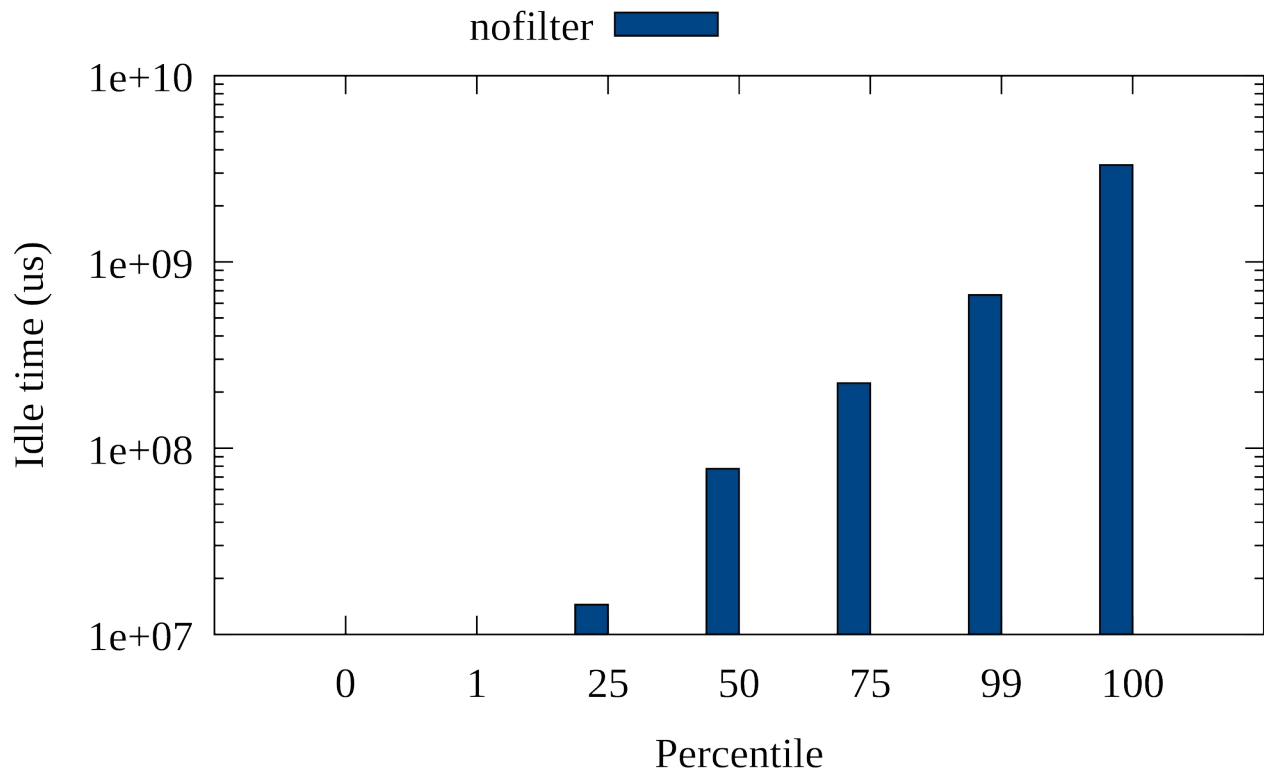


DAMON_STAT: Recommended Way For System-wide Access Monitoring

- Kernel module running DAMON for the entire physical address space
- Use intervals auto-tuning with the suggested auto-tune parameters
- Extract Idle time percentile
 - distribution of per-byte memory idle times (time the byte was not accessed)
 - P75 idle time 2minutes: 75 percent of the memory was accessed at least once in last 2 minutes; rest 25 percent of memory was not accessed at all for last 2 minutes
- Extract estimated memory bandwidth
 - Memory bandwidth estimated based on access events that captured in the last snapshot
- Recommended way for system-wide access monitoring
 - Easy to enable (`CONFIG_DAMON_STAT_DEFAULT_ENABLED=y`), aggregate, compare
 - Can be enabled/disabled at build, boot time and runtime

