

Practical Hardware Transactional vEB Trees

PPoPP'24

Mohammad Khalaji¹ Trevor Brown¹ Khuzaima Daudjee¹ Vitaly Aksenov²

1



2



Introduction

Sets and Maps

DS	Look-up	Insert	Delete	Predecessor	Successor
Hash Table	$\mathcal{O}(1)$	$\mathcal{O}(1)$	$\mathcal{O}(1)$	$\mathcal{O}(1)$	$\mathcal{O}(1)$
BST	$\mathcal{O}(\lg n)$	$\mathcal{O}(\lg n)$	$\mathcal{O}(\lg n)$	$\mathcal{O}(\lg n)$	$\mathcal{O}(\lg n)$
B-Tree	$\mathcal{O}(\log_4 n)$	$\mathcal{O}(\log_4 n)$	$\mathcal{O}(\log_4 n)$	$\mathcal{O}(\log_4 n)$	$\mathcal{O}(\log_4 n)$

Sets and Maps

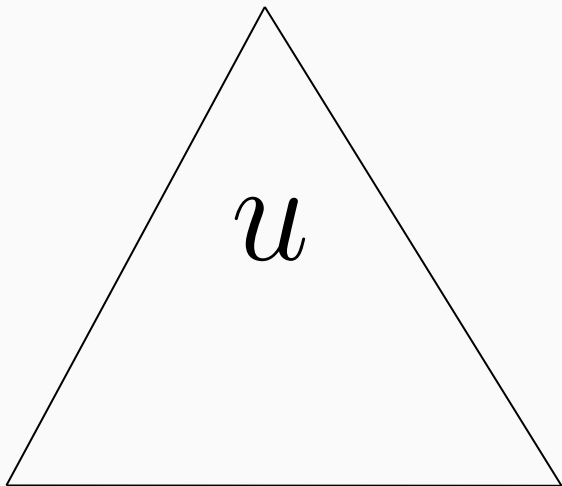
DS	Look-up	Insert	Delete	Predecessor	Successor
Hash Table	$\mathcal{O}(1)$	$\mathcal{O}(1)$	$\mathcal{O}(1)$	$\mathcal{O}(1)$	$\mathcal{O}(1)$
BST	$\mathcal{O}(\lg n)$	$\mathcal{O}(\lg n)$	$\mathcal{O}(\lg n)$	$\mathcal{O}(\lg n)$	$\mathcal{O}(\lg n)$
B-Tree	$\mathcal{O}(\log_4 n)$	$\mathcal{O}(\log_4 n)$	$\mathcal{O}(\log_4 n)$	$\mathcal{O}(\log_4 n)$	$\mathcal{O}(\log_4 n)$
vEB Tree	$\mathcal{O}(\lg \lg n)$	$\mathcal{O}(\lg \lg n)$	$\mathcal{O}(\lg \lg n)$	$\mathcal{O}(\lg \lg n)$	$\mathcal{O}(\lg \lg n)$

Sets and Maps

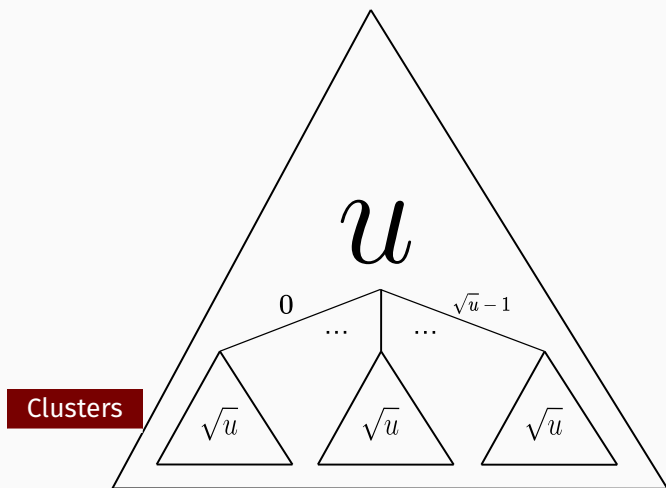
DS	Look-up	Insert	Delete	Predecessor	Successor
Hash Table	$\mathcal{O}(1)$	$\mathcal{O}(1)$	$\mathcal{O}(1)$	$\mathcal{O}(1)$	$\mathcal{O}(1)$
BST	$\mathcal{O}(\lg n)$	$\mathcal{O}(\lg n)$	$\mathcal{O}(\lg n)$	$\mathcal{O}(\lg n)$	$\mathcal{O}(\lg n)$
B-Tree	$\mathcal{O}(\log_4 n)$	$\mathcal{O}(\log_4 n)$	$\mathcal{O}(\log_4 n)$	$\mathcal{O}(\log_4 n)$	$\mathcal{O}(\log_4 n)$
vEB Tree	$\mathcal{O}(\lg \lg n)$	$\mathcal{O}(\lg \lg n)$	$\mathcal{O}(\lg \lg n)$	$\mathcal{O}(\lg \lg n)$	$\mathcal{O}(\lg \lg n)$

$\lg \lg$	
65K	4
1M	4.3
1B	4.9
1T	5.3

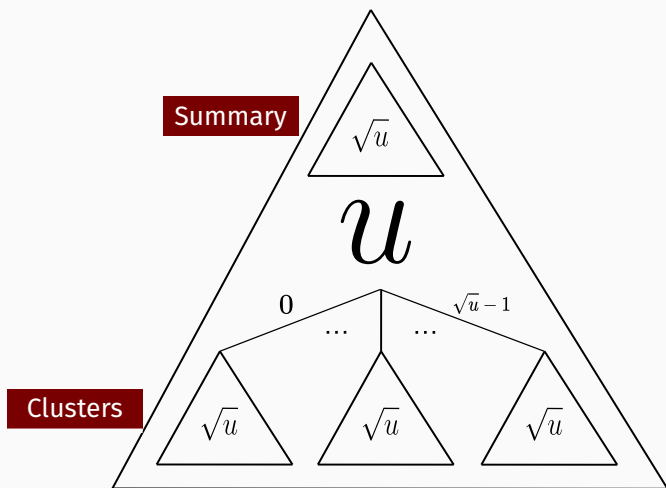
- Introduced in 1977 by Peter van Emde Boas [4]
- Fixed universe size $= 2^{lg}$: $f0;:::; 1g \rightarrow lg$ bit keys
- $\mathcal{O}(lg lg)$ complexity for all operations



