



南方科技大学
SOUTHERN UNIVERSITY OF SCIENCE AND TECHNOLOGY



Accelerating Merkle Patricia Trie with GPU

Yangshen Deng*

dengys2022@mail.sustech.edu.cn

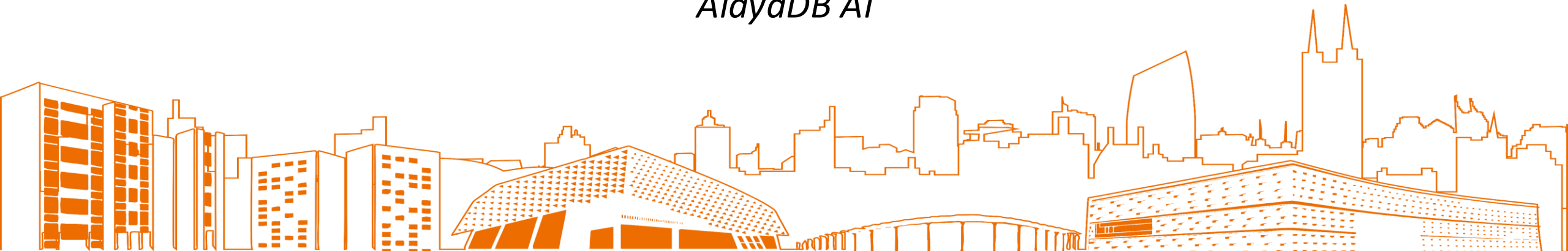
Muxi Yan*

yanmx2022@mail.sustech.edu.cn

Bo Tang

tangb3@sustech.edu.cn

Southern University of Science and Technology
AlayaDB AI





Overview



- **Motivation**
- Challenges
- Solutions
- Results
- Takeaways



Merkle Patricia Trie (MPT)

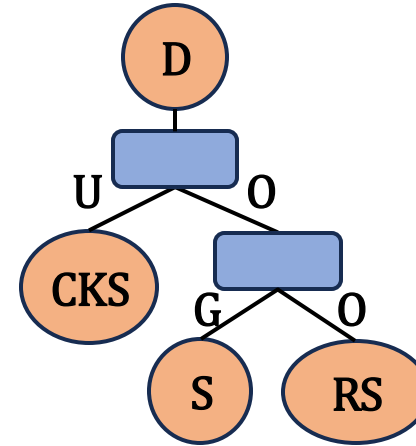


A KV index for immutable data systems

- **Blockchain:** Ethereum, FISCO BCOS, Quorum, etc.
- **Verifiable database:** Alibaba LedgerDB, Spitz, GSSE, etc.

Structure

- Patricia trie (radix trie)





Merkle Patricia Trie (MPT)

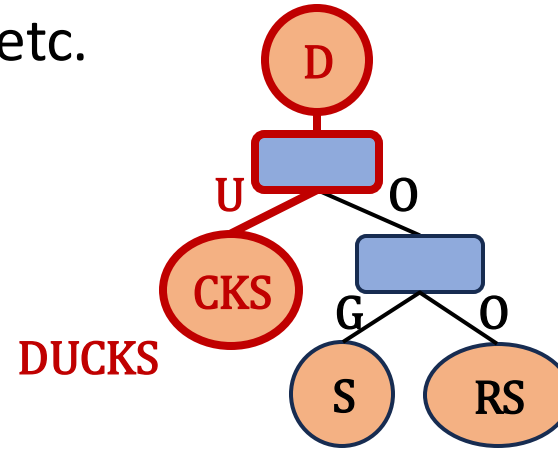


A KV index for immutable data systems

- **Blockchain:** Ethereum, FISCO BCOS, Quorum, etc.
- **Verifiable database:** Alibaba LedgerDB, Spitz, GSSE, etc.

Structure

- Patricia trie (radix trie)





Merkle Patricia Trie (MPT)

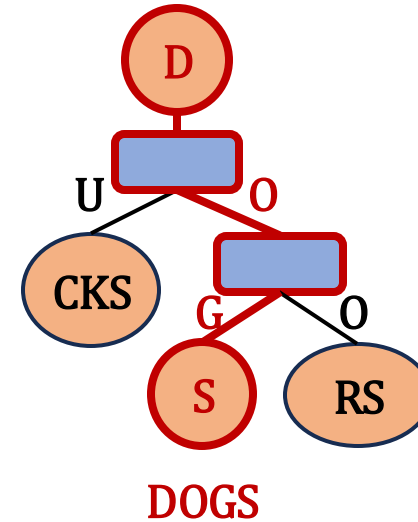


A KV index for immutable data systems

- **Blockchain:** Ethereum, FISCO BCOS, Quorum, etc.
- **Verifiable database:** Alibaba LedgerDB, Spitz, GSSE, etc.

Structure

- Patricia trie (radix trie)





Merkle Patricia Trie (MPT)

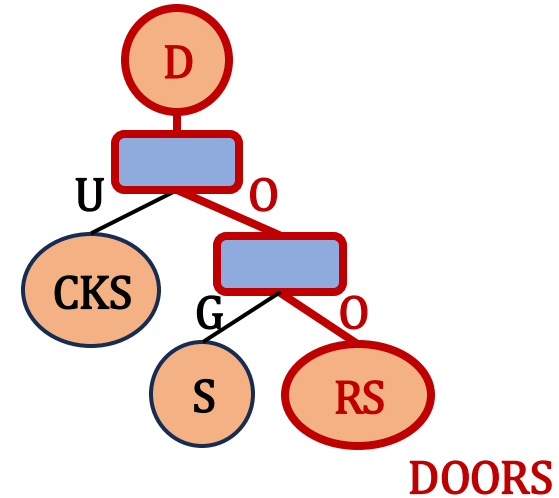


A KV index for immutable data systems

- **Blockchain:** Ethereum, FISCO BCOS, Quorum, etc.
- **Verifiable database:** Alibaba LedgerDB, Spitz, GSSE, etc.

Structure

- Patricia trie (radix trie)





Merkle Patricia Trie (MPT)

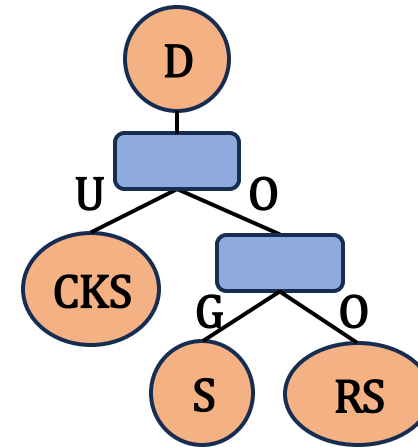


A KV index for immutable data systems

- **Blockchain:** Ethereum, FISCO BCOS, Quorum, etc.
- **Verifiable database:** Alibaba LedgerDB, Spitz, GSSE, etc.

Structure

- Patricia trie (radix trie)
- Cryptographical hashing (merkle tree)

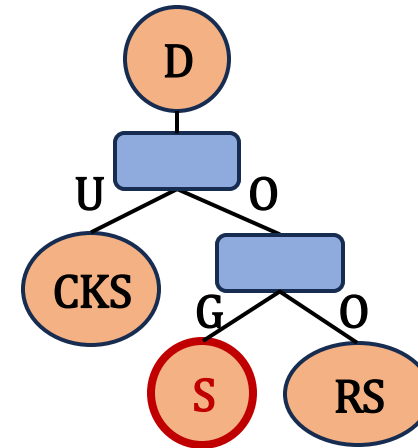


A KV index for immutable data systems

- **Blockchain:** Ethereum, FISCO BCOS, Quorum, etc.
- **Verifiable database:** Alibaba LedgerDB, Spitz, GSSE, etc.

