

Write-Aware Timestamp Tracking:

Effective and Efficient Page Replacement for Modern Hardware

Demian Vöhringer¹ Viktor Leis²

¹Friedrich-Alexander-Universität Erlangen-Nürnberg

²Technische Universität München

August 29, 2023 – VLDB

Challenges

- **Storage:** HDD → SSD
- **CPU:** single-core → multi-core

Challenges

- **Storage:** HDD → SSD
- **CPU:** single-core → multi-core

4 Goals of Page Replacement

Challenges

- **Storage:** HDD → SSD
- **CPU:** single-core → multi-core

4 Goals of Page Replacement

1. Replacement Effectiveness

Challenges

- **Storage:** HDD → SSD
- **CPU:** single-core → multi-core

4 Goals of Page Replacement

1. **Replacement Effectiveness**
2. **Write Awareness**

Challenges

- **Storage:** HDD → SSD
- **CPU:** single-core → multi-core

4 Goals of Page Replacement

1. **Replacement Effectiveness**
2. **Write Awareness**
3. **CPU Efficiency**

Challenges

- **Storage:** HDD → SSD
- **CPU:** single-core → multi-core

4 Goals of Page Replacement

1. **Replacement Effectiveness**
2. **Write Awareness**
3. **CPU Efficiency**
4. **Multi-Core Scalability**

Challenges

- **Storage:** HDD → SSD
- **CPU:** single-core → multi-core

4 Goals of Page Replacement

1. **Replacement Effectiveness**
2. **Write Awareness**
3. **CPU Efficiency**
4. **Multi-Core Scalability**

Used Strategies

- **LRU:** DB2, Oracle, InnoDB (MySQL / MariaDB), WiredTiger (MongoDB)
- **LRU2:** SQL Server
- **Clock-Sweep:** PostgreSQL

	Replacement Effectiveness	Write Awareness	CPU Efficiency	Multi-Core Scalability
random	— —	no	++	++
CLOCK	—	no	+	+
LeanEvict	—	no	+	+
LRU2	+	no	++	+
LRU	—	no	—	—
CFLRU	—	yes	—	—
LRU_WSR	—	yes	—	—
Hyperbolic	~	no	++	++
ARC	+	no	—	—

more information

more information → more insights

more information → more insights → better strategy

more information \rightarrow more insights \rightarrow better strategy

Index[i]

Timestamp[t_i]

1 2 3 4 ...

-	-	-	-	...
---	---	---	---	-----

Tracking Array (per page)

more information \rightarrow more insights \rightarrow better strategy

i
Access at $t = 0$

1	2	3	4	...
0	-	-	-	...

Tracking Array (per page)

more information \rightarrow more insights \rightarrow better strategy

i

Access at $t = 8$

1 2 3 4 ...

8	0	-	-	...
---	---	---	---	-----

Tracking Array (per page)

more information \rightarrow more insights \rightarrow better strategy

i
Access at $t = 15$

1	2	3	4	...
15	8	0	-	...

Tracking Array (per page)

more information \rightarrow more insights \rightarrow better strategy

i
Access at $t = 42$

1	2	3	4	...
42	15	8	0	...

Tracking Array (per page)

How to use it to evict a page?

Page Value [PV]: high for good pages, low for bad ones

How to use it to evict a page?

Page Value [PV]: high for good pages, low for bad ones
approximate current access frequency

How to use it to evict a page?

Page Value [PV]: high for good pages, low for bad ones
approximate current access frequency

$$PV_{LRU}(t_{now}) = \frac{1}{age} = \frac{1}{t_{now} - t_1}$$

$$PV_{LRU2}(t_{now}) = \frac{1}{t_{now} - t_2}$$

$$PV_{LFU}(t_{now}) = len(\text{Tracking Array})$$

$$\textit{Subfrequency}[SF]_i(t_{now}) := \frac{i}{t_{now} - t_i}$$

$$\textit{Subfrequency}[SF]_i(t_{now}) := \frac{i}{t_{now} - t_i}$$

$$SF_1(t_{now}) := \frac{1}{t_{now} - t_1}$$

$$SF_2(t_{now}) := \frac{2}{t_{now} - t_2}$$

$$\textit{Subfrequency}[SF]_i(t_{now}) := \frac{i}{t_{now} - t_i}$$

$$SF_1(t_{now}) := \frac{1}{t_{now} - t_1}$$

$$SF_2(t_{now}) := \frac{2}{t_{now} - t_2}$$

i
 t_i

1	2	3	4
42	15	8	0

Evaluated at $t_{now} := 50$

$$\textit{Subfrequency}[SF]_i(t_{now}) := \frac{i}{t_{now} - t_i}$$

$$SF_1(t_{now}) := \frac{1}{t_{now} - t_1}$$

$$SF_2(t_{now}) := \frac{2}{t_{now} - t_2}$$

i	1	2	3	4
t_i	42	15	8	0
$t_{now} - t_i$	8	35	42	50

Evaluated at $t_{now} := 50$

$$\textit{Subfrequency}[SF]_i(t_{now}) := \frac{i}{t_{now} - t_i}$$

$$SF_1(t_{now}) := \frac{1}{t_{now} - t_1}$$

$$SF_2(t_{now}) := \frac{2}{t_{now} - t_2}$$

i	1	2	3	4
t_i	42	15	8	0
$t_{now} - t_i$	8	35	42	50
$SF_i(t_{now})$	1/8	2/35	3/42	4/50

Evaluated at $t_{now} := 50$

$$\textit{Subfrequency}[SF]_i(t_{now}) := \frac{i}{t_{now} - t_i}$$

$$SF_1(t_{now}) := \frac{1}{t_{now} - t_1}$$

$$SF_2(t_{now}) := \frac{2}{t_{now} - t_2}$$

i	1	2	3	4
t_i	42	15	8	0
$t_{now} - t_i$	8	35	42	50
$SF_i(t_{now})$	1/8	2/35	3/42	4/50
$SF_i(t_{now}) \approx$	0.13	0.06	0.07	0.08