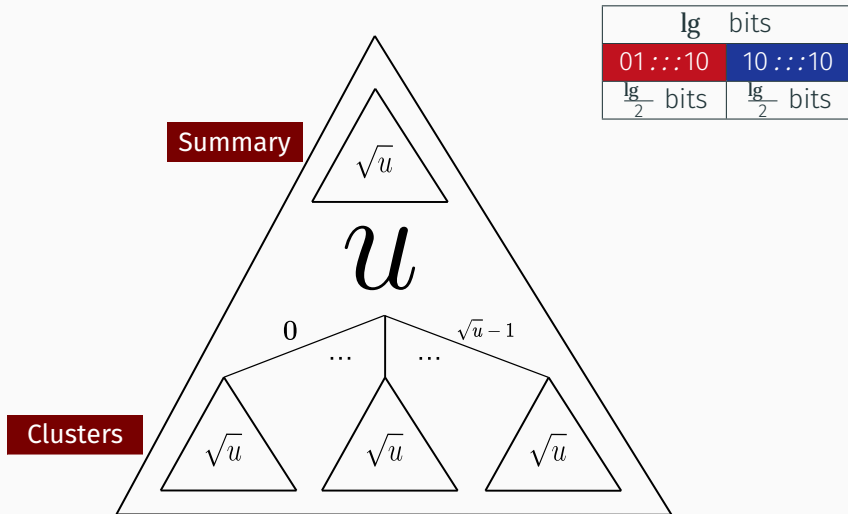
vEB Tree Structure



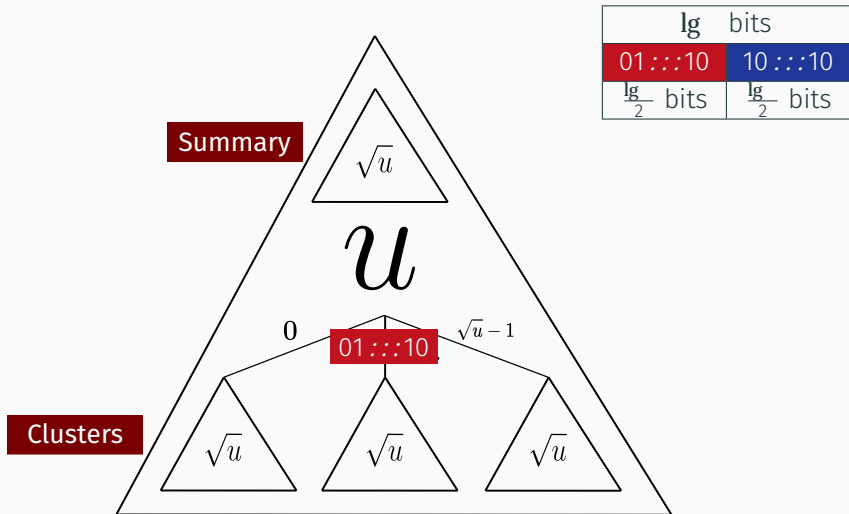
vEB Tree Structure



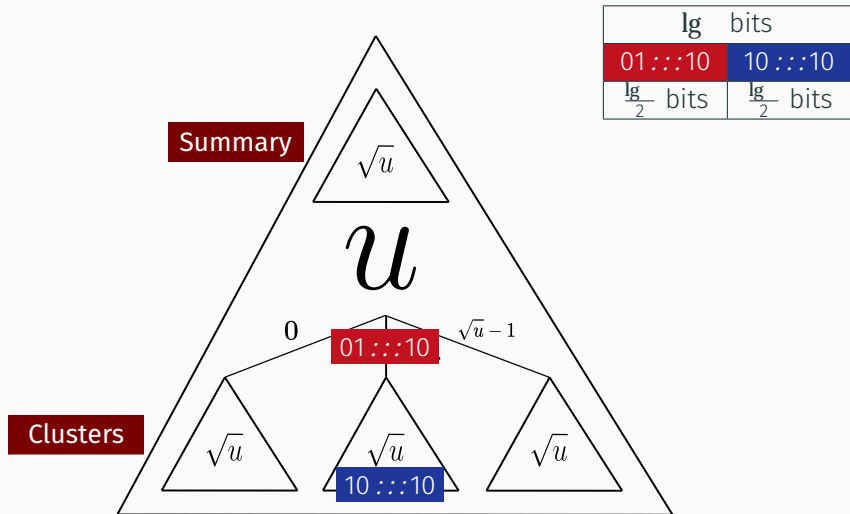
vEB Insert Example



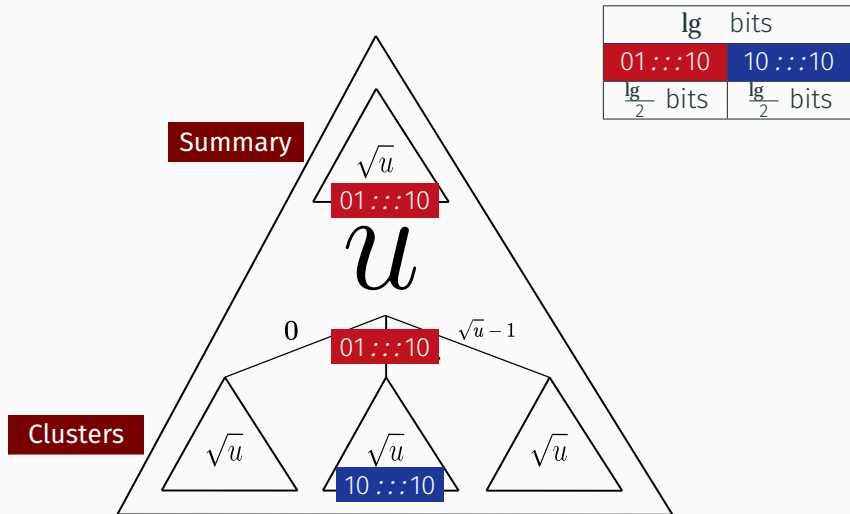
vEB Insert Example



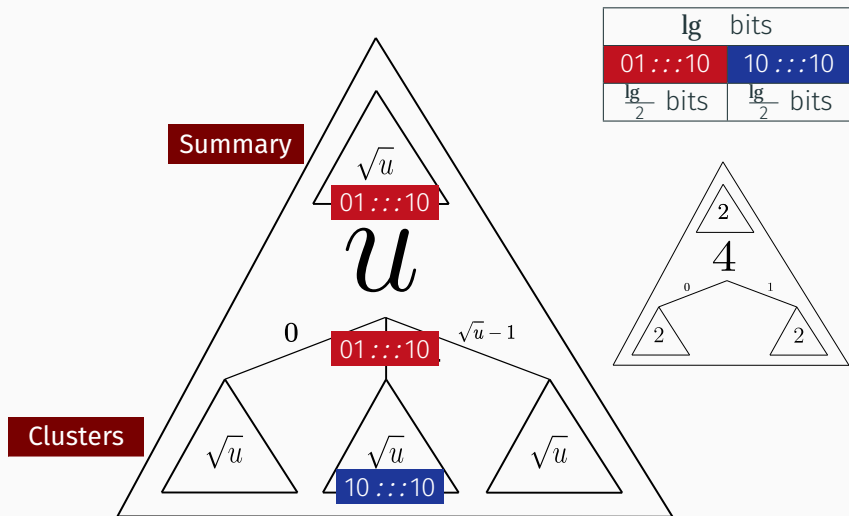
vEB Insert Example



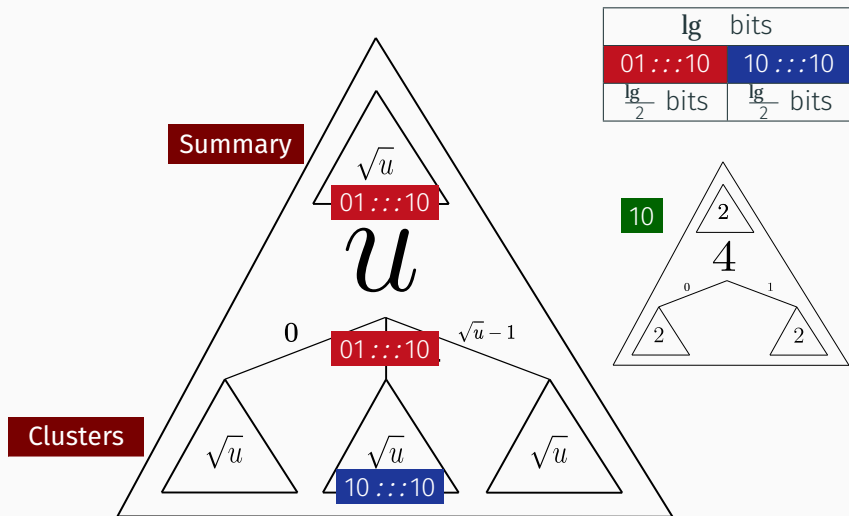
vEB Insert Example



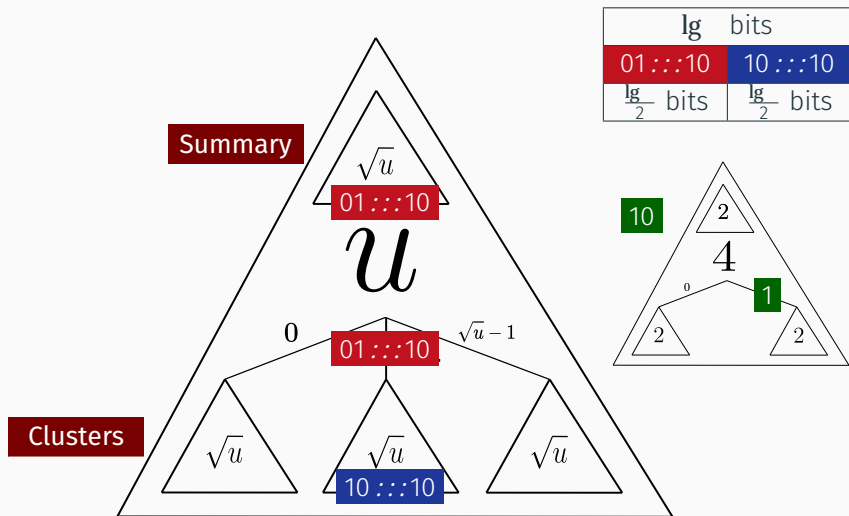
vEB Insert Example



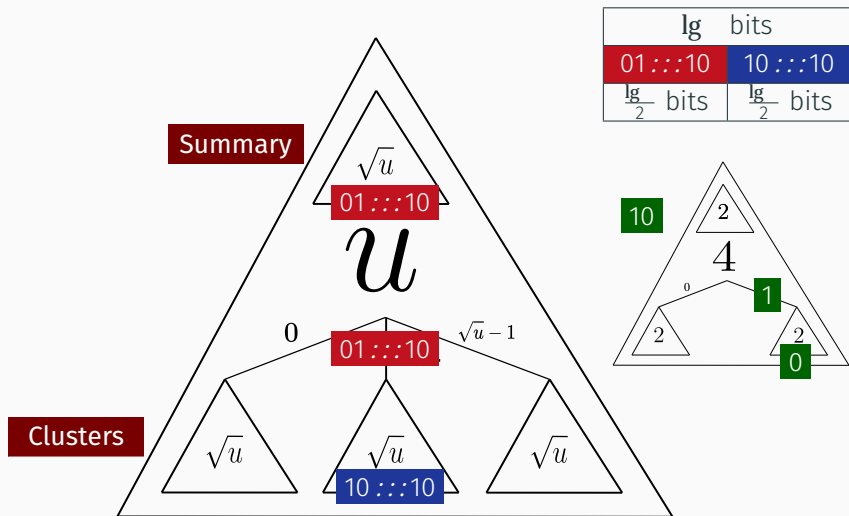
vEB Insert Example



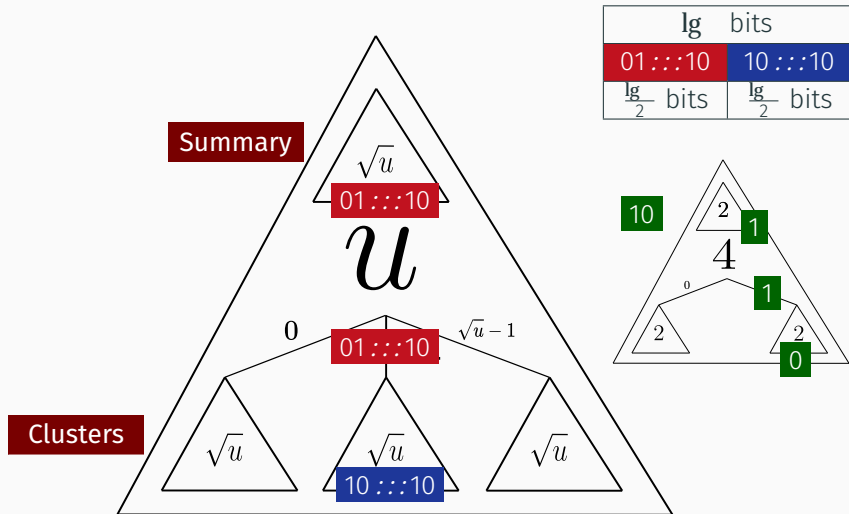
vEB Insert Example



vEB Insert Example



vEB Insert Example



Hardware Transactional Memory (HTM)

- Atomic blocks of code
 - Take effect as a single indivisible unit (Commit)
 - Or rolled back with no effect on shared memory (Abort)
- Hardware level → Efficient

Hardware Transactional Memory (HTM)

- Atomic blocks of code
 - Take effect as a single indivisible unit (Commit)
 - Or rolled back with no effect on shared memory (Abort)
- Hardware level → Efficient

ag fgTghf 0 RkUXZ\afiff.

Hardware Transactional Memory (HTM)

- Atomic blocks of code
