

On Scalable Integrity Checking for Secure Cloud Disks

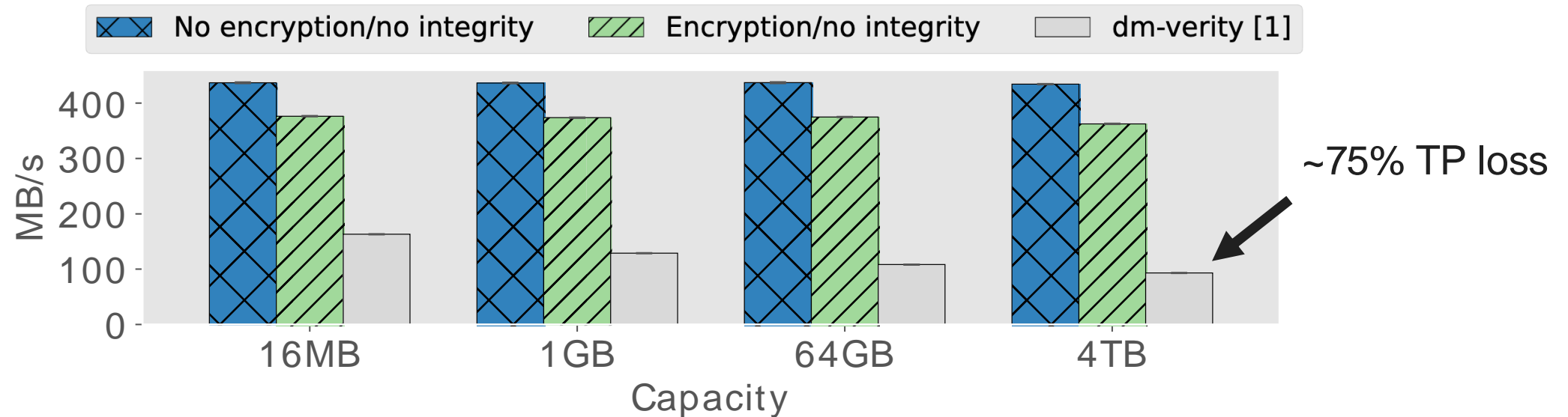
Quinn Burke, Ryan Sheatsley, Rachel King, Owen Hines,
Michael Swift, Patrick McDaniel





Integrity is costly in practice

Aggregate R/W throughput (dm-verity)



Experiment params:

Cache=10%, Read ratio=1%, I/O size=32KB,
Workload=Zipf(2.5), Threads=1, I/O
depth=32



Our work

We developed optimized integrity data structures that exploit patterns in workloads to reduce integrity costs.



Research questions:




1. What is the root cause of integrity overheads?
2. Can we model an *optimal* integrity data structure?
3. Can we realize (or approximate) the optimal in a real system?

Our work

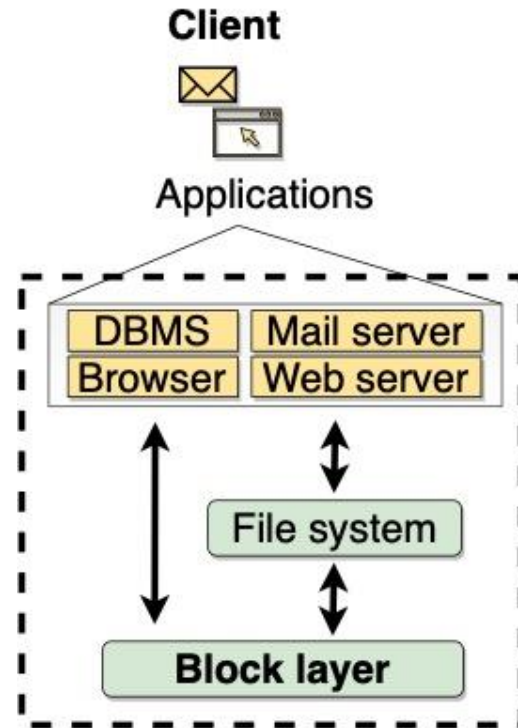
We developed optimized integrity data structures that exploit patterns in workloads to reduce integrity costs.



Research questions:

-  1. What is the root cause of integrity overheads?
 - A: Hashing (CPU costs)
-  2. Can we model an *optimal* integrity data structure?
 - A: Yes, with a priori knowledge
-  3. Can we realize (or approximate) the optimal in a real system?
 - A: Yes, by learning workload patterns

Background: Integrity checking in practice

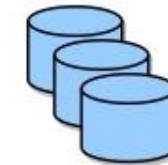


- - Trust boundary

read/write
block



**Cloud Storage
Devices**

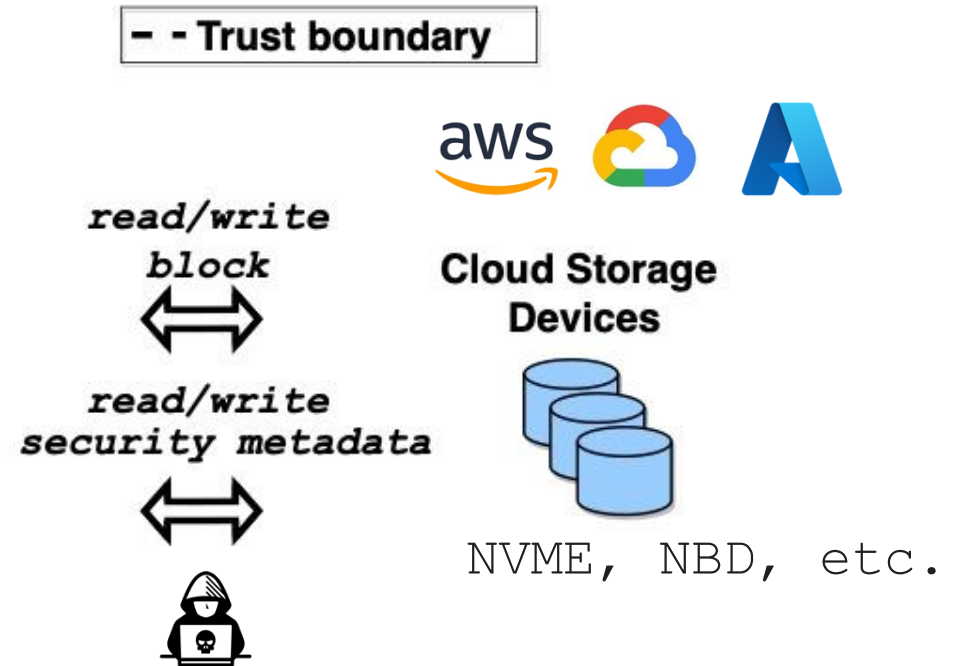
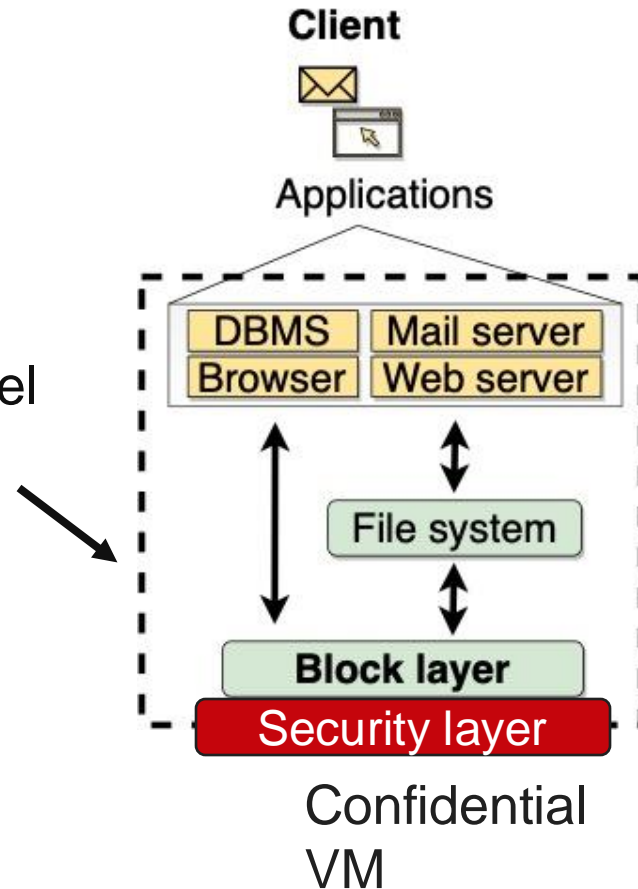