Evaluated at $t_{now} := 50$

$$\textit{Page Value}[PV](t_{now}) := \max_i SF_i(t_{now})$$

$$\textit{Page Value}[PV](t_{now}) := \max_i SF_i(t_{now})$$

i	1	2	3	4
$SF_i(50) \approx$	0.13	0.06	0.07	0.08

$$PV(50) = SF_1(50) \approx 0.13$$

$$PV_{LRU}(t_{now}) = \frac{1}{age} = \frac{1}{t_{now} - t_1}$$

$$PV_{LRU2}(t_{now}) = \frac{1}{t_{now} - t_2}$$

$$PV_{LFU}(t_{now}) = len(\text{Tracking Array})$$

$$PV_{WATT}(t_{now}) = \max_i \frac{i}{t_{now} - t_i}$$

$$PV_{LRU}(t_{now}) = \frac{1}{age} = \frac{1}{t_{now} - t_1} = SF_1(t_{now})$$

$$PV_{LRU2}(t_{now}) = \frac{1}{t_{now} - t_2}$$

$$PV_{LFU}(t_{now}) = len(\text{Tracking Array})$$

$$PV_{WATT}(t_{now}) = \max_i \frac{i}{t_{now} - t_i}$$

$$PV_{LRU}(t_{now}) = \frac{1}{age} = \frac{1}{t_{now} - t_1} = SF_1(t_{now})$$

$$PV_{LRU2}(t_{now}) = \frac{1}{t_{now} - t_2} \equiv SF_2(t_{now})$$

$$PV_{LFU}(t_{now}) = len(\text{Tracking Array})$$

$$PV_{WATT}(t_{now}) = \max_i \frac{i}{t_{now} - t_i}$$

$$PV_{LRU}(t_{now}) = \frac{1}{age} = \frac{1}{t_{now} - t_1} = SF_1(t_{now})$$

$$PV_{LRU2}(t_{now}) = \frac{1}{t_{now} - t_2} \equiv SF_2(t_{now})$$

$$PV_{LFU}(t_{now}) = len(\text{Tracking Array}) = \max_i \frac{i}{const} =: SF_{max}^*(t_{now})$$

$$PV_{WATT}(t_{now}) = \max_i \frac{i}{t_{now} - t_i}$$

⚡ Write Awareness

⚡ Memory Consumption

⚡ Bursts

⚡ Scans and Onetime Accesses

⚡ Finding Eviction Candidates

⚡Write Awareness

$$PV(t_{now}) := PV_{access}^*(t_{now}) + write_weight \cdot PV_{write}^*(t_{now})$$

⚡Memory Consumption

⚡Bursts

⚡Scans and Onetime Accesses

⚡Finding Eviction Candidates

⚡Write Awareness

$$PV(t_{now}) := PV_{access}^*(t_{now}) + write_weight \cdot PV_{write}^*(t_{now})$$

⚡Memory Consumption

limit Tracking Array size

⚡Bursts

⚡Scans and Onetime Accesses

⚡Finding Eviction Candidates

⚡Write Awareness

$$PV(t_{now}) := PV_{access}^*(t_{now}) + write_weight \cdot PV_{write}^*(t_{now})$$

⚡Memory Consumption

limit Tracking Array size

⚡Bursts

group accesses to epochs

⚡Scans and Onetime Accesses

⚡Finding Eviction Candidates

⚡Write Awareness

$$PV(t_{now}) := PV_{access}^*(t_{now}) + write_weight \cdot PV_{write}^*(t_{now})$$

⚡Memory Consumption

limit Tracking Array size

⚡Bursts

group accesses to epochs

⚡Scans and Onetime Accesses

dampening of $SF_1(t_{now}) = SF_1^*(t_{now})/10$

⚡Finding Eviction Candidates

⚡Write Awareness

$$PV(t_{now}) := PV_{access}^*(t_{now}) + write_weight \cdot PV_{write}^*(t_{now})$$