 - Take effect as a single indivisible unit (Commit)
 - Or rolled back with no effect on shared memory (Abort)
- Hardware level → Efficient

`\ag fgTghf 0 RkUXZ\afifl.`

`\Y fifgTghf 00 RK58: <ARFG4EG87fl n`

`GeTafTVg\baT_ VbW! !!`

`RkXaWifl.`

p

Hardware Transactional Memory (HTM)

- Atomic blocks of code
 - Take effect as a single indivisible unit (Commit)
 - Or rolled back with no effect on shared memory (Abort)
- Hardware level → Efficient

$\backslash ag \text{ fgTghf } 0 \text{ RkUXZ} \backslash a \text{ fiff} .$

$\backslash Y \text{ fiffgTghf } 00 \text{ RK58: } < \text{ARFG4EG87fl } n$

$\text{GeTafTVg} \backslash \text{baT}_ \text{ VbW\!} !!$

$\text{RkXaWiff} .$

p

$\text{X}_f \text{X } n$

$9T_ \text{UTV}^ \wedge \text{ VbW\!} !!$

p

Transactional Lock Elision (TLE)

Transactional Lock Elision (TLE)

\Y firKUXZ\afifl 00 RK58: <ARFG4EG87fl n

Fast Path

GeTafTVg\ba UbW!!!

RkXaWifl.

p

Transactional Lock Elision (TLE)

\Y fiRkUXZ\afifl 00 RK58: <ARFG4EG87fl n

Fast Path

GeTafTVg\ba UbW!!!

RkXaWifl.

p

X_fX n

_bV^fifl.

GeTafTVg\ba UbW!!!

ha_bV^fifl.

p

Fallback Path

Transactional Lock Elision (TLE)