Structure

- Patricia trie (radix trie)
- Cryptographical hashing (merkle tree)



Root hash =

$H(S)$



Merkle Patricia Trie (MPT)

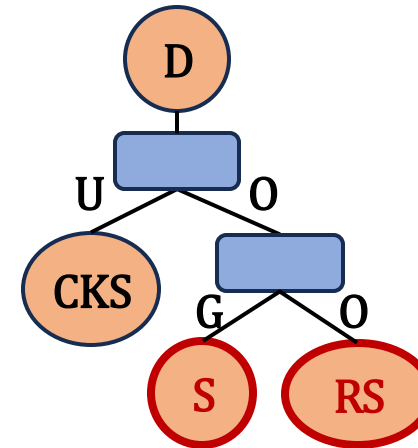


A KV index for immutable data systems

- **Blockchain:** Ethereum, FISCO BCOS, Quorum, etc.
- **Verifiable database:** Alibaba LedgerDB, Spitz, GSSE, etc.

Structure

- Patricia trie (radix trie)
- Cryptographical hashing (merkle tree)



$$\text{Root hash} = H(S) \quad H(\mathbf{RS})$$



Merkle Patricia Trie (MPT)

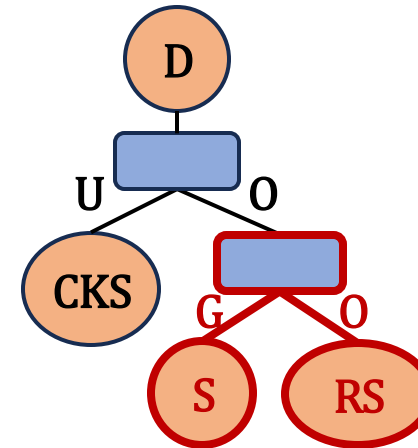


A KV index for immutable data systems

- **Blockchain:** Ethereum, FISCO BCOS, Quorum, etc.
- **Verifiable database:** Alibaba LedgerDB, Spitz, GSSE, etc.

Structure

- Patricia trie (radix trie)
- Cryptographical hashing (merkle tree)



Root hash =

$$H(\mathbf{G}, H(\mathbf{S}), \mathbf{O}, H(\mathbf{RS}))$$



Merkle Patricia Trie (MPT)

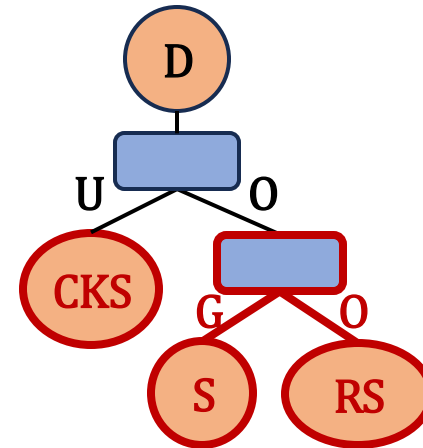


A KV index for immutable data systems

- **Blockchain:** Ethereum, FISCO BCOS, Quorum, etc.
- **Verifiable database:** Alibaba LedgerDB, Spitz, GSSE, etc.

Structure

- Patricia trie (radix trie)
- Cryptographical hashing (merkle tree)



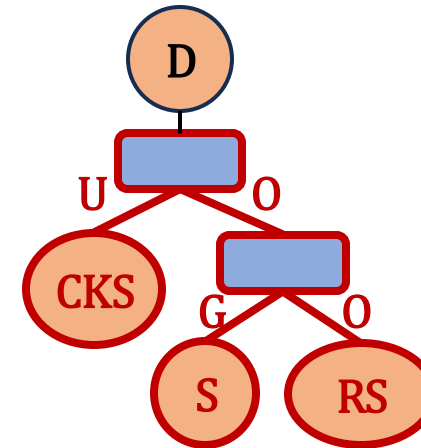
$$\text{Root hash} = H(\mathbf{CKS}) \quad H(\mathbf{G}, H(\mathbf{S}), \mathbf{O}, H(\mathbf{RS}))$$

A KV index for immutable data systems

- **Blockchain:** Ethereum, FISCO BCOS, Quorum, etc.
- **Verifiable database:** Alibaba LedgerDB, Spitz, GSSE, etc.

Structure

- Patricia trie (radix trie)
- Cryptographical hashing (merkle tree)



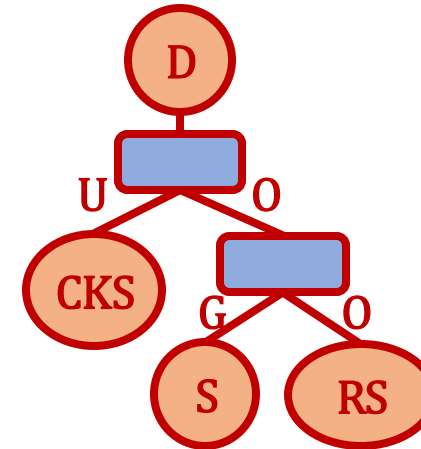
$$\text{Root hash} = H(\mathbf{U}, H(\mathbf{CKS}), \mathbf{O}, H(\mathbf{G}, H(\mathbf{S}), \mathbf{O}, H(\mathbf{RS})))$$

A KV index for immutable data systems

- **Blockchain:** Ethereum, FISCO BCOS, Quorum, etc.
- **Verifiable database:** Alibaba LedgerDB, Spitz, GSSE, etc.

Structure

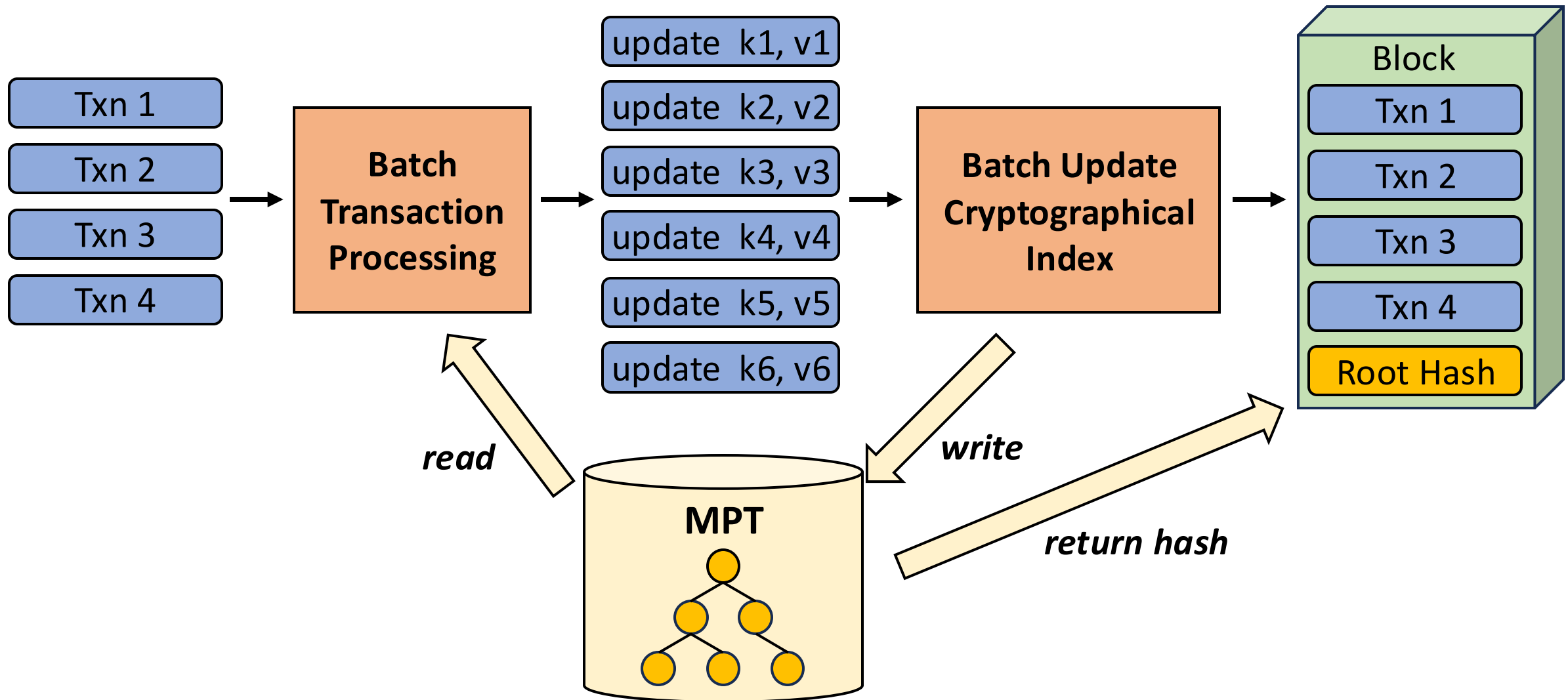
- Patricia trie (radix trie)
- Cryptographical hashing (merkle tree)