⚡Memory Consumption

limit Tracking Array size

⚡Bursts

group accesses to epochs

⚡Scans and Onetime Accesses

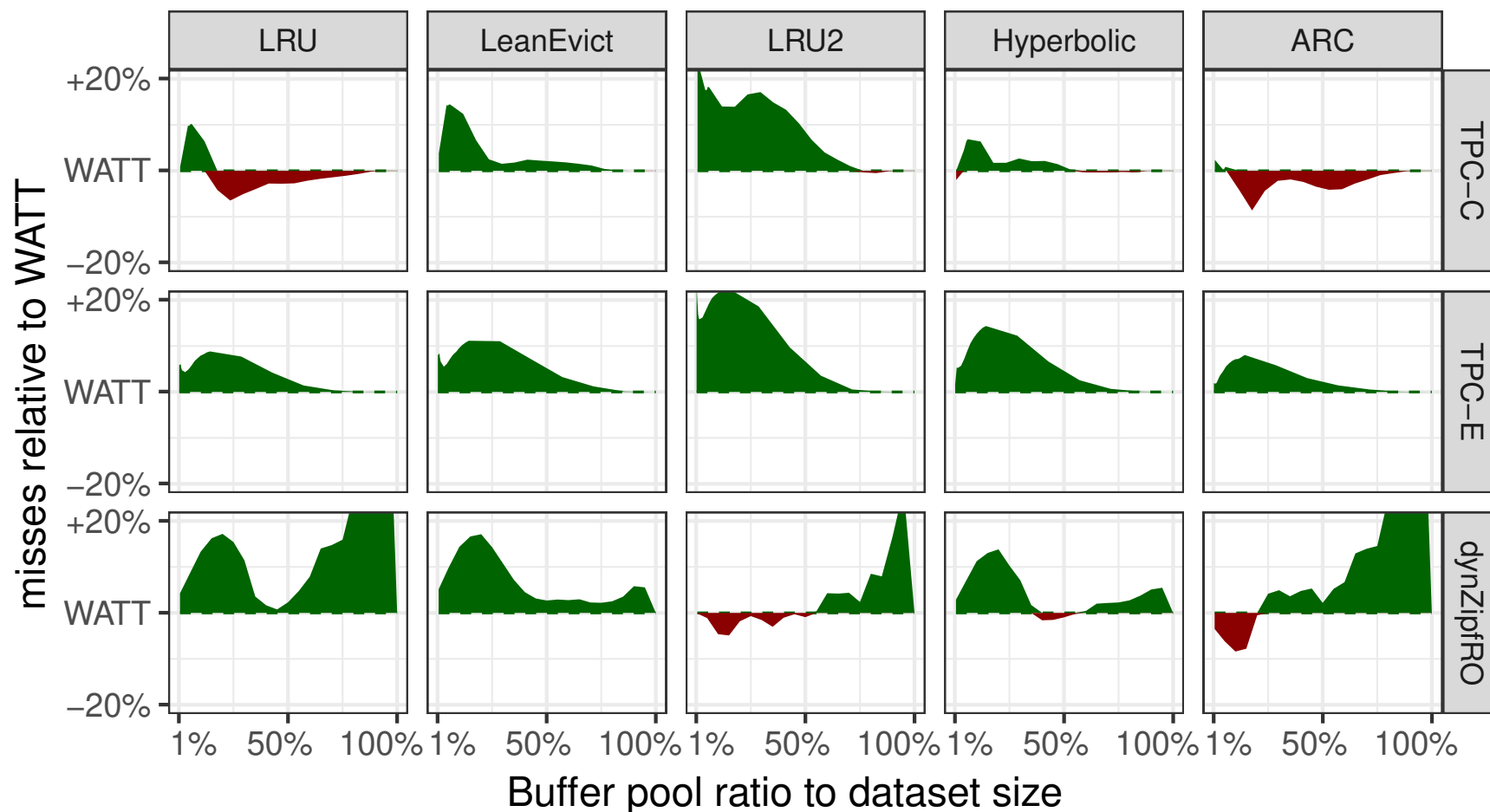
dampening of $SF_1(t_{now}) = SF_1^*(t_{now})/10$

⚡Finding Eviction Candidates

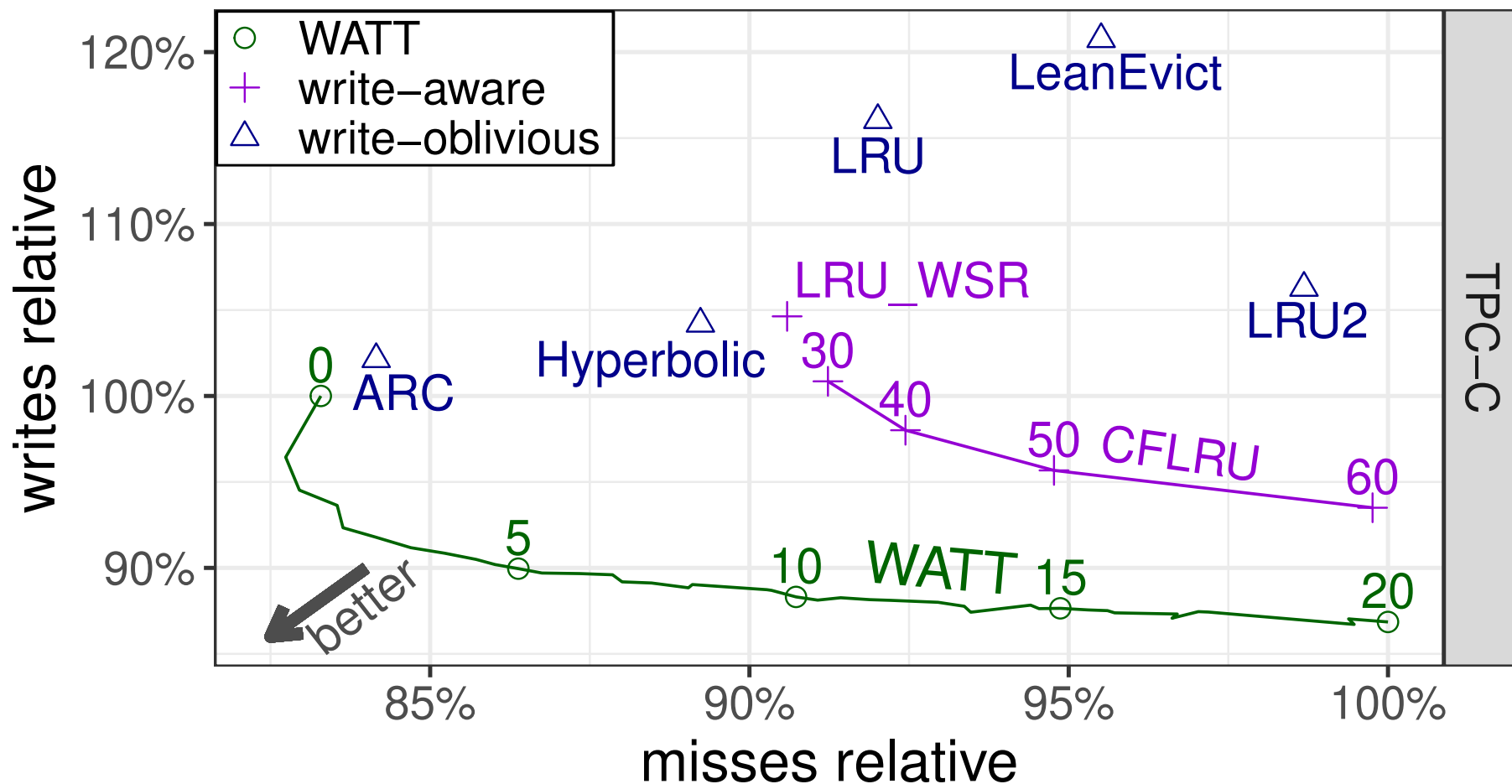
Random Sampling

Replacement Effectiveness (Goal 1)