$\backslash Y$ $f_i R_k U X Z \backslash a f i f l$ $00 R K 58: < A R F G 4 E G 8 7 f l$ n
 $\backslash Y$ $_ b V ^ \wedge \backslash f$ $T _ e X T W [X _ W$
 $T U b e g f i f l .$

Fast Path

$G e T a f T V g \backslash b a$ $U b W ! ! !$

$R k X a W i f l .$

p

$X _ f X n$

$_ b V ^ \wedge f i f l .$

$G e T a f T V g \backslash b a$ $U b W ! ! !$

$h a _ b V ^ \wedge f i f l .$

p

Fallback Path

Transactional Lock Elision (TLE)

$\backslash Y$ $f_i R k U X Z \backslash a f i f l$ $00 R K 58: < A R F G 4 E G 8 7 f l$ n
 $\backslash Y$ $_ b V ^ \wedge \backslash f$ $T _ e X T W [X _ W$
 $T U b e g f i f l .$

Fast Path

$G e T a f T V g \backslash b a$ $U b W ! ! !$ G_2

$R k X a W i f l .$

P

$X_f X n$

$_ b V ^ \wedge f i f l .$ G

$G e T a f T V g \backslash b a$ $U b W ! ! !$

$h a _ b V ^ \wedge f i f l .$

P

Fallback Path

Transactional Lock Elision (TLE)

$\backslash Y$ *fiRkUXZ\afifl 00 RK58: <ARFG4EG87fl n*
 $\backslash Y$ *_bV^ \f T_eXTW [X_W*
TUbegfifl.

Fast Path

GeTafTVg\ba UbW!!! G_2

RkXaWifl.

P

X_fX n

_bV^fifl.

Fallback Path

GeTafTVg\ba UbW!!! G_1

ha_bV^fifl.