NVME, NBD, etc.



Background: Integrity checking in practice

Application isolation enforced at hardware-level (Intel TDX, ARM CCA, etc.)



Security layer

interface:

`update()` security metadata when block is written

`verify()` security metadata when block is read

Background: Security metadata



Integrity guarantees:

MACs → Authenticity

Merkle tree → Freshness

`write()`:

- Compute new block MAC, update tree with MAC, return to caller

`read()`:

- Check MAC against block, verify MAC in tree, return to caller



Background: Security metadata

Integrity guarantees:

MACs → Authenticity

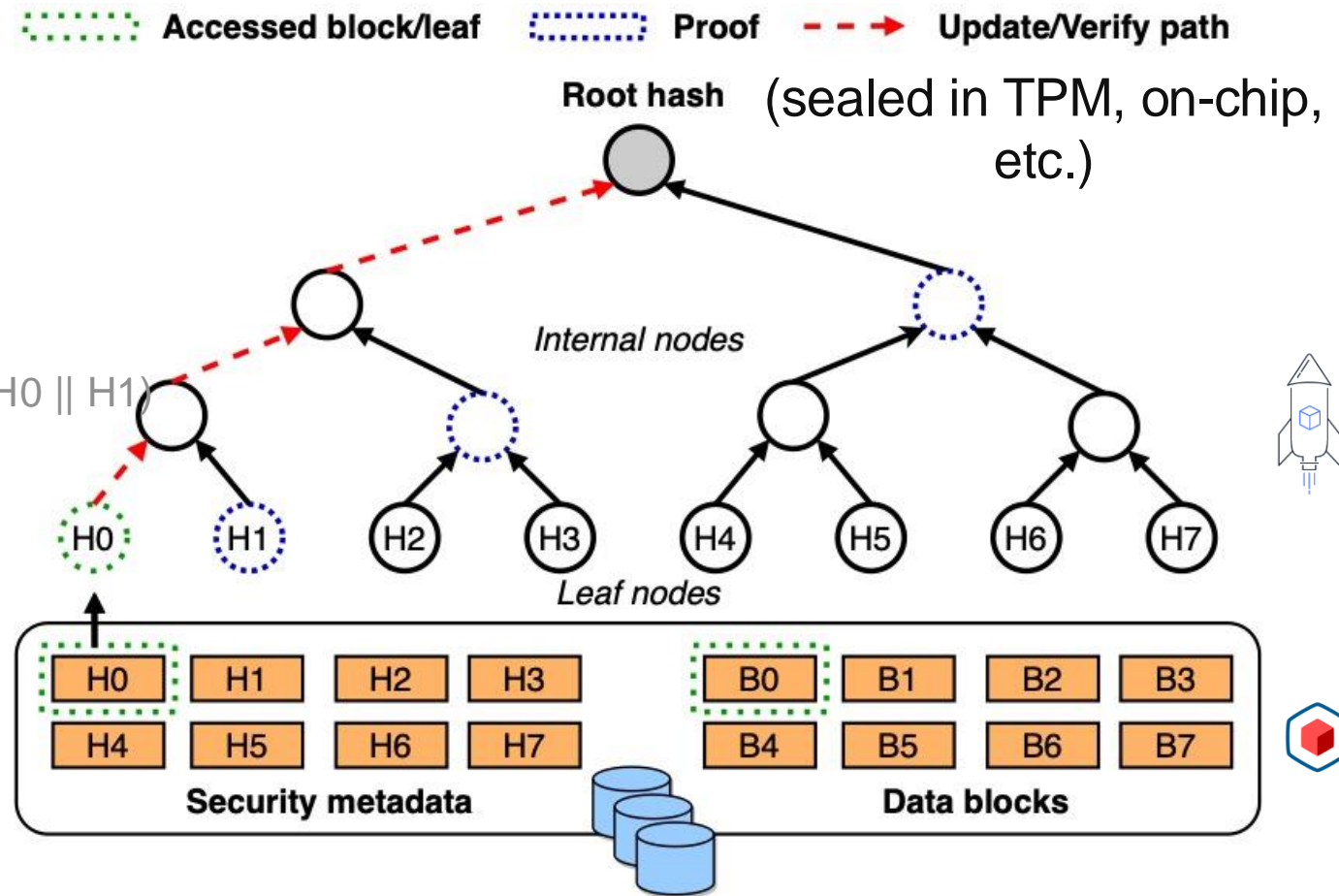
Merkle tree → Freshness

write():

- Compute new block MAC, update tree with MAC, return to caller

read():

- Check MAC against block, verify MAC in tree, return to caller





Root cause: Merkle tree traversals



Integrity guarantees:

MACs → Authenticity

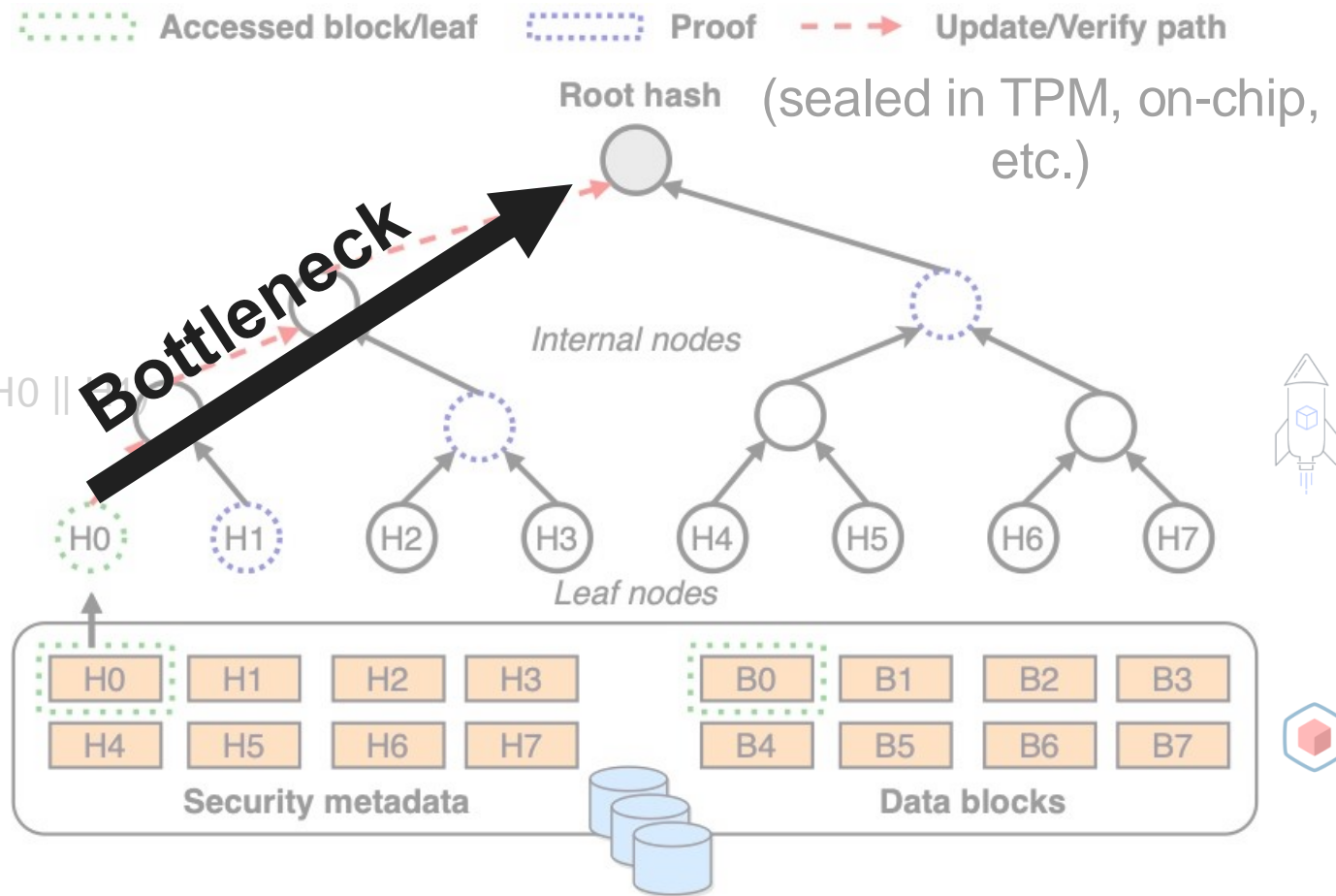
Merkle tree → Freshness

write():

- Compute new block MAC, update tree with MAC, return to caller

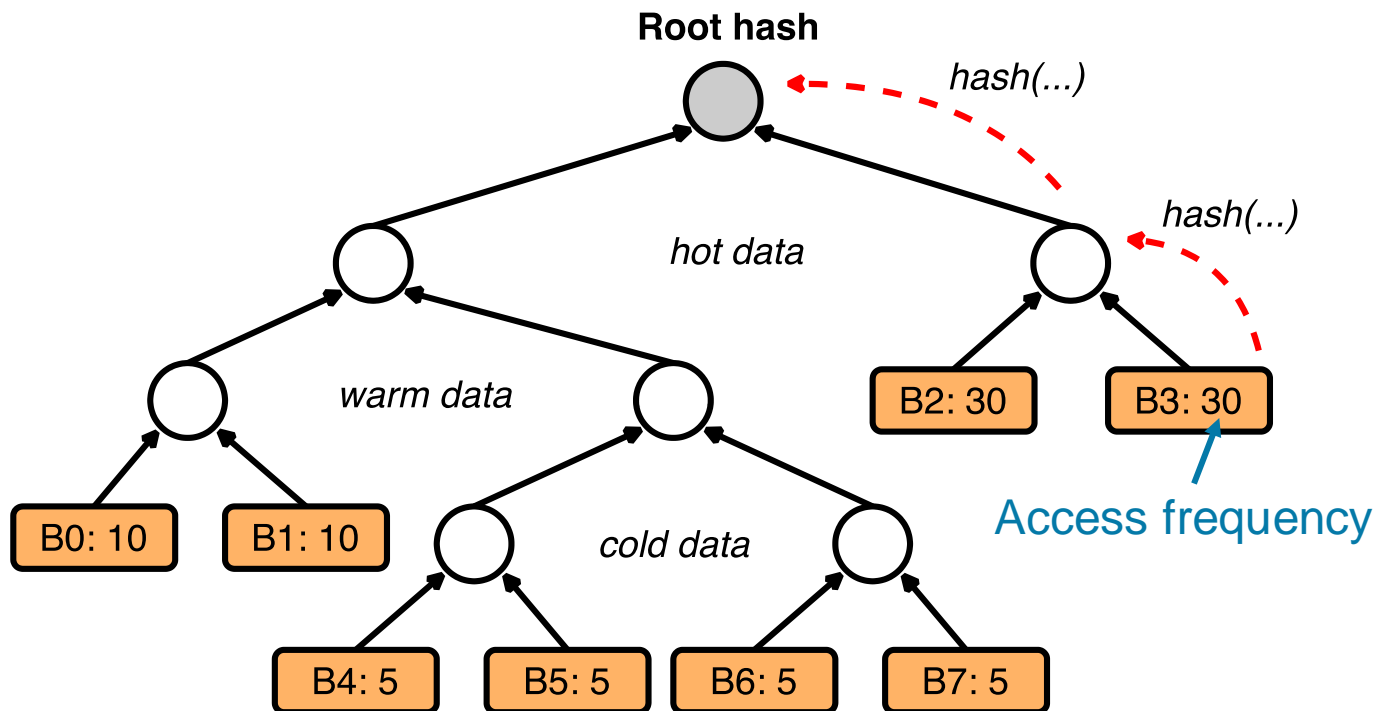
read():