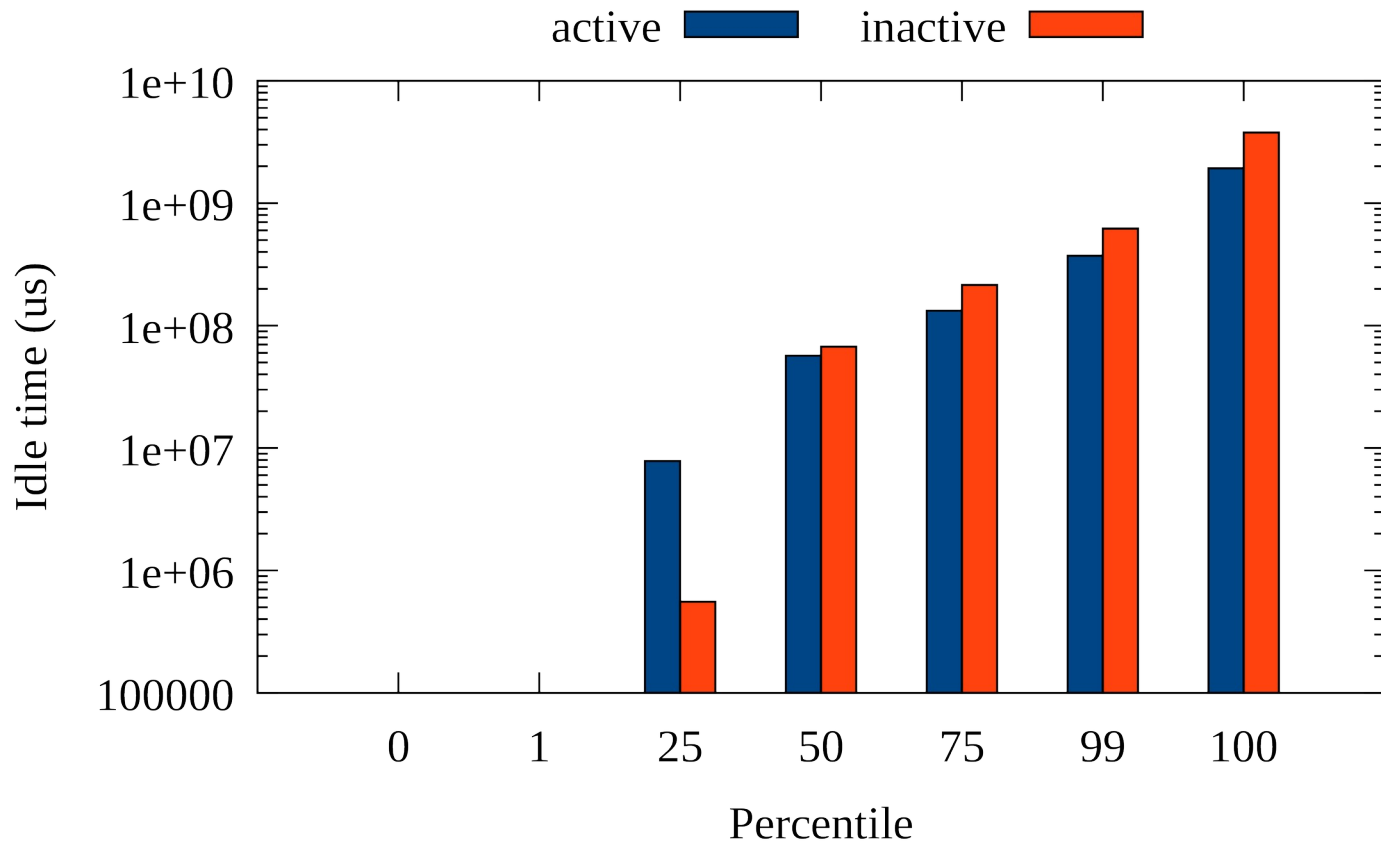
Results on a Real Workload: Auto-tuned Total Memory Idle Time Percentiles