$$\text{Root hash} = H(\mathbf{D}, H(\mathbf{U}, H(\mathbf{CKS}), \mathbf{O}, H(\mathbf{G}, H(\mathbf{S}), \mathbf{O}, H(\mathbf{RS}))))$$



Immutable Databases: e.g., Blockchain



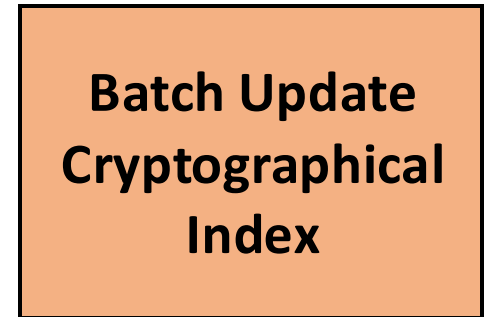
The bottleneck: MPT update

- Heavy writes to index structure

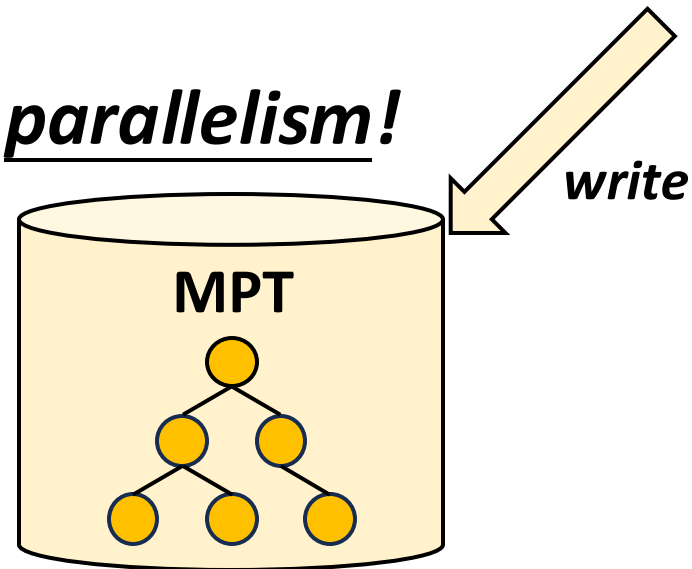
Observation 1: The writes are batched

- Heavy hash computation

Observation 2: Compute-sensitive and easy to parallelize



GPU processes in batch with massive parallelism!





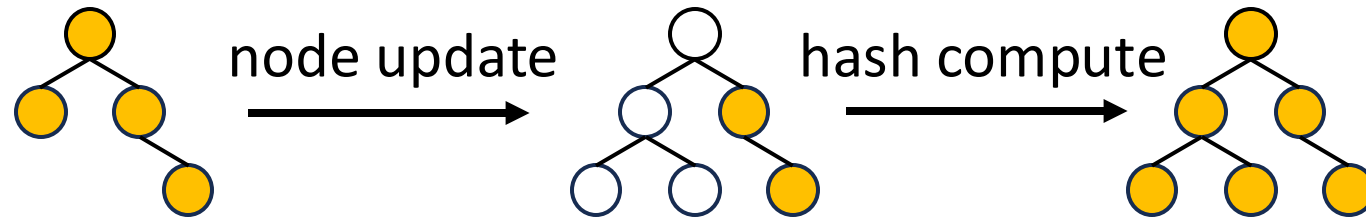
Overview



- Motivation
- **Challenges**
- Solutions
- Results
- Takeaways

Separate into two phases

- **Node update:** concurrent modifications to index structure.
- **Hash compute:** read-write dependency in concurrent hash updates.





Overview



- Motivation
- Challenges
- **Solutions**
- Results
- Takeaways



Solution



Node Update

- PhaseNU
- LockNU

Hash Compute

- PhaseHC



Node Update

- PhaseNU (lock-free)
- LockNU (lock-base)

Hash Compute

- PhaseHC



Solution



Node Update

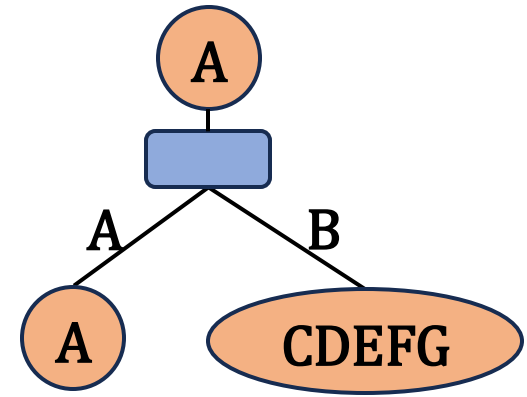
- PhaseNU
- LockNU

Hash Compute

- PhaseHC

Key Idea: Proactively expanding nodes and eliminating locks.

Without PhaseNU:



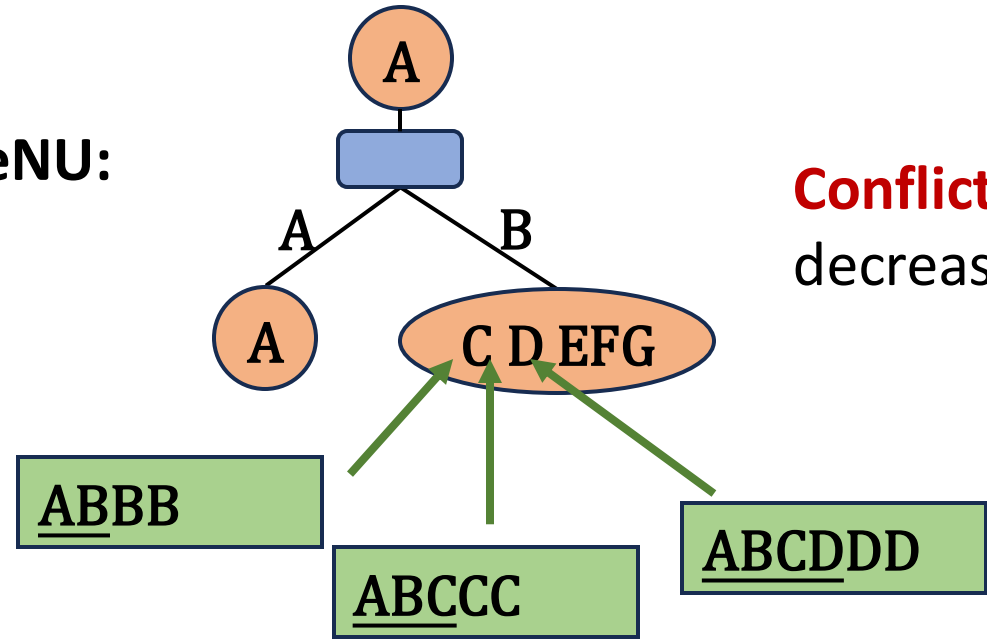
ABBB

ABCCC

ABCD DD

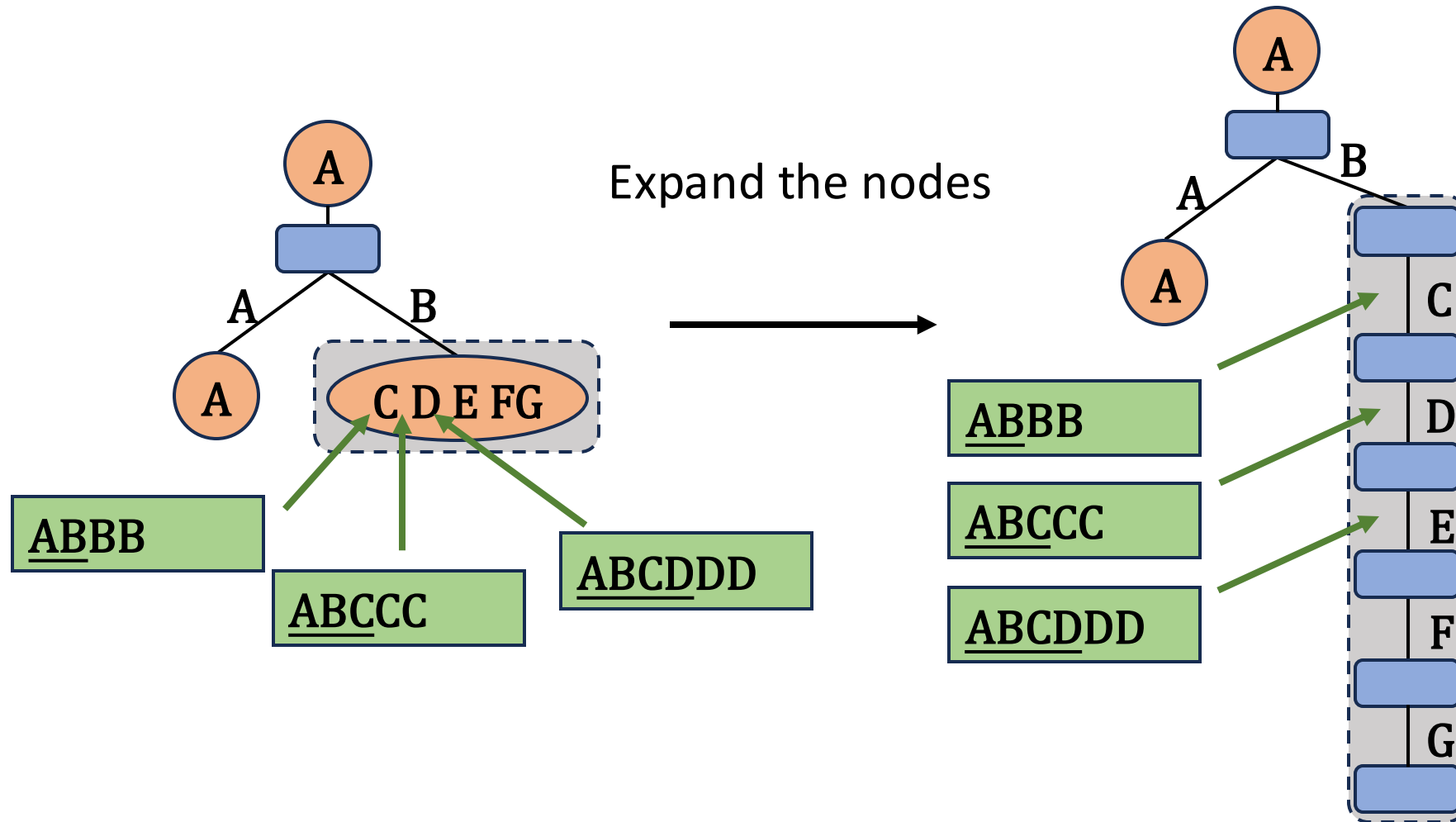
Key Idea: Proactively expanding nodes and eliminating locks.

Without PhaseNU:

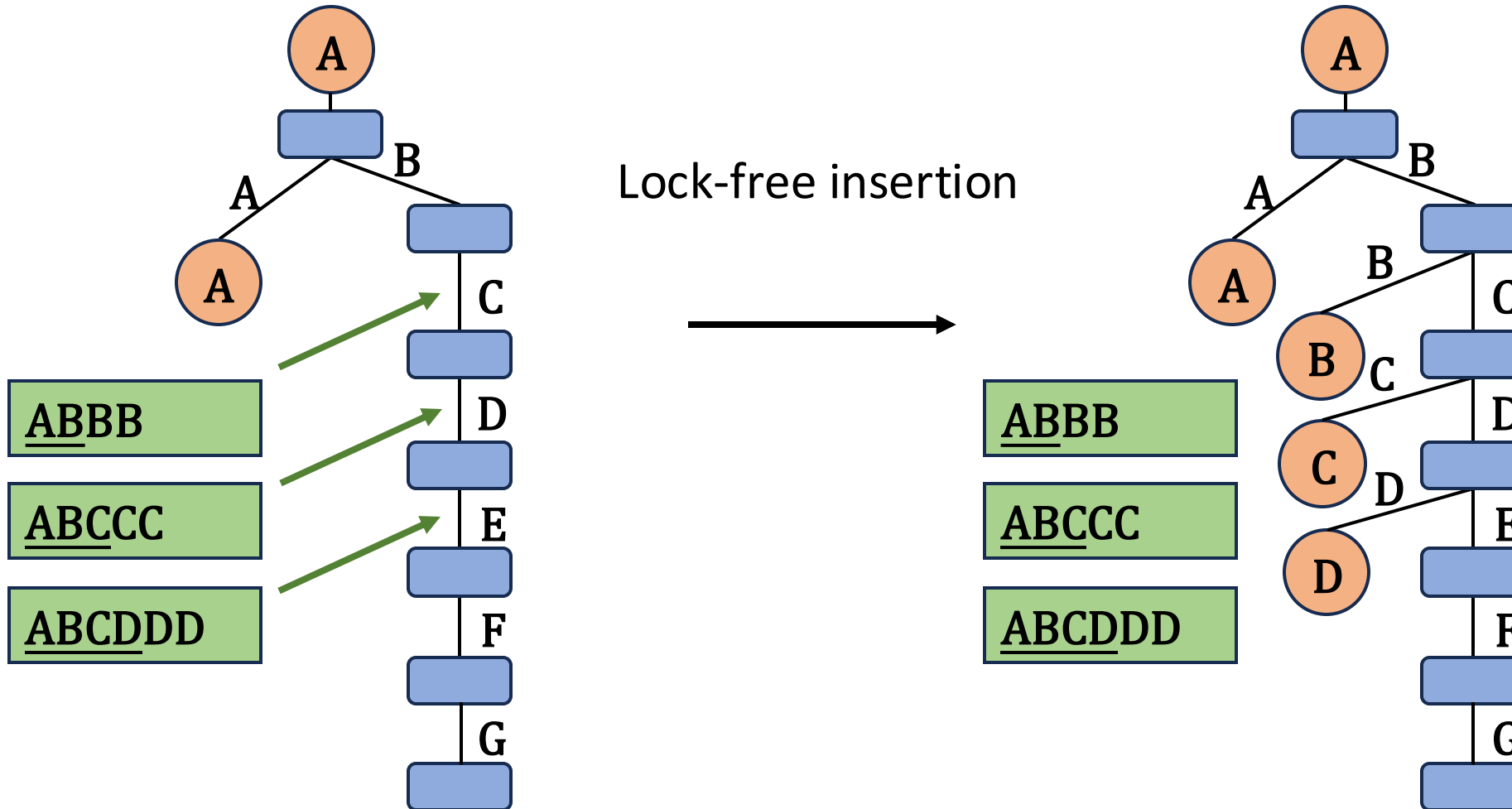


Conflict: parallelism degree decreases from 3 to 1

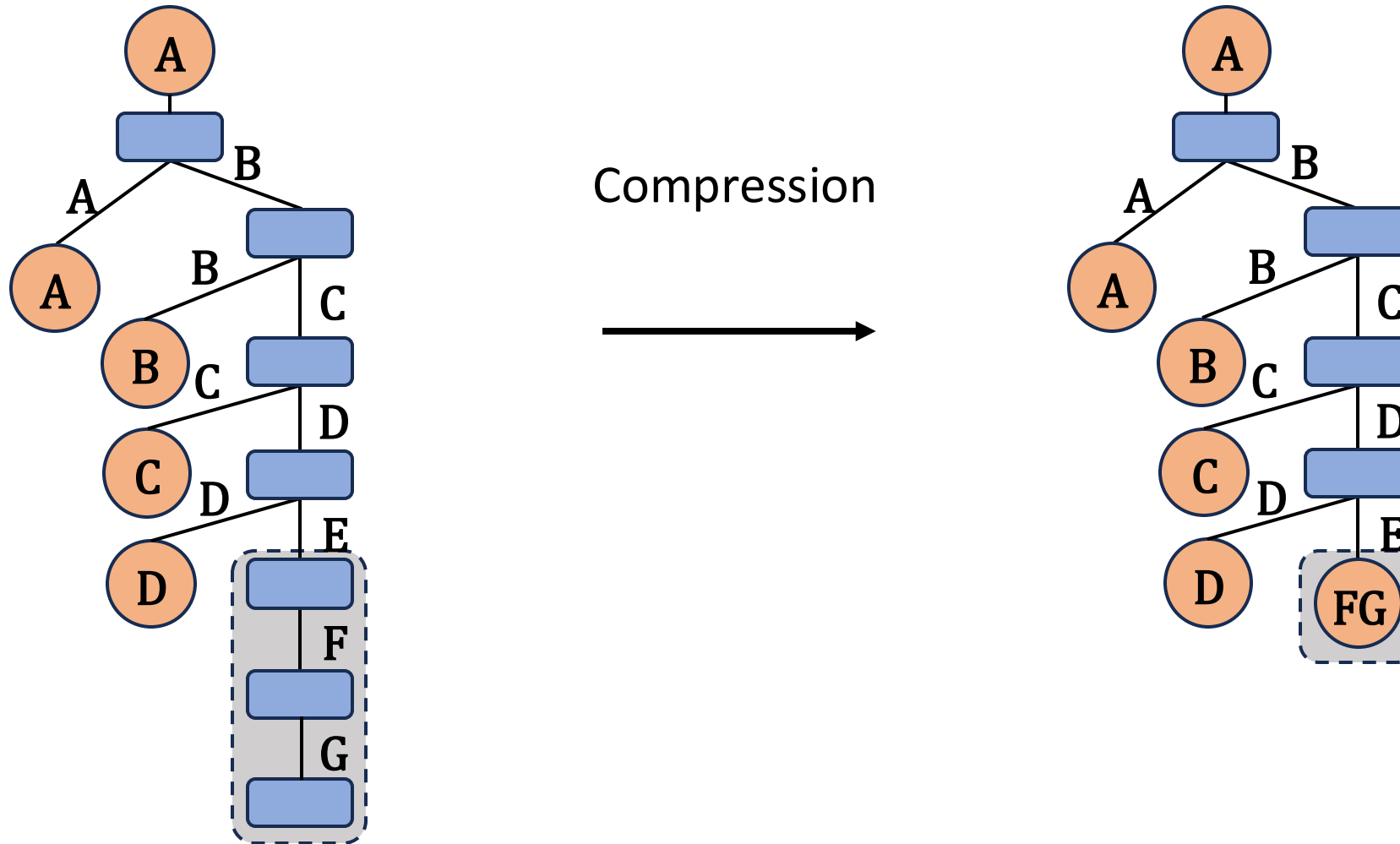
Key Idea: Proactively expanding nodes and eliminating locks.



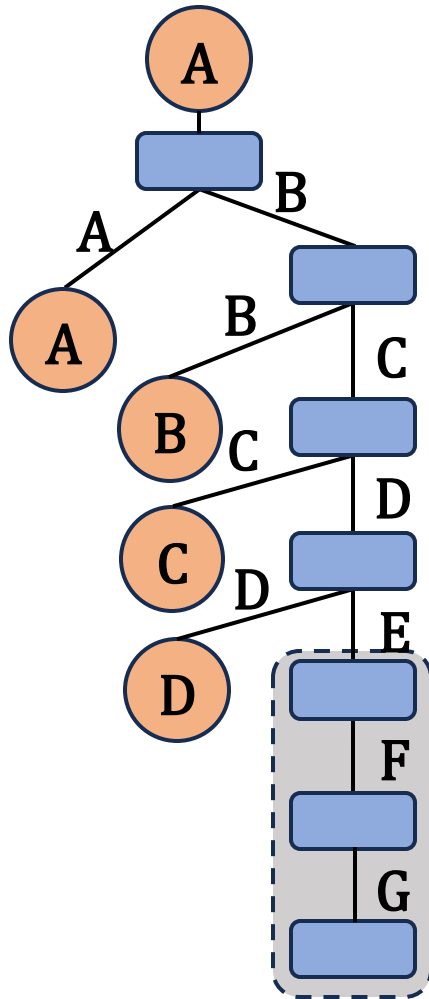
Key Idea: Proactively expanding nodes and eliminating locks.



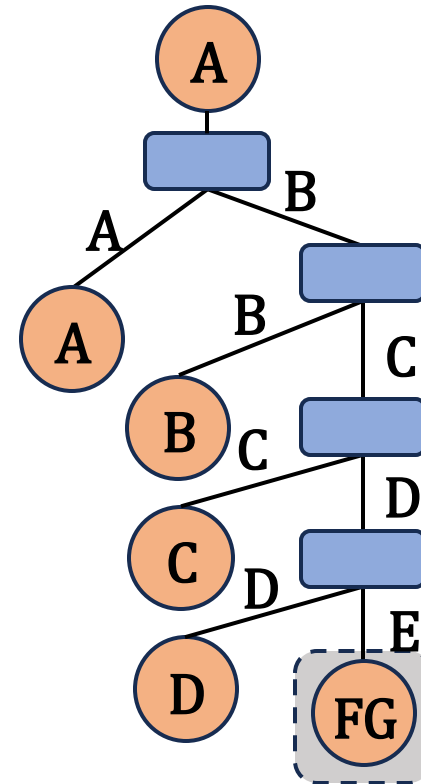
Key Idea: Proactively expanding nodes and eliminating locks.