- Check MAC against block, verify MAC in tree, return to caller



Opportunity: Exploit reference locality

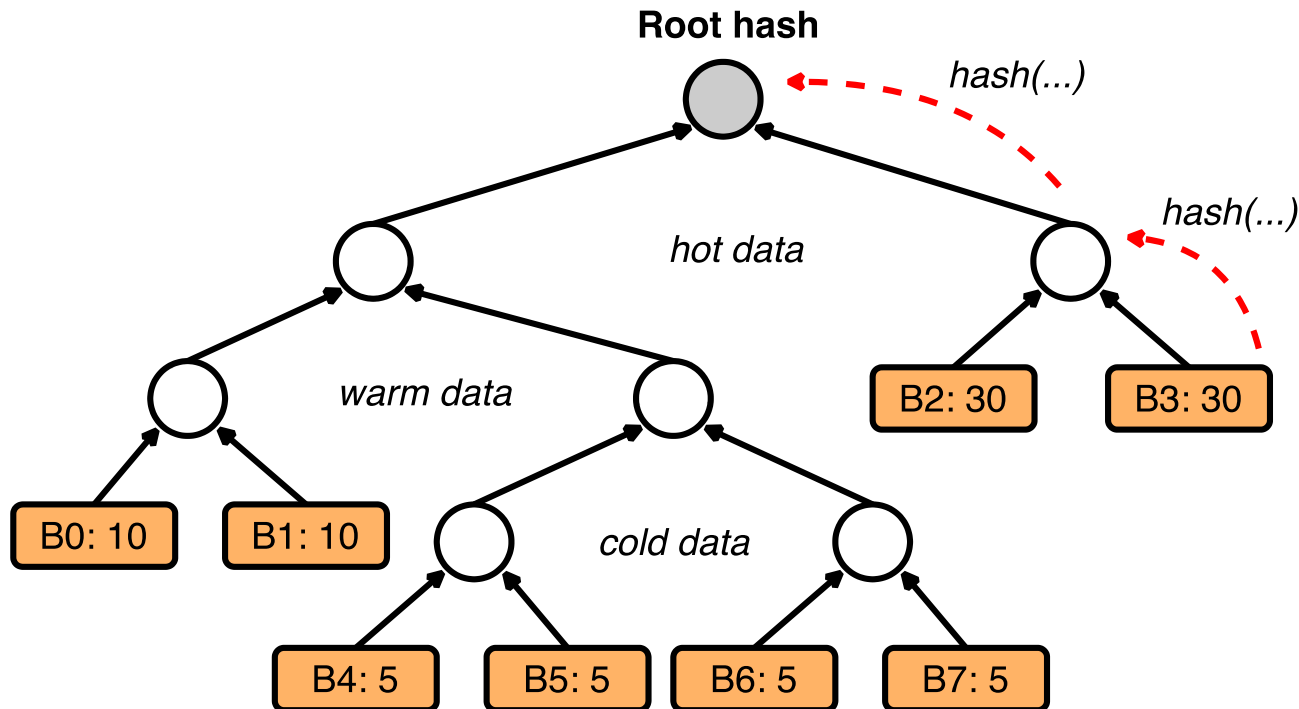
- - -> hashes computed during update/verification



- Allow the tree to become **unbalanced**
- Keep **hot leaves** (i.e., block MACs) closer to root
- Leads to lower hashing costs for hotter data (and overall, by weighted average)

Opportunity: Exploit reference locality

- - - → hashes computed during update/verification



Prior works:

1. Focused on offline integrity checks, performance not a priority, balanced trees suffice
2. Focused on slow media or small capacity



Our work:

1. Optimal tree definition
2. Optimal tree approximation (heuristic)

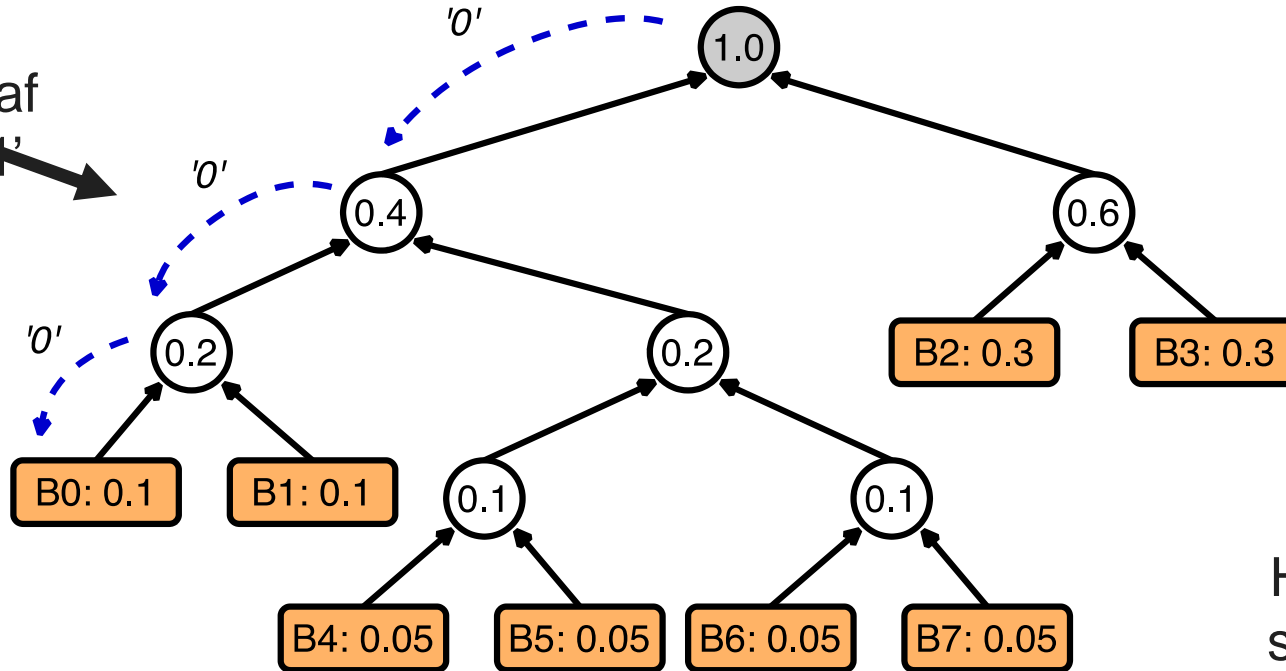


Modelling optimal performance



---> bits in codeword ($c_{B0} = '000'$)

Encoding: traverse a path to the symbol leaf and append a '0' or '1'