- Small hot memory, exponentially increasing idle time (long tail of cold pages)



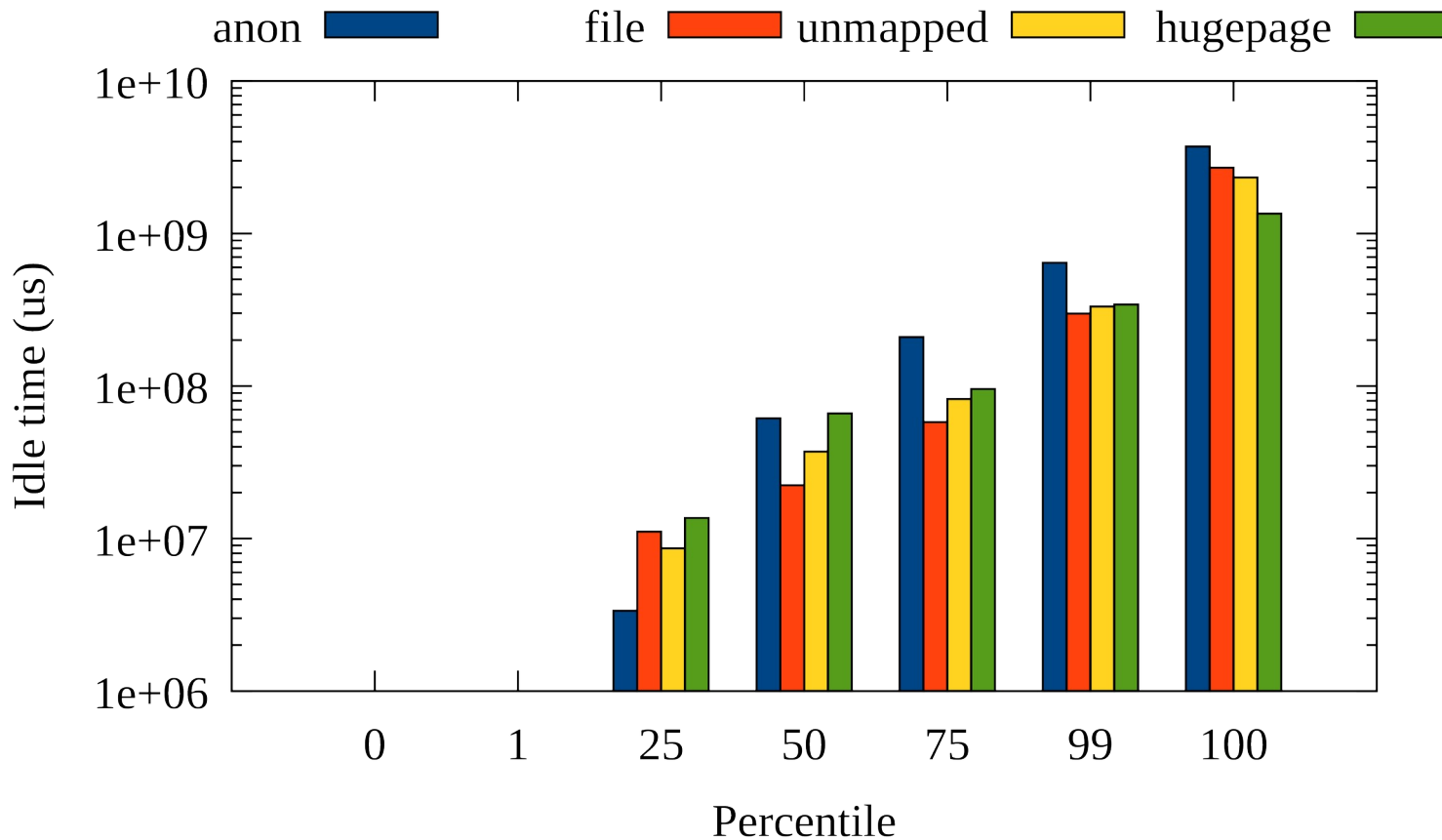
Results on a Real Workload: Active vs Inactive Pages Idle Time Breakdown