P

Transactional Lock Elision (TLE)

\Y fiRkUXZ\afifl 00 RK58: <ARFG4EG87fl n
_bV^ \f T_eXTW [X_W
TUbegfifl.

Fast Path

GeTafTVg\ba UbW!!!

RkXaWifl. G_2

P
X_fX n

_bV^fifl.

Fallback Path

GeTafTVg\ba UbW!!! G_1

ha_bV^fifl.

P

Transactional Lock Elision (TLE)

$\backslash Y$ $f_i R_k U X Z \backslash a f i f l$ $00 R K 58: < A R F G 4 E G 8 7 f l$ n
 $\backslash Y$ $_ b V ^ \wedge \backslash f$ $T_e X T W$ $[X _ W$
 $T U b e g f i f l .$

Fast Path

$G e T a f T V g \backslash b a$ $U b W ! ! !$

$R k X a W i f l .$

P

$X_f X n$ G_2

$_ b V ^ \wedge f i f l .$

Fallback Path

$G e T a f T V g \backslash b a$ $U b W ! ! !$ G_1

$h a _ b V ^ \wedge f i f l .$

P

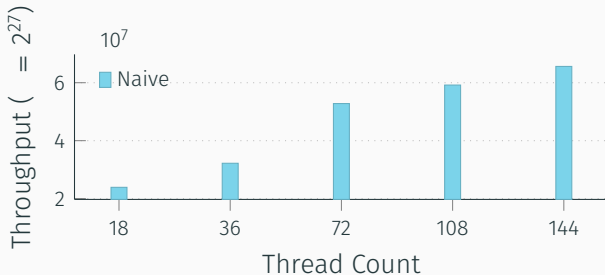
Methodology

- Concurrent vEB trees with recursive summaries
- Full support for successor and predecessor queries
- Starting with a naive vEB set implementation
- Use simple building blocks to end up with a performant vEB map

First Step: Naive

Idea

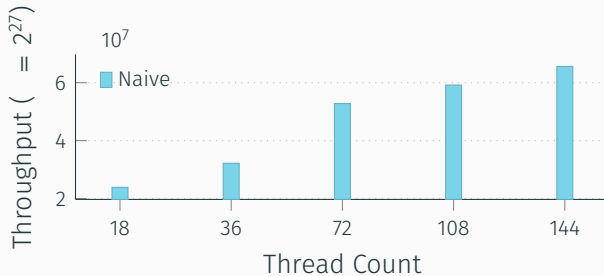
Wrap sequential vEB operations in TLE blocks



First Step: Naive

Idea

Wrap sequential vEB operations in TLE blocks



Problem

Memory usage does not depend on the number of keys!

Improve Memory Footprint: Dynamic

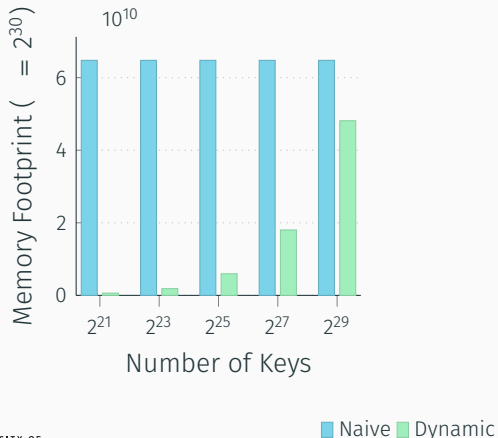
Can we save memory?

In the Naive variation, the entire tree is pre-allocated

Improve Memory Footprint: Dynamic

Can we save memory?

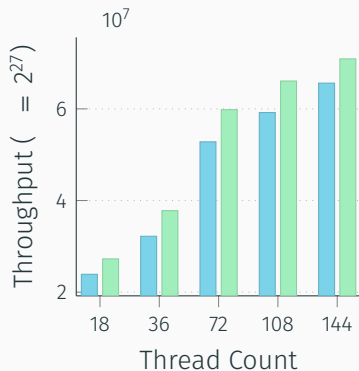
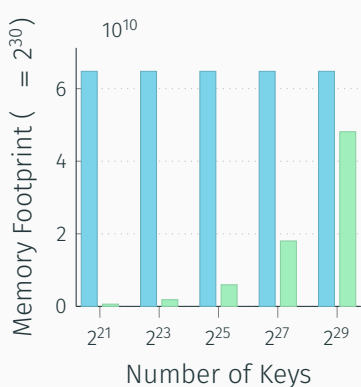
In the Naive variation, the entire tree is pre-allocated



Improve Memory Footprint: Dynamic

Can we save memory?

In the Naive variation, the entire tree is pre-allocated



■ Naive ■ Dynamic

Room For Improvement: Cut-Off



Data

2^{27}

Room For Improvement: Cut-Off



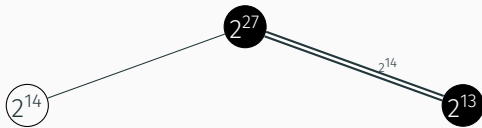
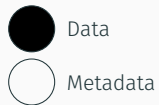
Data