Replacement Effectiveness (Goal 1)



Write Awareness (Goal 2)




Implementation

- implemented in  leanstore

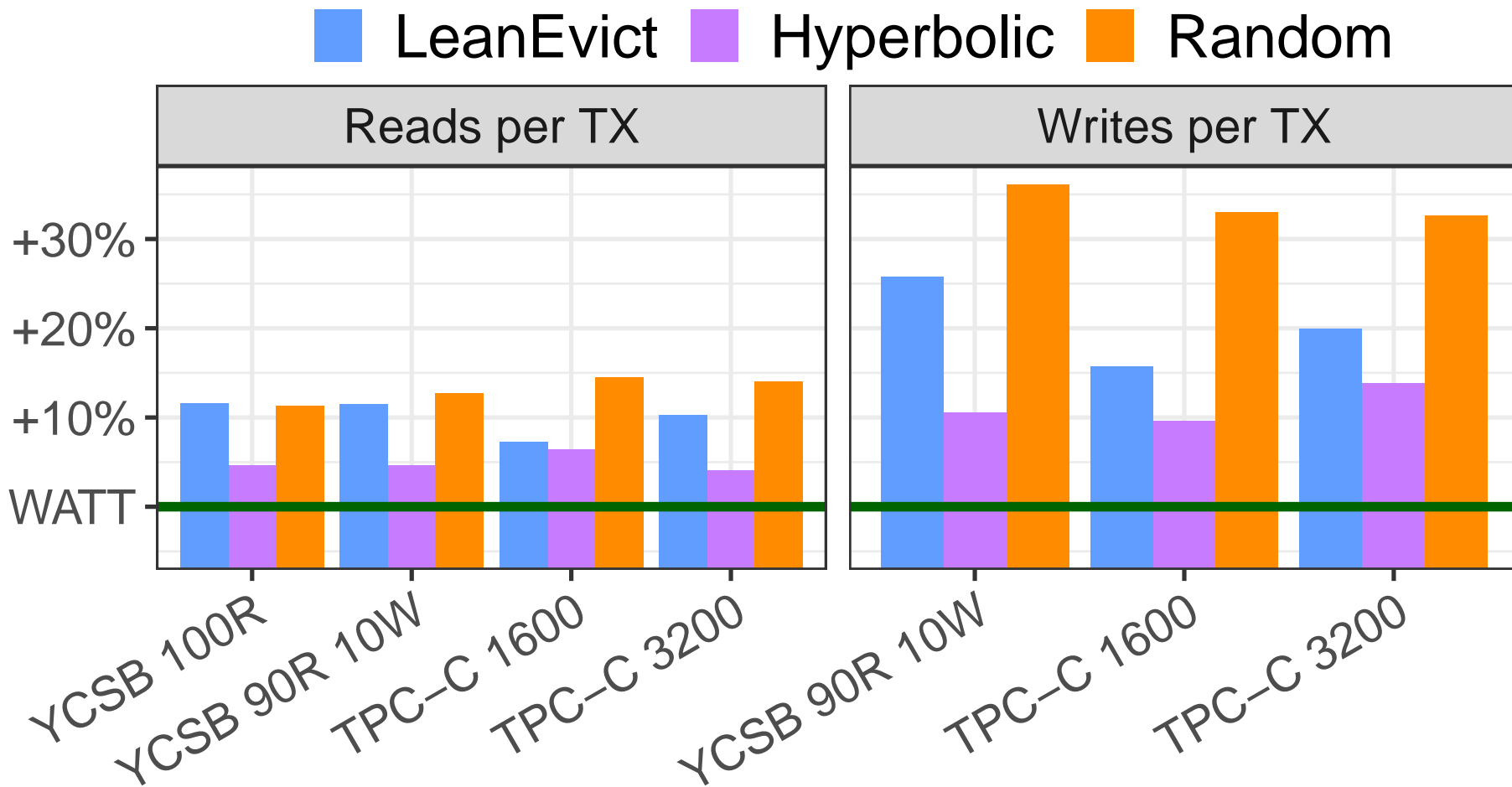
- implemented in  leanstore
- lock free access to Tracking Array

- implemented in  leanstore
- lock free access to Tracking Array
- SIMD vectorization to calculate Page Value

- implemented in  leanstore
- lock free access to Tracking Array
- SIMD vectorization to calculate Page Value
- Prefetching of Tracking Arrays for Page Value calculation

Replacement Effectiveness (Goal 1)