Formally:

$$\arg \min_C \sum_{i=1}^n w_i |c_i|, c_i \in C$$

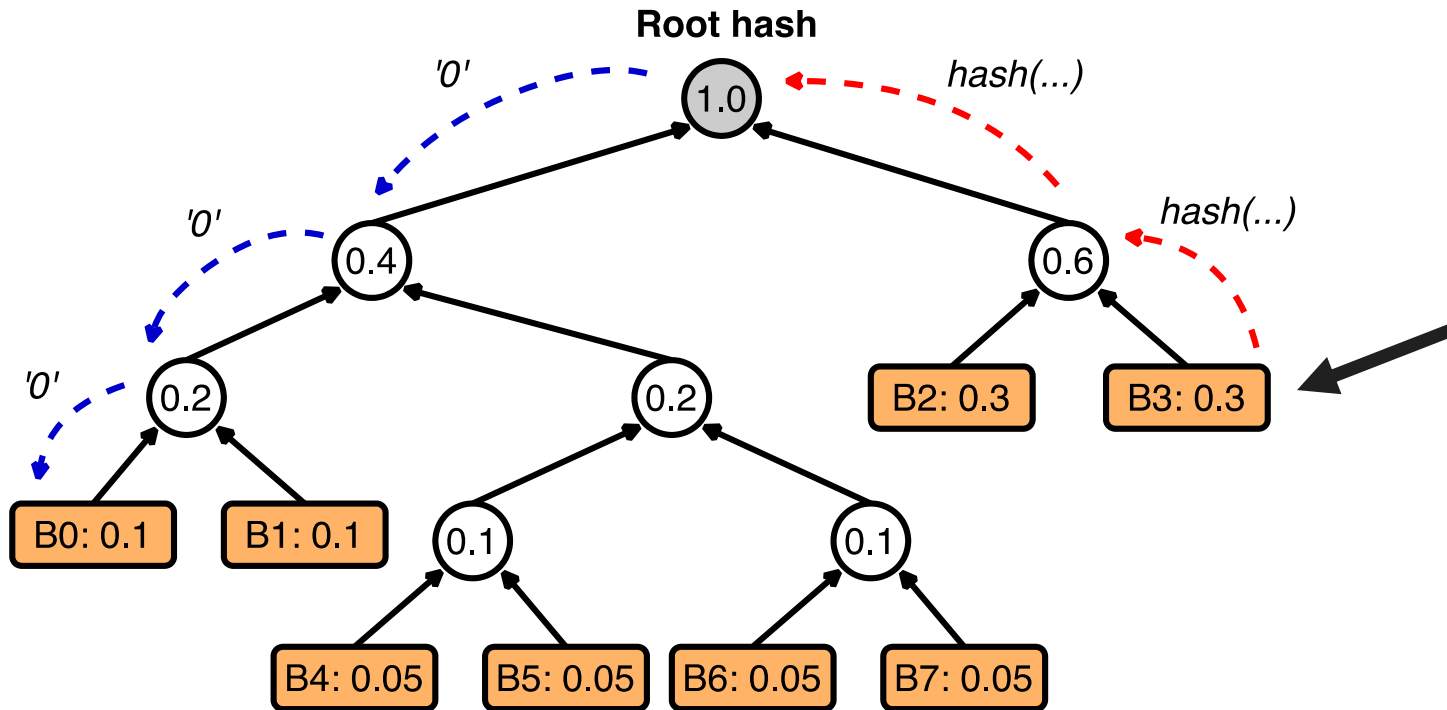


Hot symbols are assigned shorter codewords (paths)

Access probabilities								
Symbol/Block	$B0$	$B1$	$B2$	$B3$	$B4$	$B5$	$B6$	$B7$
Probability	0.1	0.1	0.3	0.3	0.05	0.05	0.05	0.05

Modelling a hash tree as a prefix tree

- - - - - ➔ hashes computed during update/verification
- - - - - ➔ bits in codeword ($c_{B0} = '000'$)



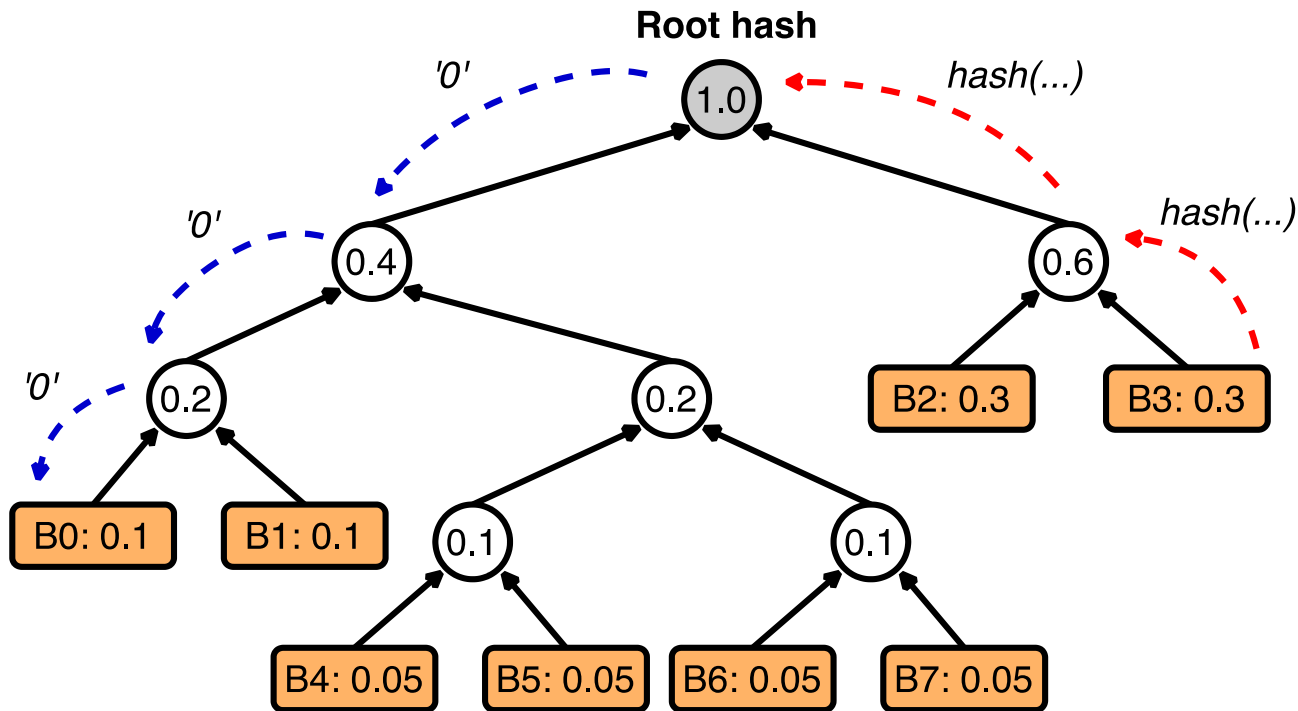
- Frequently accessed MACs (i.e., blocks) have shorter paths
- Shorter paths = lower hashing costs (for hot data and overall)