Metadata



Room For Improvement: Cut-Off



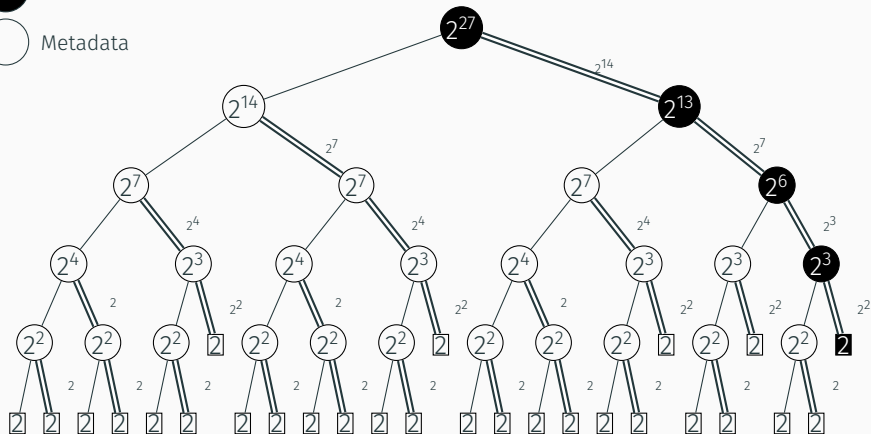
Room For Improvement: Cut-Off



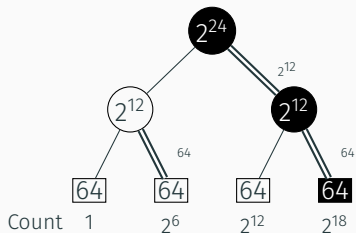
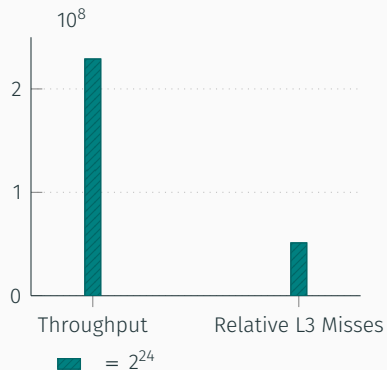
Data



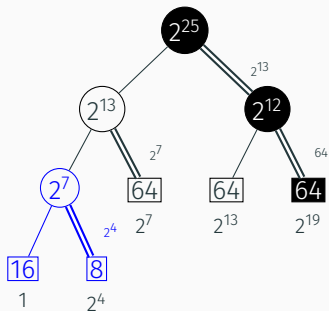
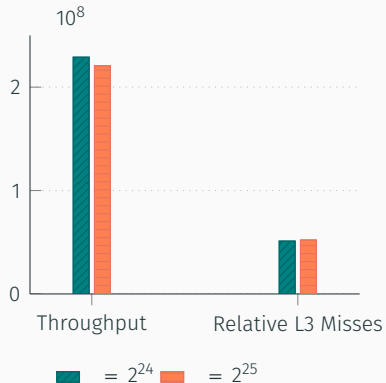
Metadata



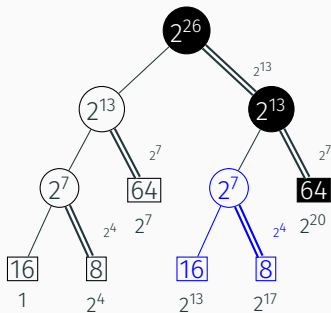
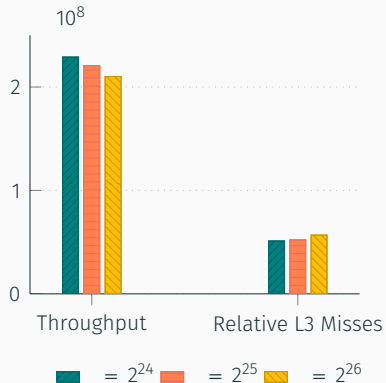
Problem with Cut-Off



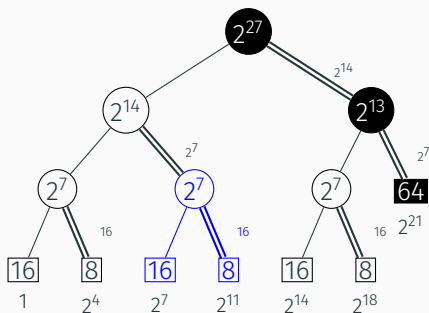
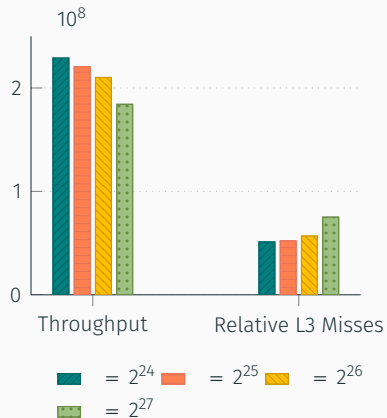
Problem with Cut-Off



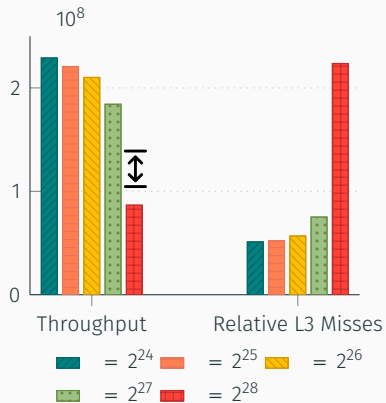
Problem with Cut-Off



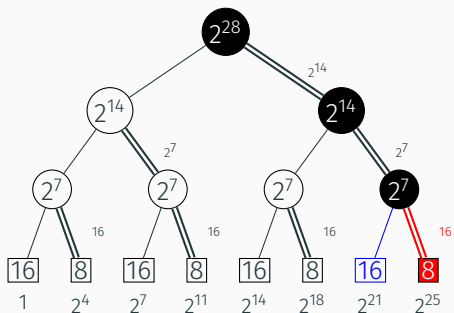
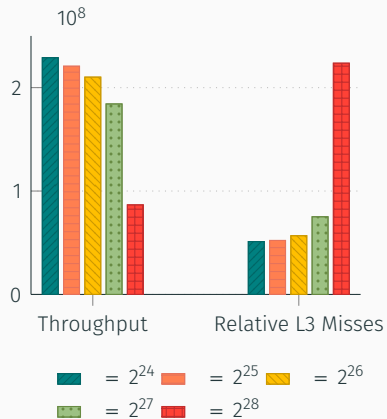
Problem with Cut-Off