Limitation: Complexity increases with the length of key.



Compression





Solution



Node Update

- PhaseNU
- **LockNU**

Hash Compute

- PhaseHC



Key Idea: Optimistic lock coupling, inspired by Viktor Leis.^[1]

[1] Viktor Leis, Florian Scheibner, Alfons Kemper, and Thomas Neumann. 2016. The ART of practical synchronization. In DaMoN '16.



Key Idea: Optimistic lock coupling, inspired by Viktor Leis.^[1]

- Read with **optimistic read lock**
 - Get the **version** when lock
 - Verify the **version** when unlock
- Write with **pessimistic write lock**

[1] Viktor Leis, Florian Scheibner, Alfons Kemper, and Thomas Neumann. 2016. The ART of practical synchronization. In DaMoN '16.



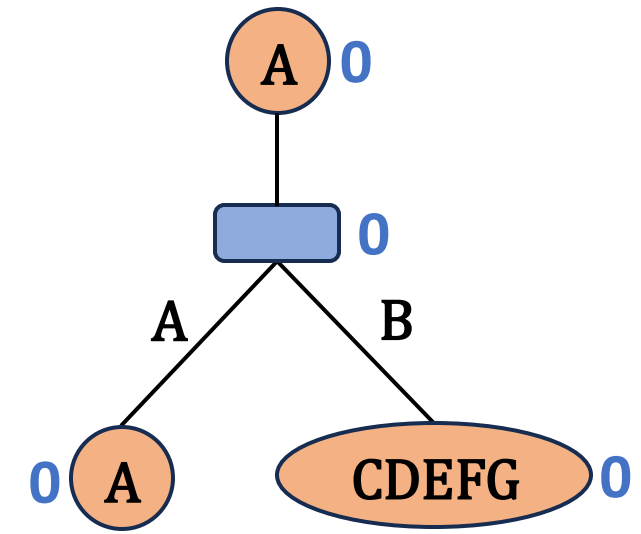
LockNU: An Example of Inserting

ABCDEEE



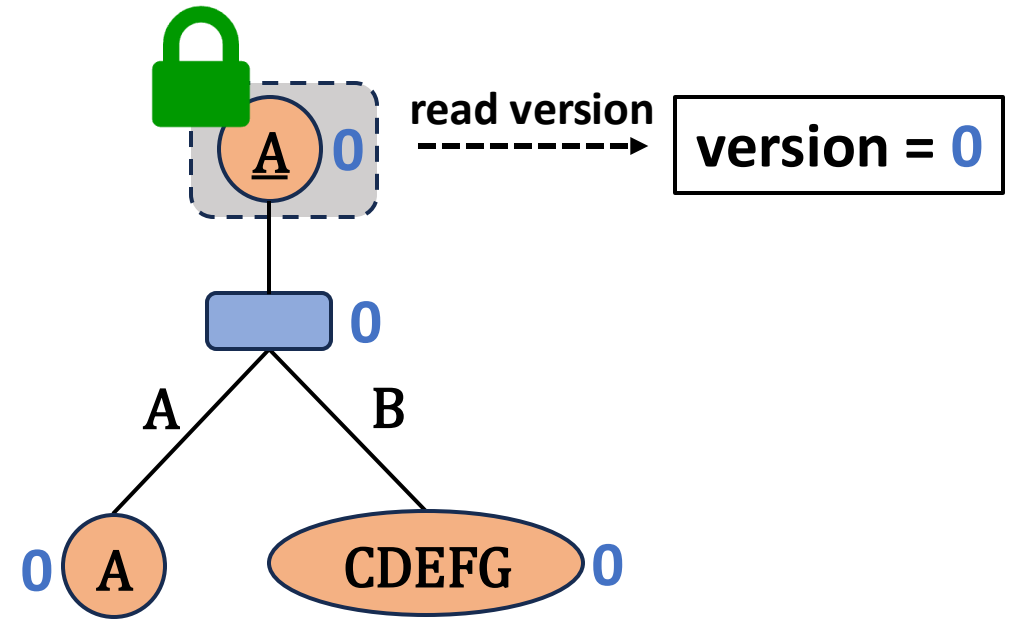
Key Idea: Optimistic lock coupling.

- Read with **optimistic read lock**
 - Get the **version** when lock
 - Verify the **version** when unlock
- Write with **pessimistic write lock**



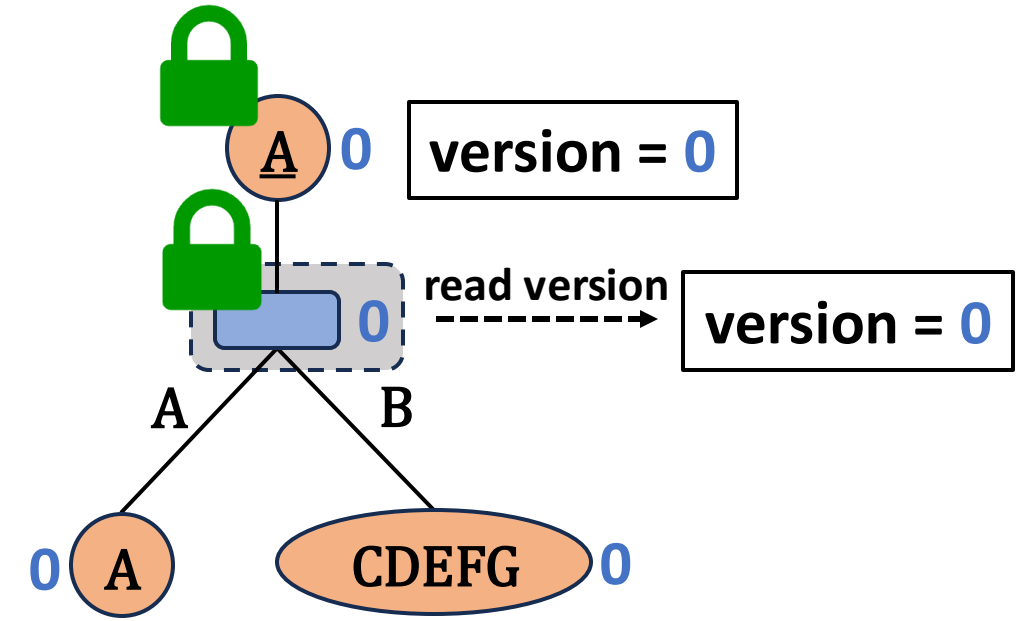
Key Idea: Optimistic lock coupling.

- Read with **optimistic read lock**
 - Get the **version** when lock
 - Verify the **version** when unlock
- Write with **pessimistic write lock**



Key Idea: Optimistic lock coupling.