- Active pages have rooms to be more hot than inactive (ideally, p100 of active should < p0 of inactive)



Results on a Real Workload: Per Page Type Idle Time Breakdown

- You can check if your workload has expected access pattern



Getting Started

Availability

- Merged into the mainline from v5.15
 - Intervals auto-tuning is available from v6.15
- Backported and [enabled](#) on major Linux distro kernels
 - Major distros: Alma, Amazon, Android, Arch, CentOS, Debian, Fedora, Oracle, ...
- DAMON user-space tool (damo) is available on major packaging systems
 - Arch, Debian, Fedora, [PyPi](#), ...

Interfaces

- Kernel [API](#): Recommended for kernel programmers
- DAMON sysfs [interface](#): Recommended for user-space program development
- DAMON user-space [tool](#): Recommended for general usages from user-space
- DAMON [modules](#): Recommended for specific usages
 - DAMON_STAT for simple system-wide monitoring

Community: For Questions, Help, Patch Reviews

- Public mailing [list](https://lore.kernel.org/damon) (<https://lore.kernel.org/damon>)
- Bi-weekly virtual [meetup](#)
 - Occasional/regular private meetings on demand
- Project [website](https://damonitor.github.io) (<https://damonitor.github.io>)
 - Starting point for DAMON users and developers
- Not used to mail-based code review? Try [hkml](#)
 - Developed and maintained for DAMON and Linux kernel developers
- The future of DAMON is open and up to you
 - “Prefer random evolution over intelligent design”

Summary: That's DAMON

- DAMON is a kernel subsystem
 - For data access monitoring and access-aware system operations
- DAMON can be safely controlled and self-tunable
 - Try again with those features if you didn't, particularly for monitoring intervals auto-tuning
- The future is open and up to the community
 - Make your selfish voice

Questions?

- You can also ask questions anytime to
 - sj@kernel.org
 - Public mailing [list](https://lore.kernel.org/daemon) (<https://lore.kernel.org/daemon>)
 - Bi-weekly virtual [meetup](#)
 - Occasional/regular private meetings on demand
 - Project [website](https://damonitor.github.io) (<https://damonitor.github.io>)

Backup Slides

Controlled and Auto-tuned Access-aware System Operation Performance

- Parsec3/splash2x.fft
- Page out regions that not accessed for ≥ 5 seconds, up to 1GiB/sec, using up to 100ms/sec, aiming 10ms/sec memory pressure stall

	Runtime	RSS
Baseline	50.489s	10.005 GiB
+DAMOS-reclaim	120s	4.955 GiB
+Quota	51.772s	8.527 GiB
+Goal	49.741s	9.721 GiB

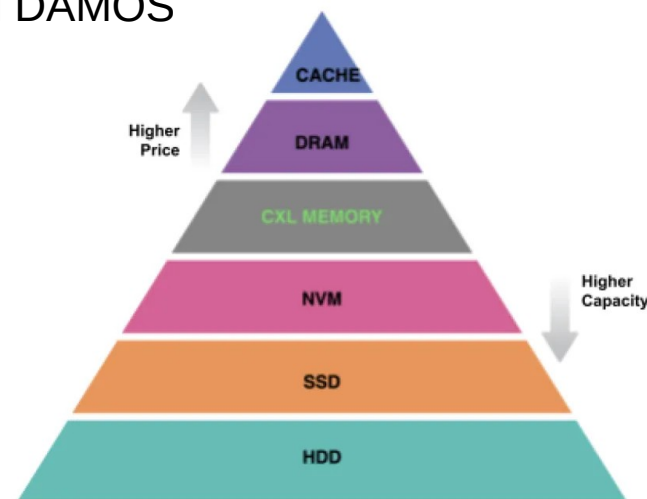
Real-world DAMON Use Cases: Proactive Reclamation and CXL Memory Tiering