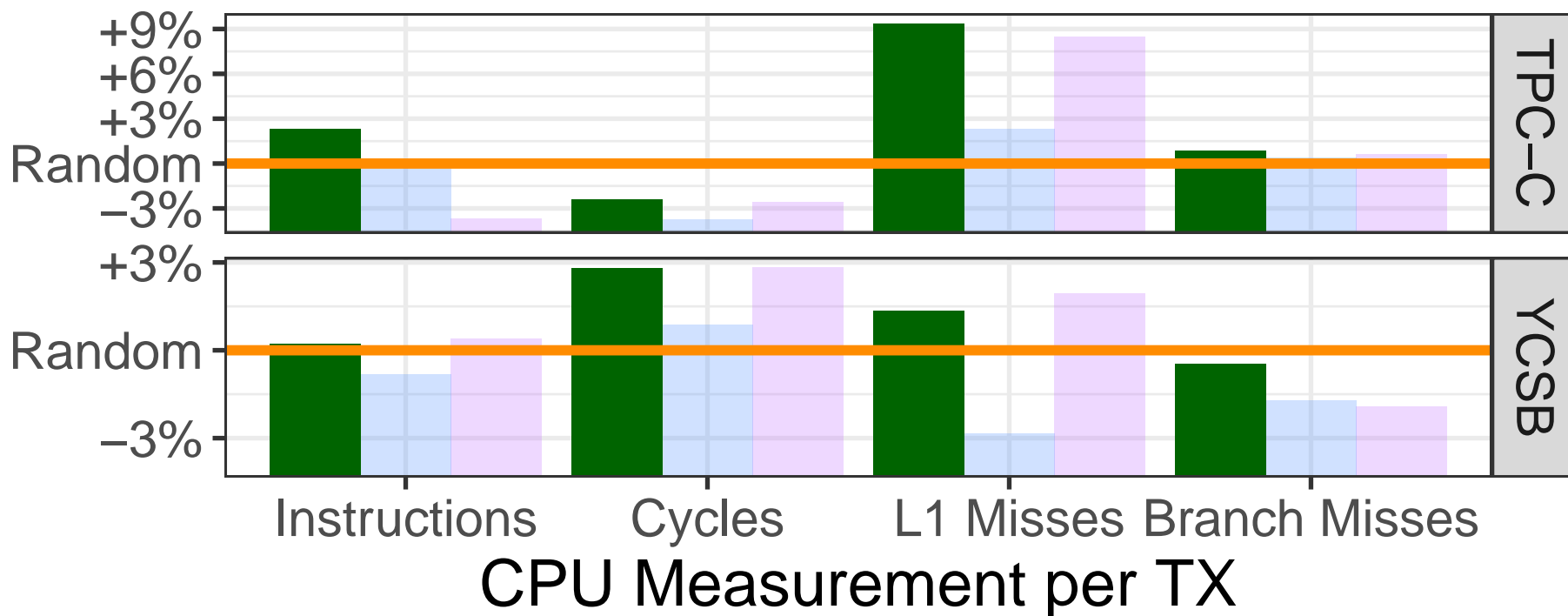
Write Awareness (Goal 2)



CPU Efficiency (Goal 3)