Access probabilities

Symbol/Block	$B0$	$B1$	$B2$	$B3$	$B4$	$B5$	$B6$	$B7$
Probability	0.1	0.1	0.3	0.3	0.05	0.05	0.05	0.05

Optimal trees require a priori knowledge

- - - - - ➔ hashes computed during update/verification
- - - - - ➔ bits in codeword ($c_{B0} = '000'$)



Condition:

Optimal for a **known**, fixed access probability distribution

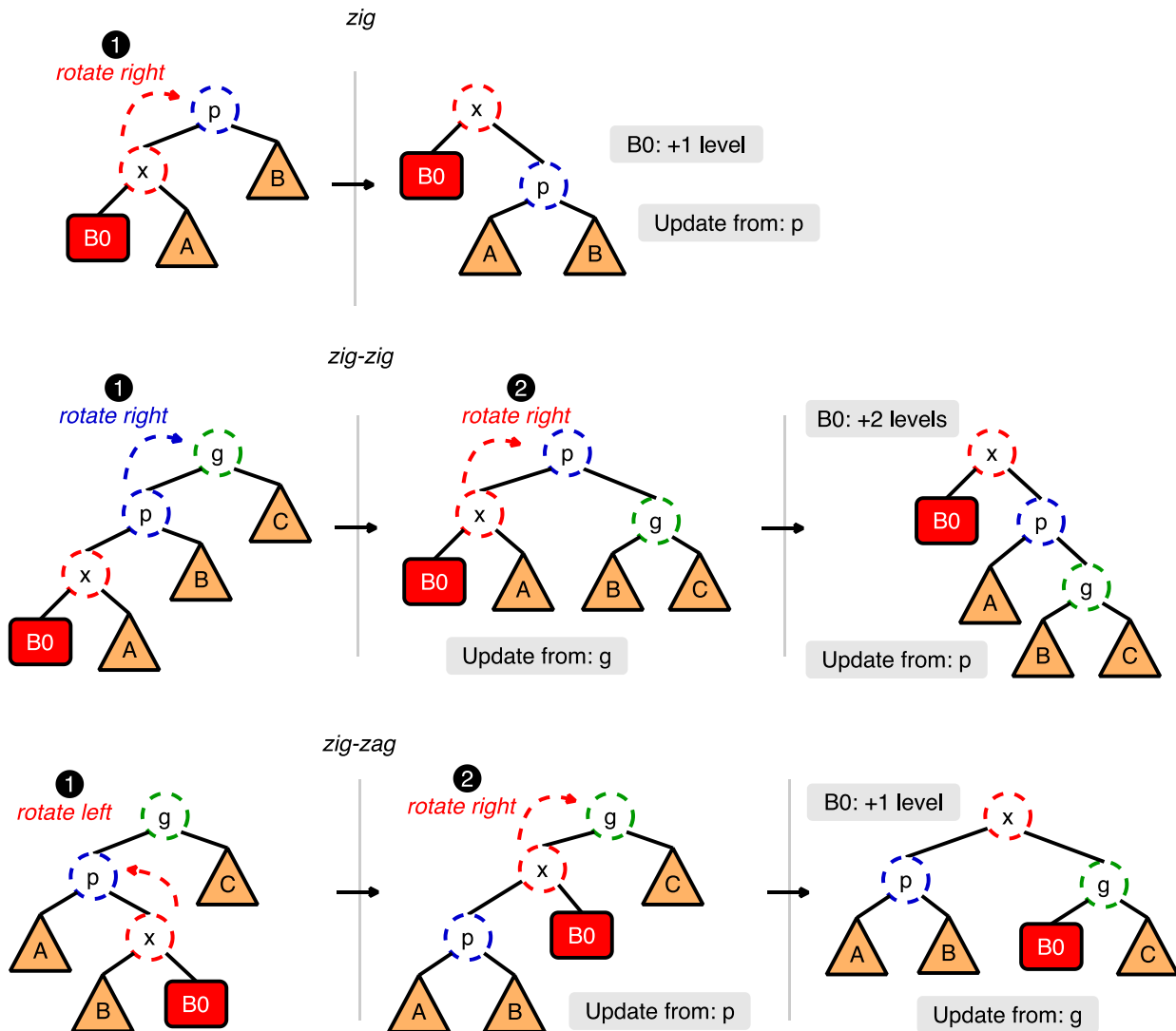
Utility:

Establishes an upper bound on performance (i.e., an optimal tree oracle)

Access probabilities

Symbol/Block	$B0$	$B1$	$B2$	$B3$	$B4$	$B5$	$B6$	$B7$
Probability	0.1	0.1	0.3	0.3	0.05	0.05	0.05	0.05

Learning workload patterns: *Dynamic Merkle Trees (DMTs)*



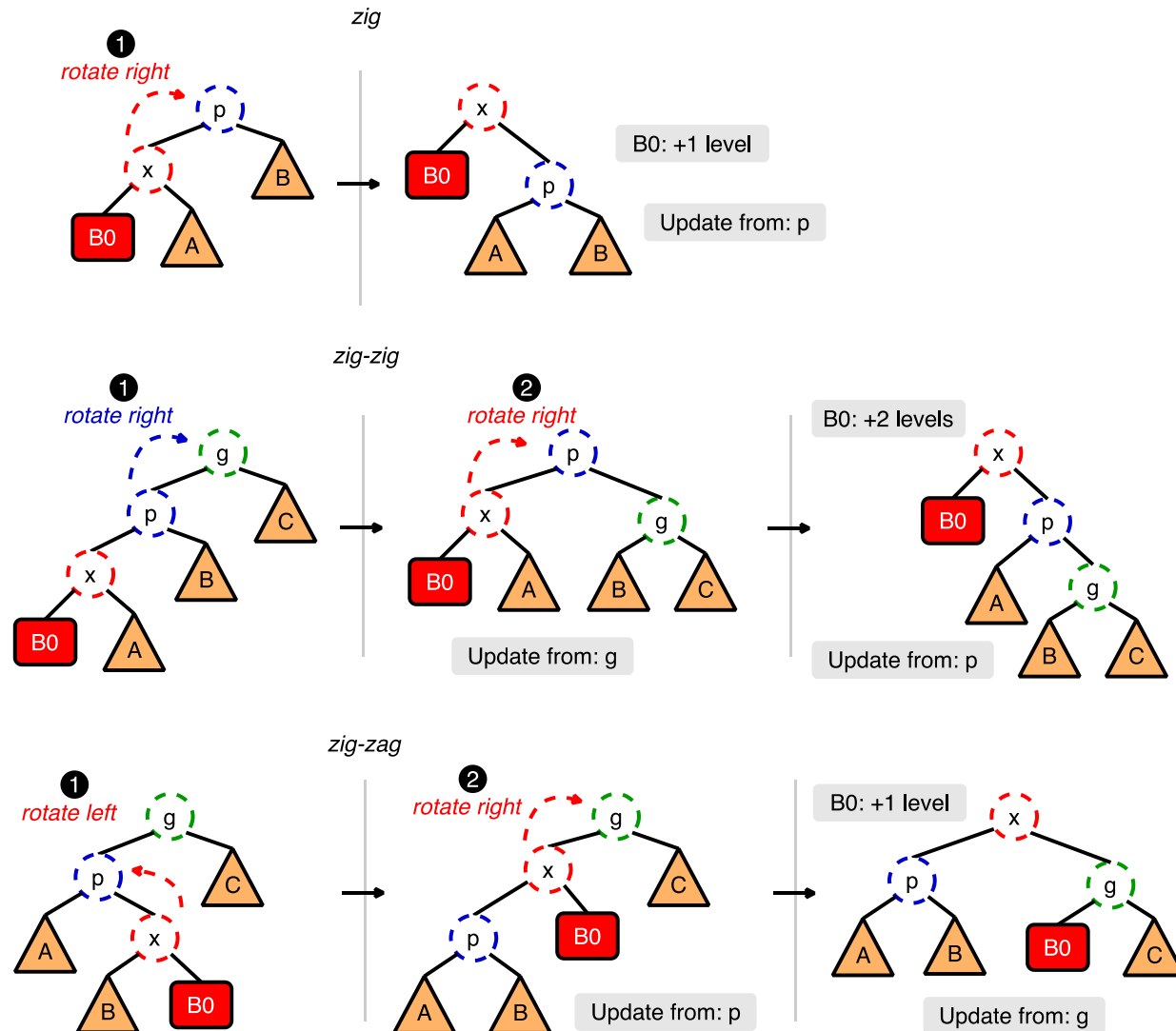
Mechanism:

- When a MAC (i.e., block) is accessed, rotate it higher in the tree

Splay trees:

- Unbalanced tree structure
- No a priori knowledge
- Start from any initial state
- Quickly adapts to new regions of interest

DMT operation



Heuristic parameters: splay probability p and splay distance d

Technical approach:

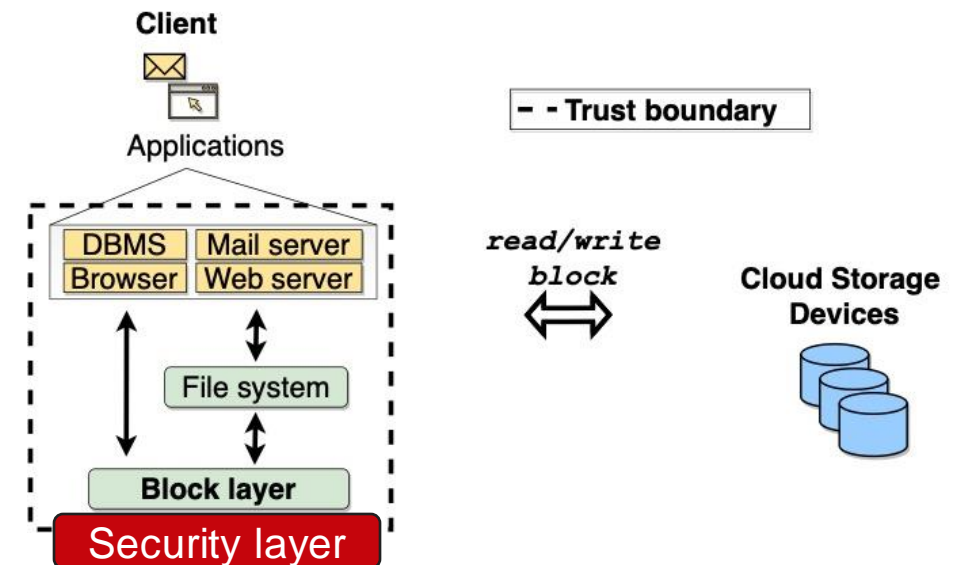
- 1 Track node hotness
- 2 Rotate nodes up the tree based on their hotness
- 3 Maintain tree invariants

Implementation & Experiment setup

- **Implementation:** As a block device driver, 5K LoC (C++) using the BDUS framework
- **Platform:** Runs on standard Linux systems
- **Setup:** AWS EC2 i4i.8xlarge instances with local NVMe disks
- **Workloads:** Fio-generated Zipfian workloads, an Alibaba trace dataset, and Filebench

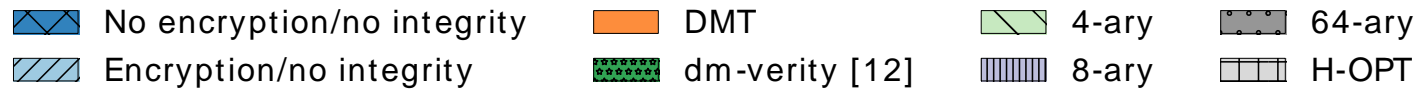
- **Evaluation questions (more in the paper):**