Proactive Reclamation

- Reactive reclamation: Reclaim cold memory when memory pressure happens
- Proactively reclamation: Reclaim cold memory before memory pressure
- Benefit 1: Reduce memory footprint without performance degradation
- Benefit 2: Minimize degradation from direct reclamation
- Known usages: [Google](#), [Meta](#), and [Amazon](#)
 - Each company uses its own implementation for its usage
- AWS uses DAMOS-based implementation since 2022

CXL Memory Tiering

- CXL-tiered memory: Put CXL memory between DRAM and NVM
 - Pros: Higher capacity with lower price (higher efficiency)
- Challenge: Dynamic placement of pages (CXL mem is slower than DRAM)
- DAMON-based approach: Place hot pages on DRAM node, Place cold pages on CXL node
- SK hynix developed their CXL memory SDK (HMSDK) using DAMOS
 - Reports ~12.9% speed up



Architectures

Execution Model: Kernel Thread per Requests

- “struct damon_ctx”: Data structure for DAMON user input/output containing
 - User requests: target address space, address range, intervals, DAMOS schemes
 - Operation results: access snapshot, DAMOS stats
- “kdamond”: DAMON worker thread
 - Create one kdamond per “damon_ctx”
 - In future, could support multiple “damon_ctx” per kdamond
 - In future, could separate DAMOS to another thread (maybe useful for cgroup charging)
 - Allows async DAMON execution and multiple kdamonds (CPUs) scaling

Extensible Layers

User-space
Tools

DAMO

datop

DAMON API User
Kernel Modules

General-purpose User ABI

Special-purpose Modules

DAMON_SYSFS

DAMON_DBGFS

DAMON_RECLAIM

DAMON_LRU_SORT

DAMON_WSS

DAMON Application Programming Interface

DAMON

DAMOS

Adaptive Regions Adjustment

Action and Pattern

Region-based Sampling

Quotas and Prioritization

Access Frequency Monitoring

Feedback-based auto-tuning

Advanced Regions Adjustment

Watermarks

Parameters Auto-tuning

Filters

DAMON Operations Set Registration Interface

Operations Set

paddr

vaddr

Read/write-only

NUMA-cpus-only

Primitives
that DAMON
depends on

PTE/VMA/rmap, ...