- Read with **optimistic read lock**
 - Get the **version** when lock
 - Verify the **version** when unlock
- Write with **pessimistic write lock**





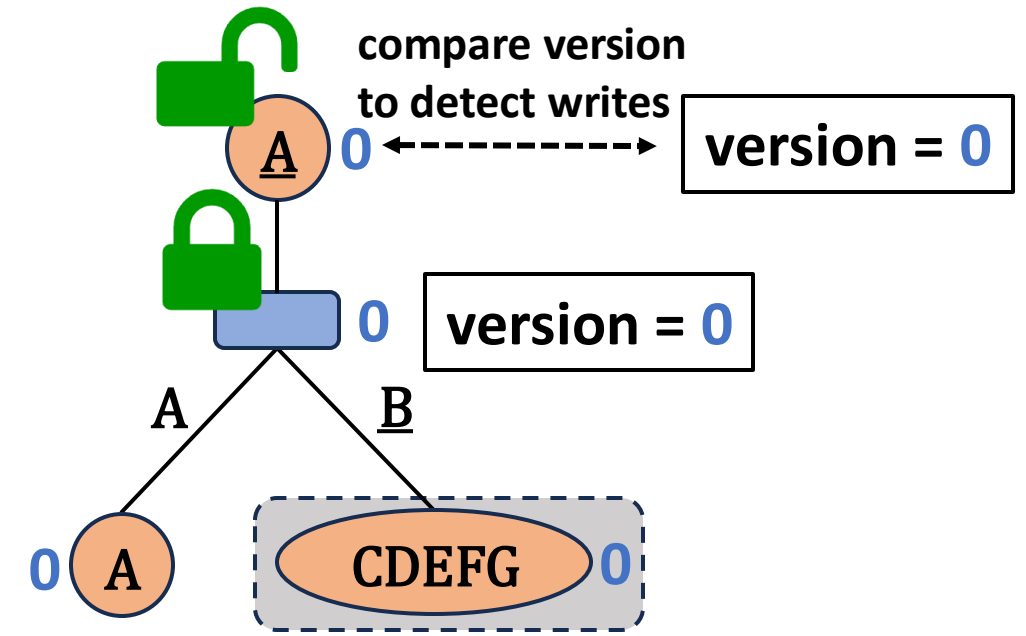
LockNU: An Example of Inserting

ABCDEEEE



Key Idea: Optimistic lock coupling.

- Read with **optimistic read lock**
 - Get the **version** when lock
 - Verify the **version** when unlock
- Write with **pessimistic write lock**





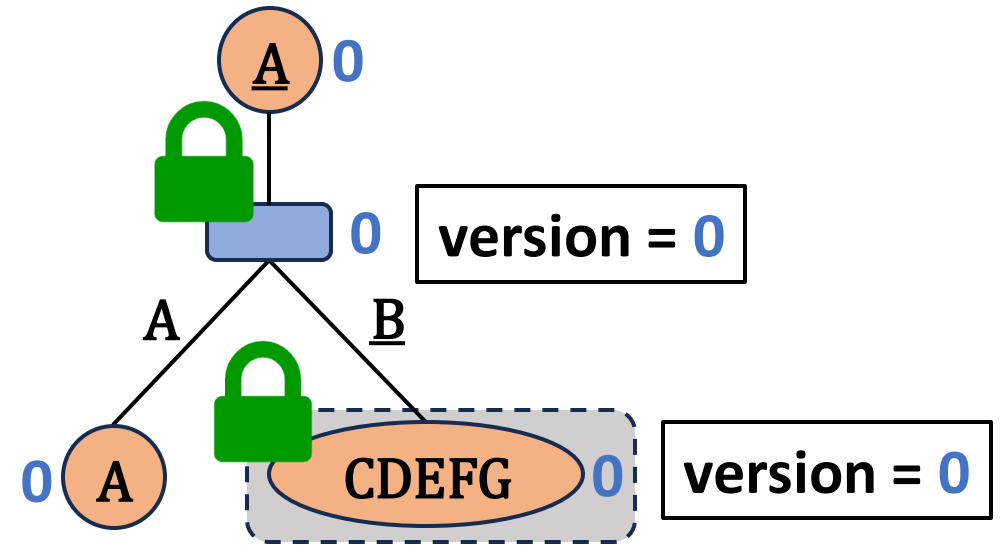
LockNU: An Example of Inserting

ABCDEEE



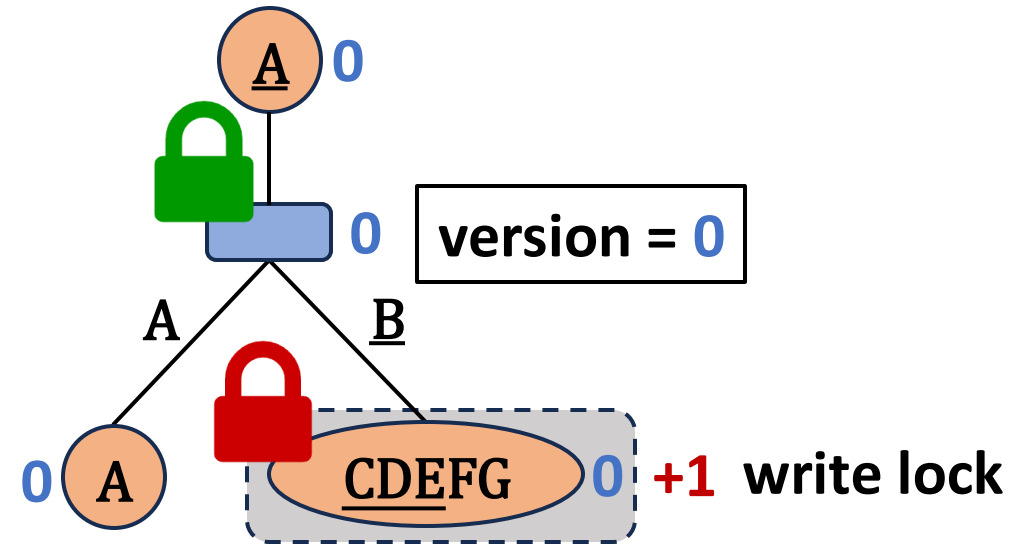
Key Idea: Optimistic lock coupling.

- Read with **optimistic read lock**
 - Get the **version** when lock
 - Verify the **version** when unlock
- Write with **pessimistic write lock**



Key Idea: Optimistic lock coupling.

- Read with **optimistic read lock**
 - Get the **version** when lock
 - Verify the **version** when unlock
- Write with **pessimistic write lock**





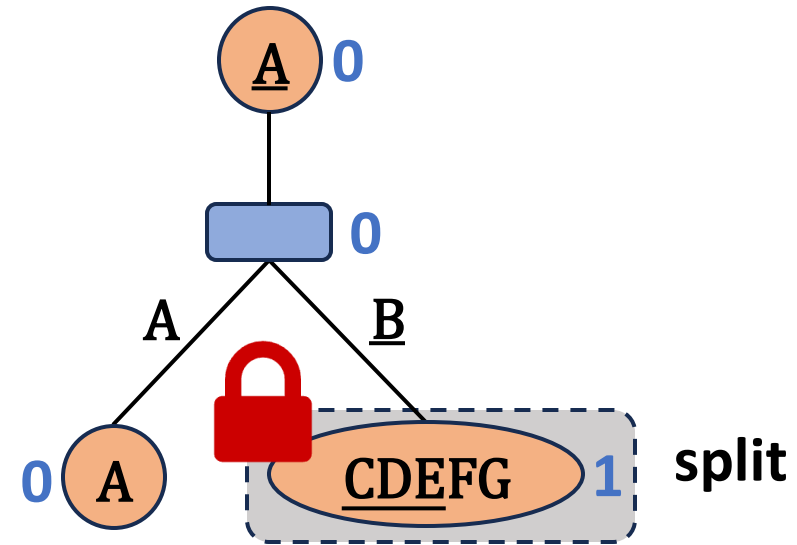
LockNU: An Example of Inserting

ABCDEEE



Key Idea: Optimistic lock coupling.