1. Can DMTs scale better with capacity?
2. Can DMTs adapt to workload changes?

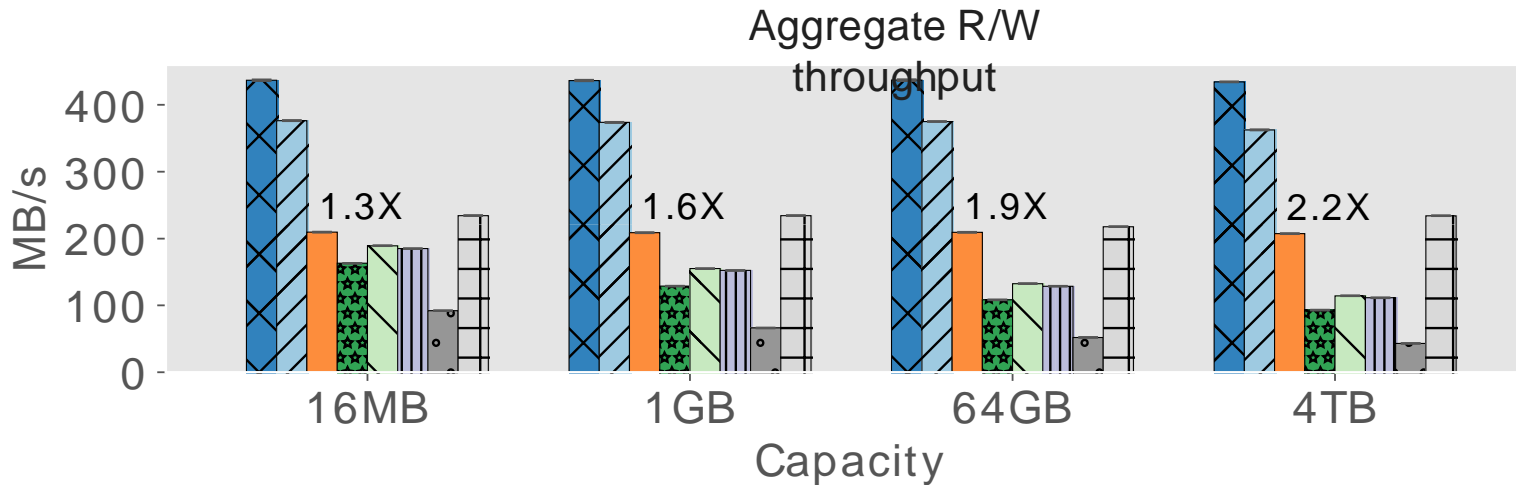




Key result: DMTs scale better with capacity

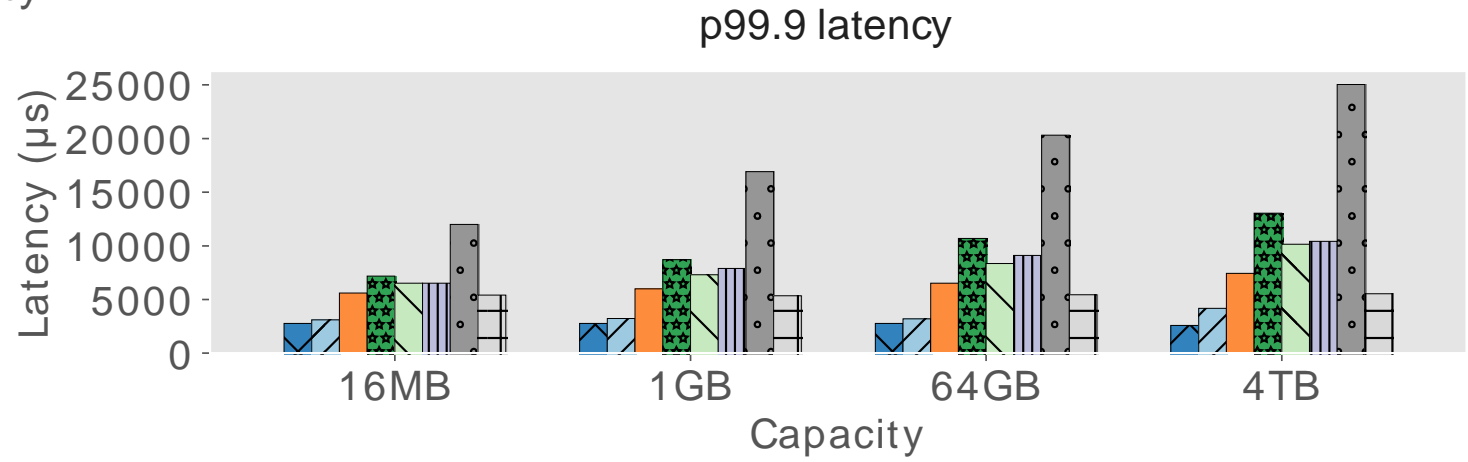


Max achievable throughput



DMTs deliver up to 2.2X speedup and consistently provide >85% of optimal (H-OPT) throughput

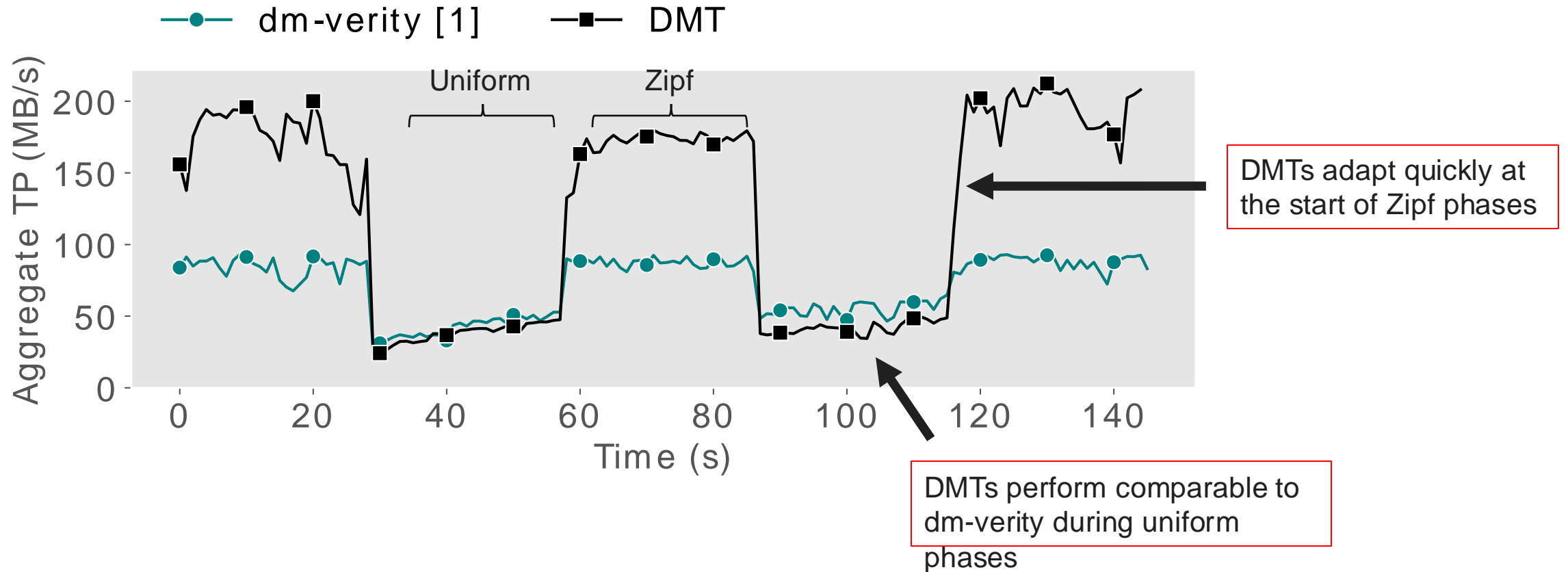
DMTs provide up to 40% lower p99.9 latency – through our design, splay costs are amortized over time



Experiment params:
Cache=10%, Read ratio=1%, I/O size=32KB,
Workload=Zipf(2.5), Threads=1, I/O depth=32



Key result: DMTs adapt quickly to workload changes





Experiment params:
Capacity=1TB, Cache=10%, Read ratio=1%, I/O size=32KB,
Workload=alternating Zipf / Uniform, Threads=1, I/O depth=32



Questions?

- Storage integrity is essential, but costly in practice
- The root cause is Merkle tree traversal (hashing) costs
- Costs can be systematically reduced by formalizing optimization objectives [1] and leveraging adaptive data structures
- DMTs can be extended to other use cases (other data structures, distributed storage deployments, etc.)
- Code: <https://github.com/MadSP-McDaniel/dmt>



 qkb@cs.wisc.edu
 <https://www.quinnburke.net>