AMD IBS

LRU State

Extensible Layers

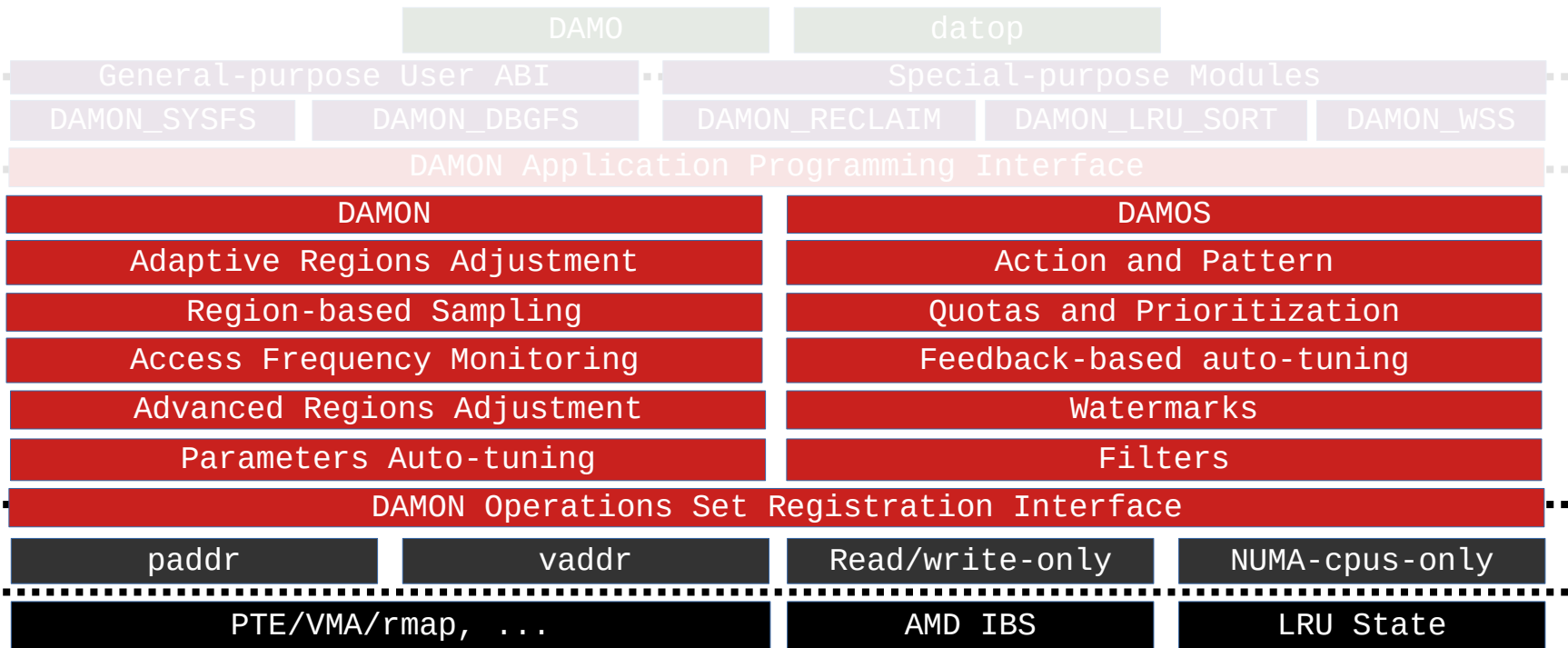
User-space
Tools

DAMON API User
Kernel Modules

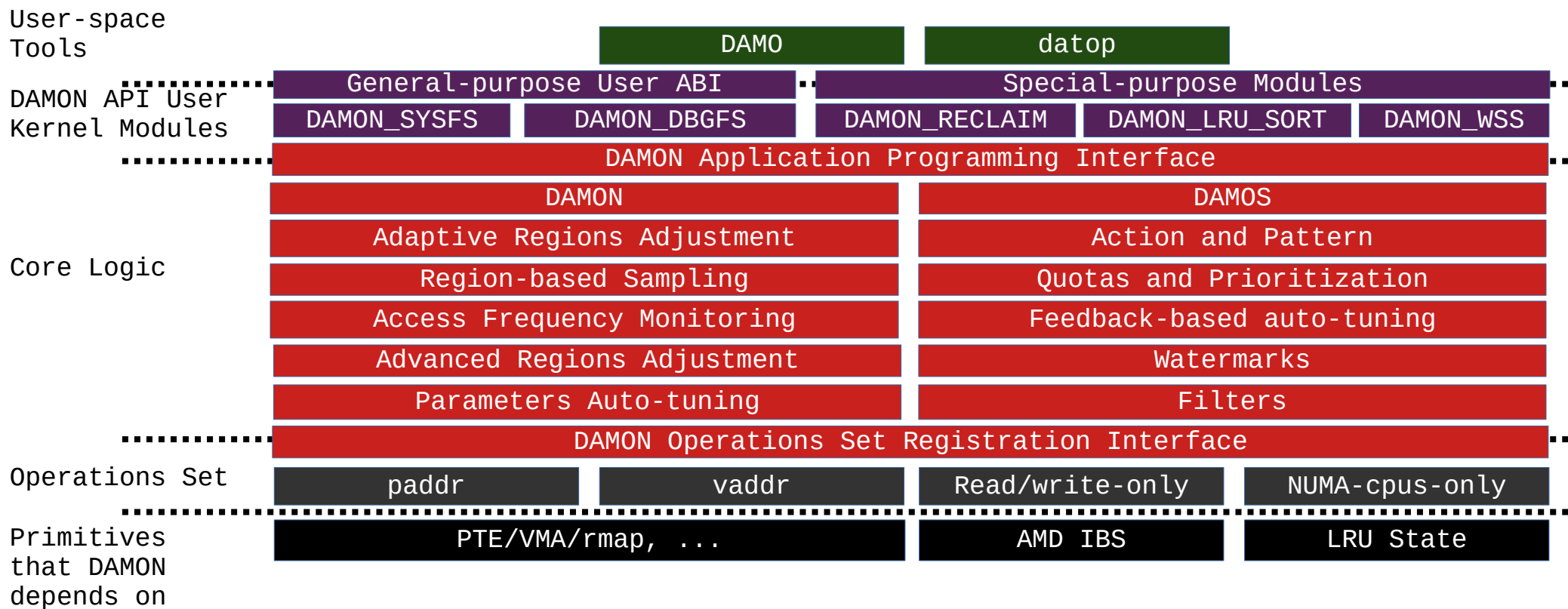
Core Logic

Operations Set

Primitives
that DAMON
depends on



Extensible Layers



DAMOS Quotas: Intuitive Aggressiveness Control

- Before applying DAMOS schemes
 - Set temperature-based priority score of each region
 - Build “priority score”: “total size of regions of the priority” histogram
 - Find lowest priority threshold for the scheme meeting the quota
 - Skip applying action to regions having lower-than-threshold priority scores
- Single snapshot and histogram iteration: $O(\leq \text{user-defined-N})$
- Quota auto-tuning: A simple proportional feedback algorithm
 - Reward metrics: Arbitrary user-input or self-retrievable metrics like memory PSI

DAMOS Filters: Fine-grained Target Selection

- Before applying DAMOS action, check the properties of region and skip action if needed
- Non-page granular (high level) filters
 - Filtered out before applying actions
 - Address ranges (e.g., NUMA nodes or Zone)
 - DAMON-defined monitoring target (e.g., process)
- Page granular (low level) filters
 - Filtered out in the middle of actions in page level
 - Anon/File-backed
 - Belonging memory cgroup
 - `page_idle()`

Pseudo-code of DAMON v5.15

```
While True:
    for region in regions:
        if region.accessed():
            region.nr_accesses += 1
    sleep(sampling_interval)
    if now() % aggregation_interval:
        merge(regions)
        user_callback(regions)
        for region in regions:
            region.nr_accesses = 0