Problem with Cut-Off

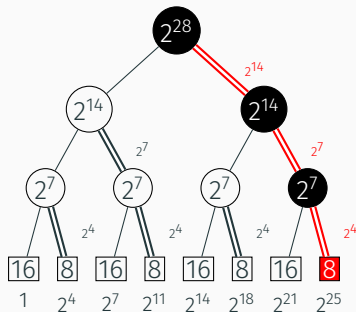


Problem with Cut-Off



Solution: New-Root

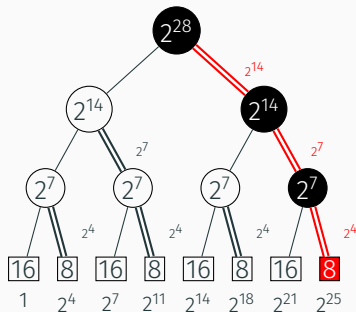
Before



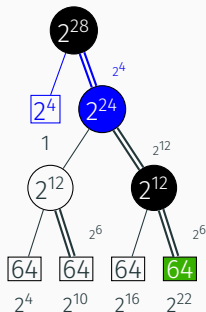
Universe Size	2^{24}	2^{25}	2^{26}	2^{27}	2^{28}	2^{29}
Data Leaf Count	2^{18}	2^{19}	2^{20}	2^{21}	2^{25}	2^{26}
Data Leaf Size	64	64	64	64	8	8

Solution: New-Root

Before

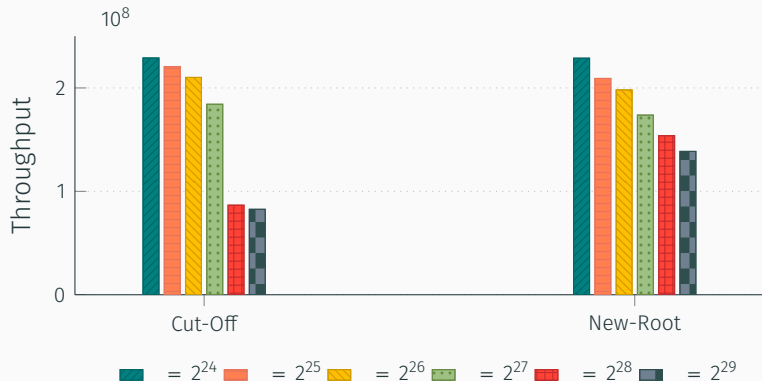


After

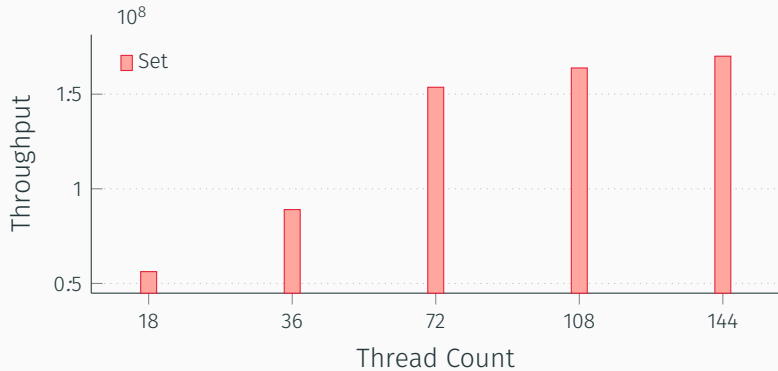


Universe Size	2^{24}	2^{25}	2^{26}	2^{27}	2^{28}	2^{29}
Data Leaf Count	2^{18}	2^{19}	2^{20}	2^{21}	2^{22}	2^{23}
Data Leaf Size	64	64	64	64	64	64