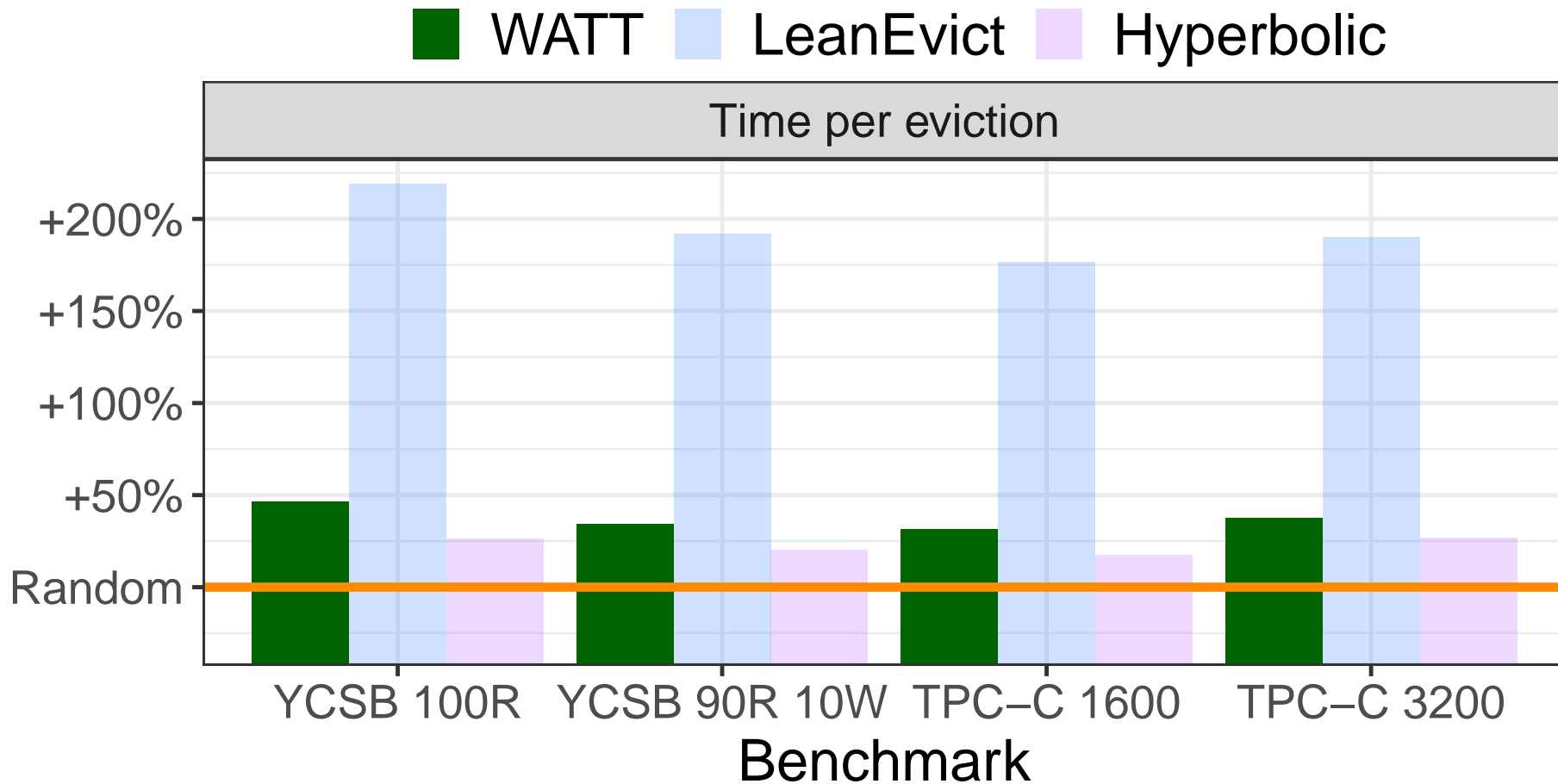
Access

■ WATT ■ LeanEvict ■ Hyperbolic

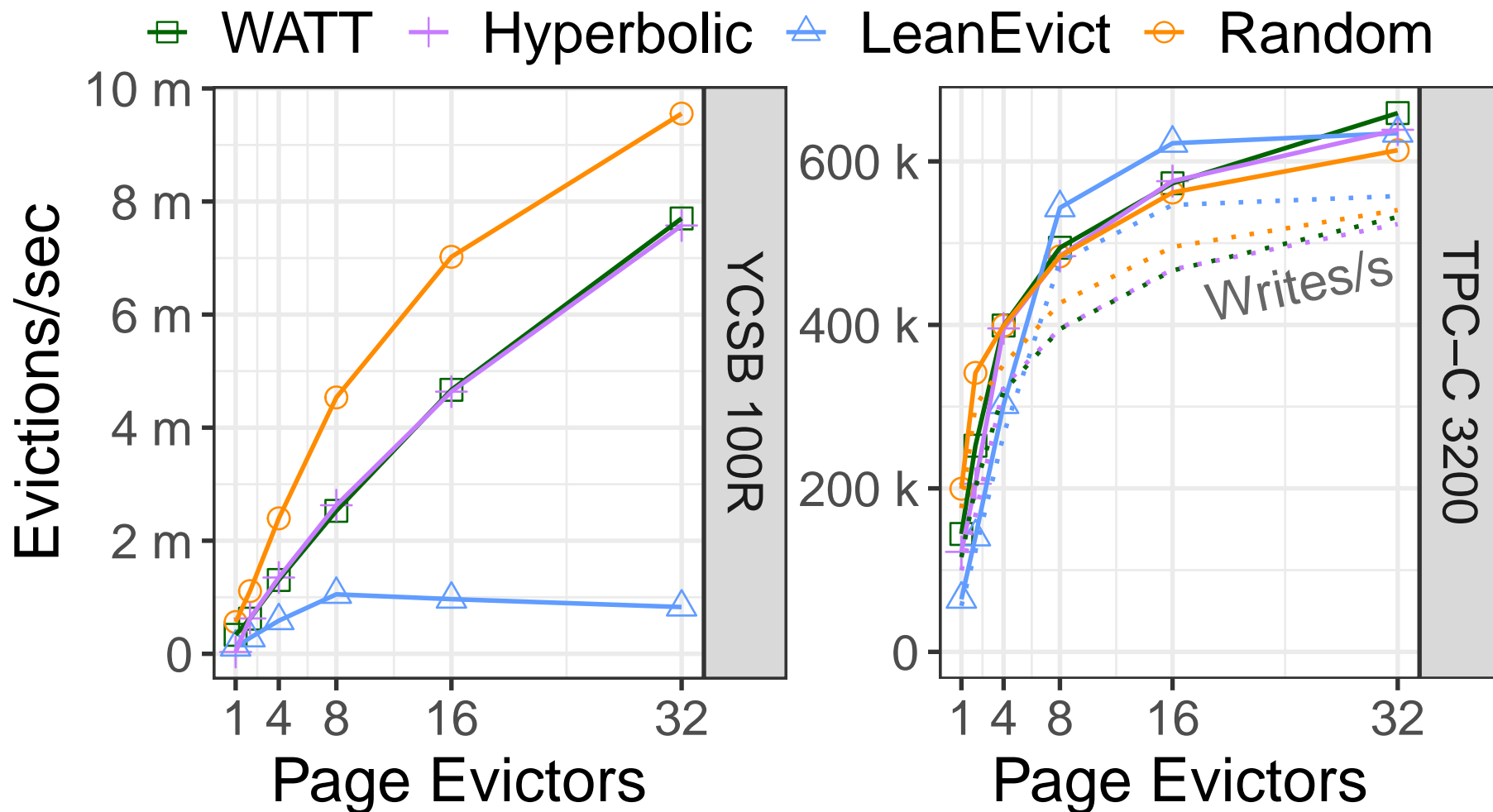


CPU Efficiency (Goal 3)

Eviction



Multi-Core Scalability (Goal 4)



4 Goals

	Replacement Effectiveness	Write Awareness	CPU Efficiency	Multi-Core Scalability
random	— —	no	++	++
CLOCK	—	no	+	+
LeanEvict	—	no	+	+
LRU2	+	no	++	+
LRU	—	no	—	—
CFLRU	—	yes	—	—
LRU_WSR	—	yes	—	—
Hyperbolic	~	no	++	++
ARC	+	no	—	—
WATT (our)	++	yes	++	++

- Write-Aware Timestamp Tracking:
Effective and Efficient Page Replacement for Modern Hardware
- 4 Goals of Page Replacement
- Paper and Artifacts: leanstore.io

Leanstore on VLDB 2023 :

D2, Tue 3:30 PM - 5:00 PM

What Modern NVMe Storage Can Do, And How To Exploit It:
High-Performance I/O for High-Performance Storage Engines
Gabriel Haas (TUM); Viktor Leis (TUM)

C8, Thu 3:30 PM - 5:00 PM

Scalable and Robust Snapshot Isolation for High-Performance Storage Engines
Adnan Alhomssi (FAU); Viktor Leis (TUM)

Write-Aware Timestamp Tracking:

Effective and Efficient Page Replacement for Modern Hardware

Demian Vöhringer¹ Viktor Leis²

¹Friedrich-Alexander-Universität Erlangen-Nürnberg

²Technische Universität München

August 29, 2023 – VLDB

More Leanstore:

D2: What Modern NVMe Storage Can Do, And How To Exploit It:
High-Performance I/O for High-Performance Storage Engines