        split(regions)
```

DAMON accuracy on Low-locality Space/Workloads

- It is proven to work on real world products for years
- Pareto principle and unconscious bias will make the pattern
 - Entropy-full situation is when the data center is doom-ed
- “age” avoid immature decision
- More [works](#) for accuracy improvement will be continued
- DAMON could be decoupled with the region-based mechanisms in future
- Let's collect data and continue discussions together

Can DAMON Extended for Non-snapshot Access Patterns?

- TL; DR: Yes, why not?
- DAMON is for any access information; Snapshot is one of the representations
- If the information/representation is useful for users, DAMON can add support
- We started discussion for Memory bandwidth visibility

Can DAMON Use features Other than PTE Accessed bits?

- The extensible layer allows it
- AMD IBS and page fault-based approaches (e.g., PTE_NONE) are on the table
- In future, if GPU provides access check feature, we can extend to use it
- Such extension would allow
 - More lightweight and precise monitoring
 - Access source, read/write-aware monitoring
 - Kernel memory access monitoring

DAMOS for Efficient and Fine-grained Data Access Monitoring

- DAMOS_STAT
 - Special action making no system change but expose the scheme-internal information
 - Let user knows which of the memory are eligible for the scheme
- With DAMOS filters, can do page level properties-based monitoring
 - “How much of >2 minutes unaccessed memory are in hugepages and belong to cgroup A?”
- With DAMOS quotas, can do overhead-controlled monitoring