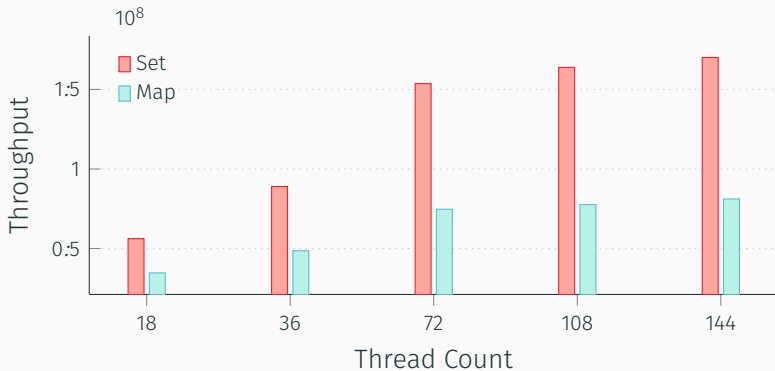
Solution: New-Root



Moving On to Maps

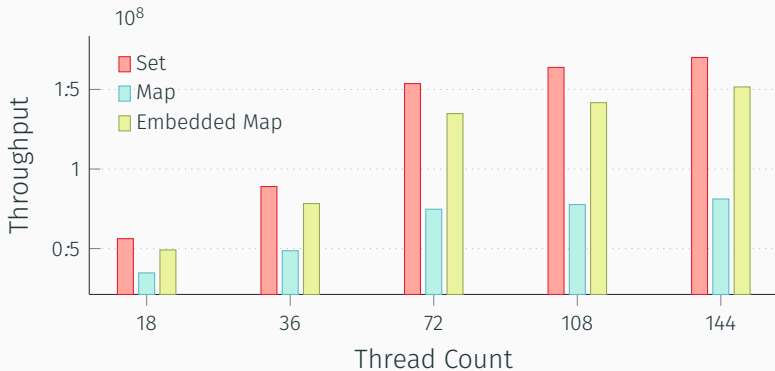


Moving On to Maps



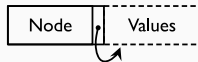
Map

Moving On to Maps



Map

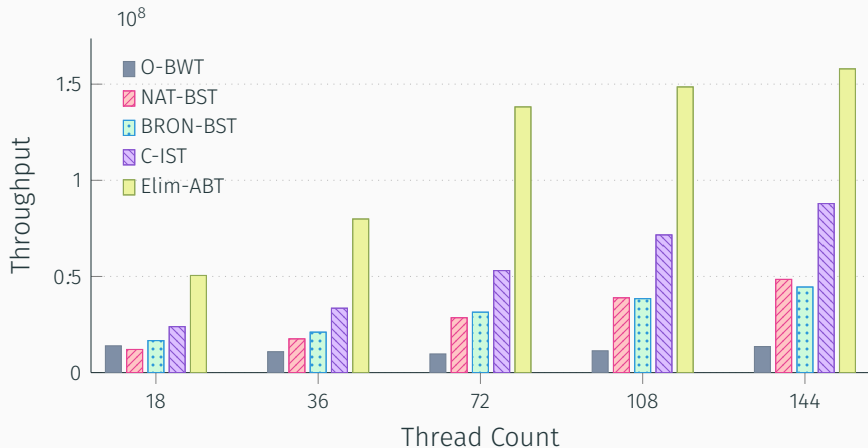
Embedded Map



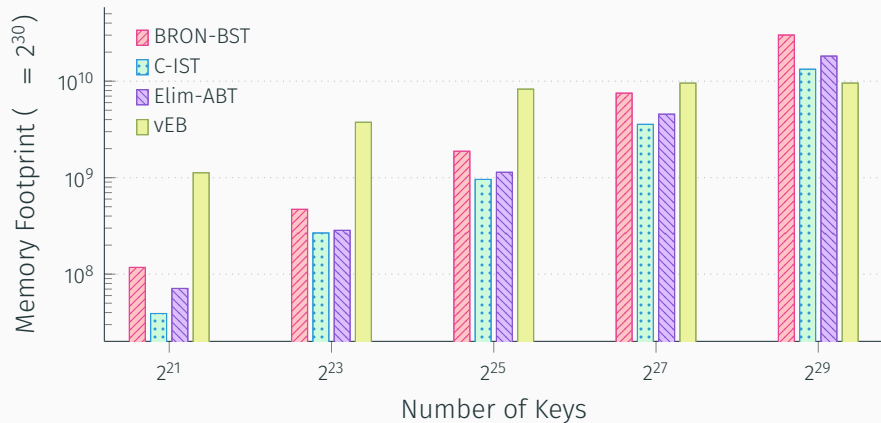
Experimental Evaluation

- Test harness: SetBench [2]
- 4 sockets 18 cores per socket = 72 cores, 144 hardware threads
- Comparison points:
 - **Elim-ABT**: A concurrent (a,b)-tree [3]
 - **C-IST**: A doubly logarithmic interpolation search tree [2]
 - **BRON-BST**: An OCC BST [1]
 - **O-BWT**: The open-source implementation of BW-Tree [5]

Performance

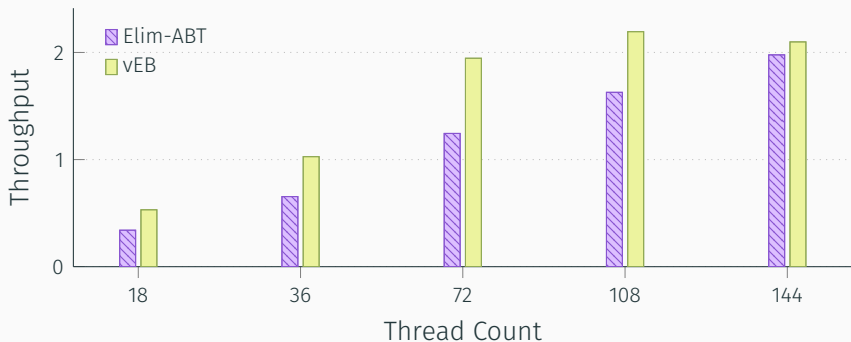


Memory Footprint



Successor Queries

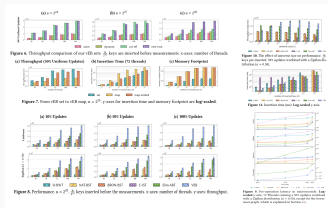
- The successor operation for Elim-ABT is non-linearizable
- Workload: 98% successor queries (Zipfian with $\alpha = 0.50$)



Summary

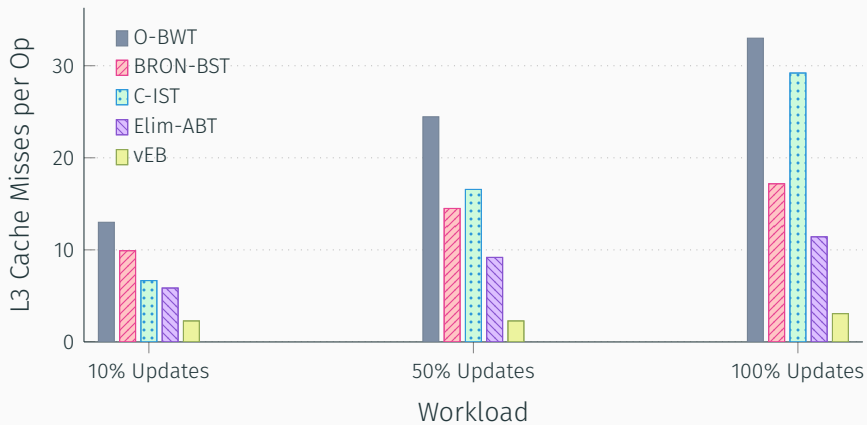
Summary

- Concurrent vEB trees utilizing HTM
- Low-overhead thread synchronization
- On average 5× faster
- More in the paper: <https://bit.ly/htmveb>
 - Zipfian workloads
 - Cache performance
 - Insertion time



Questions?

L3 Cache Misses



References ii

- [3] A. Srivastava and T. Brown.
Elimination (a, b)-trees with fast, durable updates.
In *Proceedings of the 2022 ACM Symposium on Principles and Practice of Parallel Programming*, pages 416–430. ACM, 2022.
- [4] P. van Emde Boas, R. Kaas, and E. Zijlstra.
Design and implementation of an efficient priority queue.
Communications of the ACM, 10:99–127, 1977.
- [5] Z. Wang, A. Pavlo, H. Lim, V. Leis, H. Zhang, M. Kaminsky, and D. G. Andersen.
Building a bw-tree takes more than just buzz words.
In *Proceedings of the 2018 ACM Symposium on Principles and Practice of Parallel Programming*, pages 473–488. ACM, 2018.