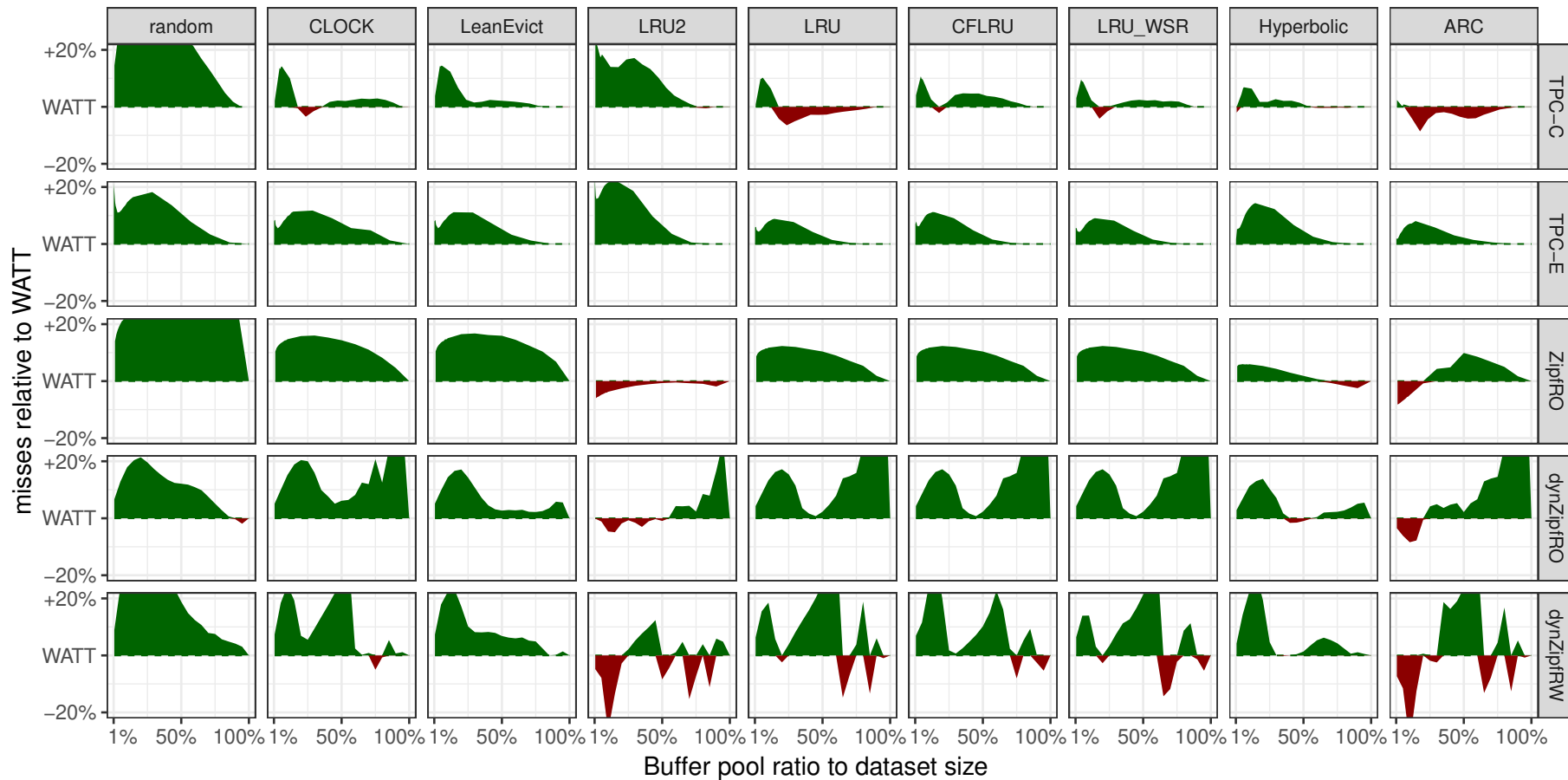
C8: Scalable and Robust Snapshot Isolation for High-Performance Storage Engines

Memory Overhead	random	CLOCK	LeanEvict	LRU2	LRU	CFLRU	LRU_WSR	Hyperbolic	ARC	WATT
Bytes*	0	$\frac{1}{8}$	8	8	24	24	$24 + \frac{1}{8}$	8	48	50
Overhead [†] [≈%]	0	0	0.2	0.2	0.6	0.6	0.6	0.2	1.2	1.2

*Using: timestamp = 4 Byte, pointer or PageID = 8 Byte, doubly-linked list = 2 pointers + 1 PageID = 24 Byte

[†]4KB pagesize

Full Comparison



```
1 atomic<uint32_t> globalTrackerTime // time (epoch)

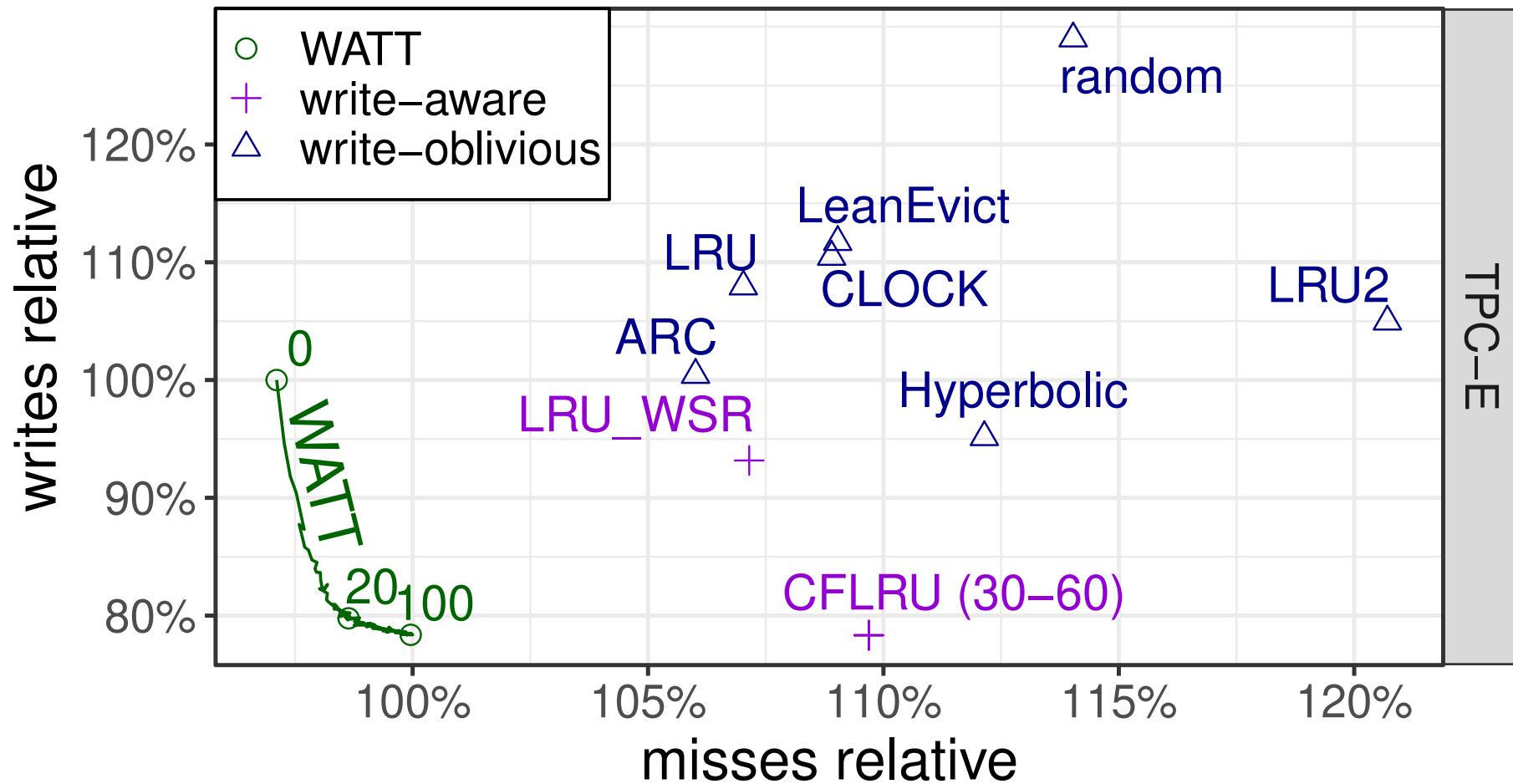
2 class PageTracker // sizeof(PageTracker) < cache_line
3     atomic<uint32_t> accessLog[8], writeLog[4]
4     atomic<uint8_t> accessHead, writeHead

5 void PageTracker::track()
6     // Compare last tracked and current epoch
7     uint8_t oldPos = accessHead.load()
8     uint32_t now = globalTrackerTime.load()
9     if (now != accessLog[oldPos])
10         // Store current epoch if they differ
11         uint8_t pos = (oldPos+1) % 8
12         accessLog[pos].store(now, memory_order_release)
13         accessHead.store(pos, memory_order_release)
```

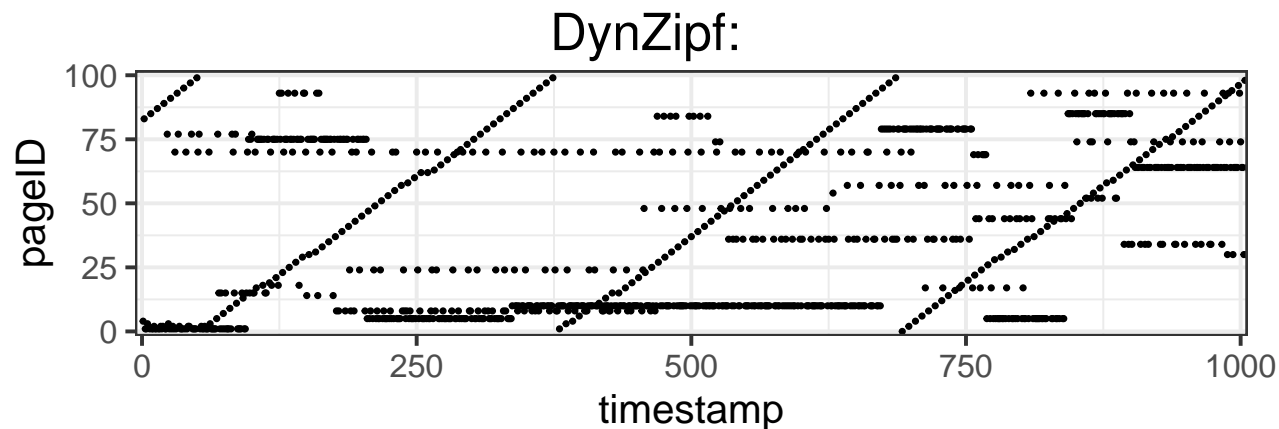
```
1 // Table of precalculated quotients i
2 // with dampening of 0.1
3 float iQuotient[] = {{0.1, 8, 7, 6, 5, 4, 3, 2},
4                      {2, 0.1, 8, 7, 6, 5, 4, 3},...}

5 float PageTracker::PVaccess()
6     uint8_t    head = accessHead.load()
7     uint32_t    now = globalTrackerTime.load()
8     __m256i     ts8 = _mm256_loadu_si256(accessLog)
9     __m256i     now8 = _mm256_set1_epi32(now)
10    __m256i     ageInt8 = _mm256_sub_epi32(now8, ts8)
11    __m256      age8 = _mm256_cvtepi32_ps(ageInt8)
12    __m256      i8 = _mm256_loadu_ps(iQuotient[head])
13    __m256      subfreq8 = _mm256_div_ps(i8, age8)
14    return      _mm256_reduce_max_ps(subfreq8)
```


Write Aware (TPC-E)



	Writes	Accesses	Pages	Hot pages [†]	Top 10 pages
TPC-C	15.6%	1 M	16,128	8	36.0%
TPC-E	5.7%	1.5 M	65,656	24	30.4%
dynZipfRO	0%	2 M	20,000	1	4.5%



^{*}Traces published: github.com/itodnerd/WATT-traces/tree/main/WATT_competition_traces

[†]more than 2% of all accesses

Benchmark	YCSB		TPC-C	
Configuration	100R	90R 10 W	1600	3200
Dataset Size	400 GB	400 GB	264 GB	588 GB