- Read with **optimistic read lock**
 - Get the **version** when lock
 - Verify the **version** when unlock
- Write with **pessimistic write lock**





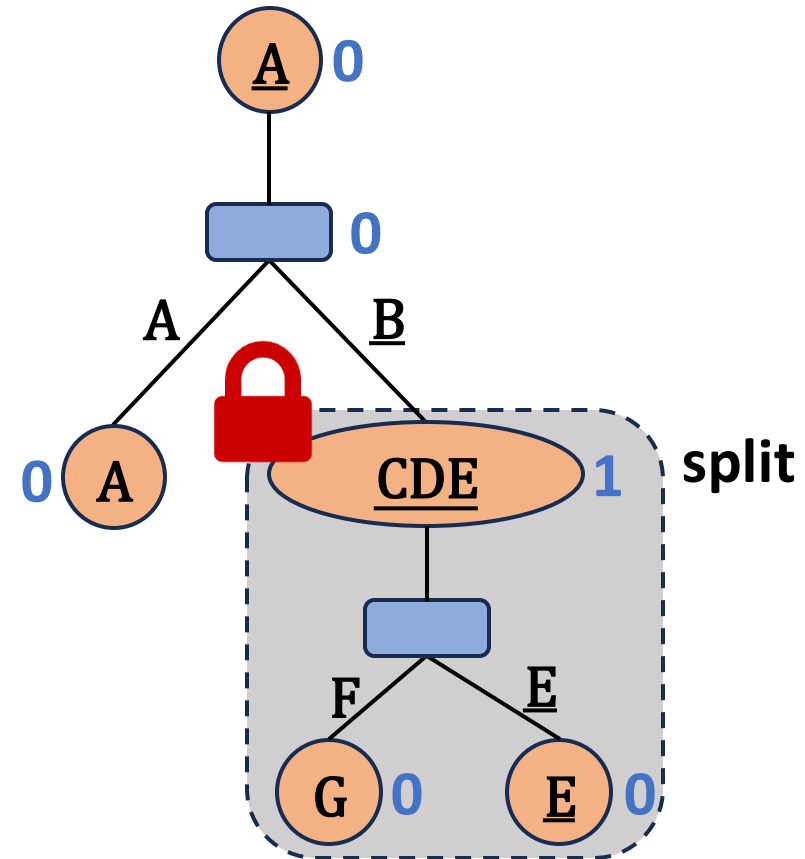
LockNU: An Example of Inserting

ABCDEEE



Key Idea: Optimistic lock coupling.

- Read with **optimistic read lock**
 - Get the **version** when lock
 - Verify the **version** when unlock
- Write with **pessimistic write lock**





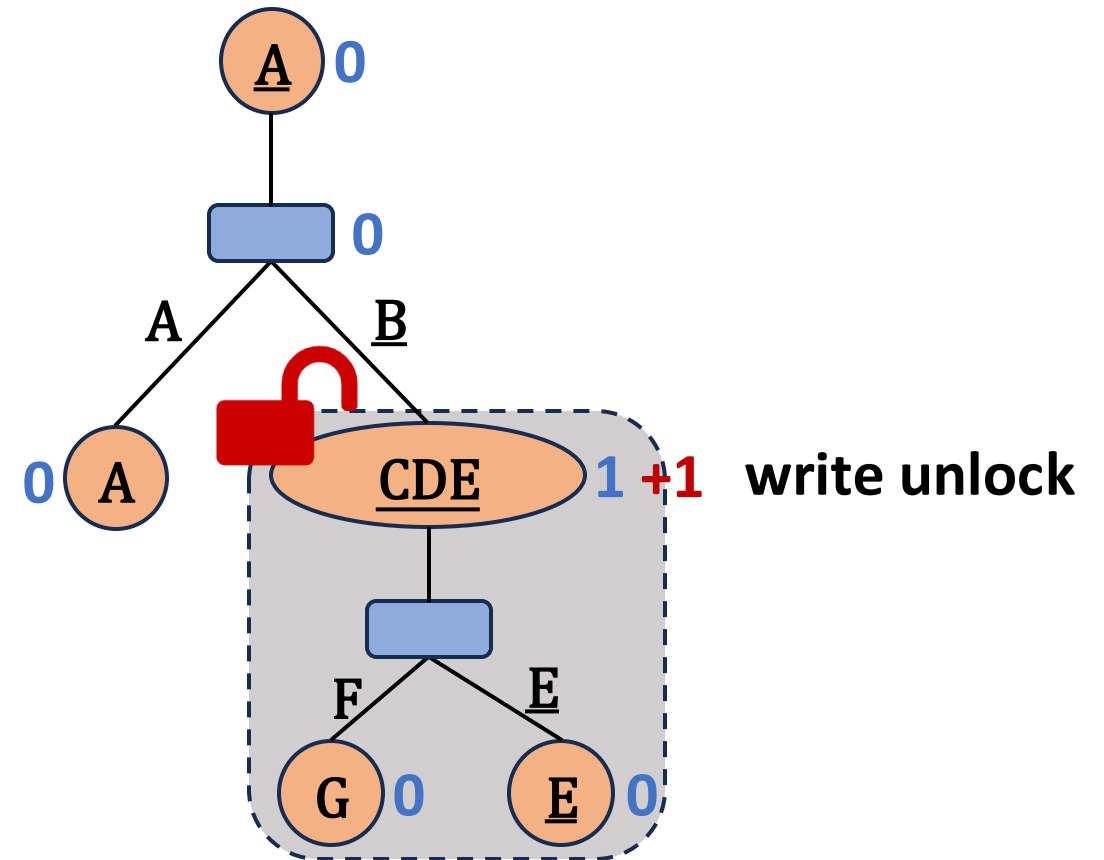
LockNU: An Example of Inserting

ABCDEEE



Key Idea: Optimistic lock coupling.

- Read with **optimistic read lock**
 - Get the **version** when lock
 - Verify the **version** when unlock
- Write with **pessimistic write lock**





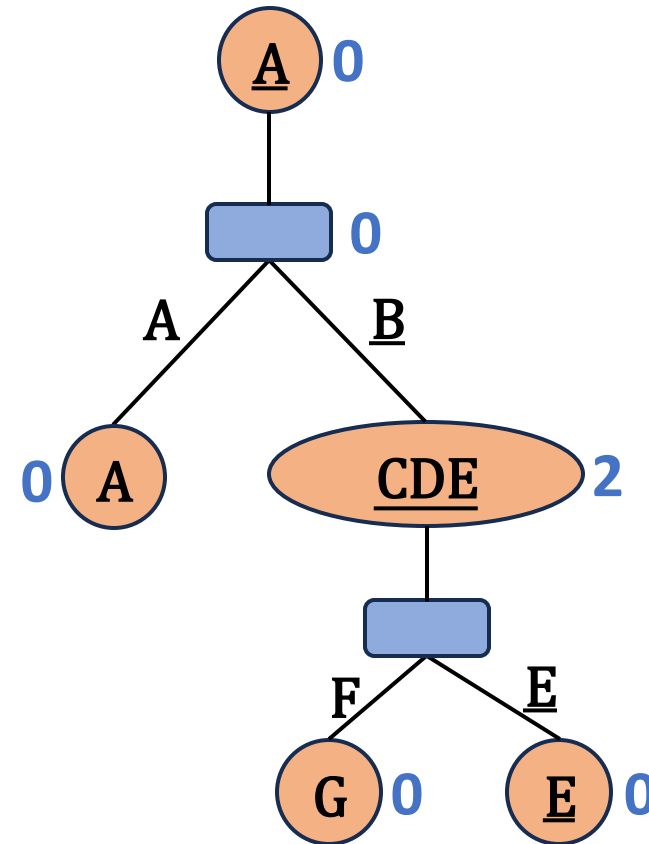
LockNU: An Example of Inserting

ABCDEEE



Key Idea: Optimistic lock coupling.

- Read with **optimistic read lock**
 - Get the **version** when lock
 - Verify the **version** when unlock
- Write with **pessimistic write lock**





Solution



Node Update

- PhaseNU
- LockNU

Hash Compute

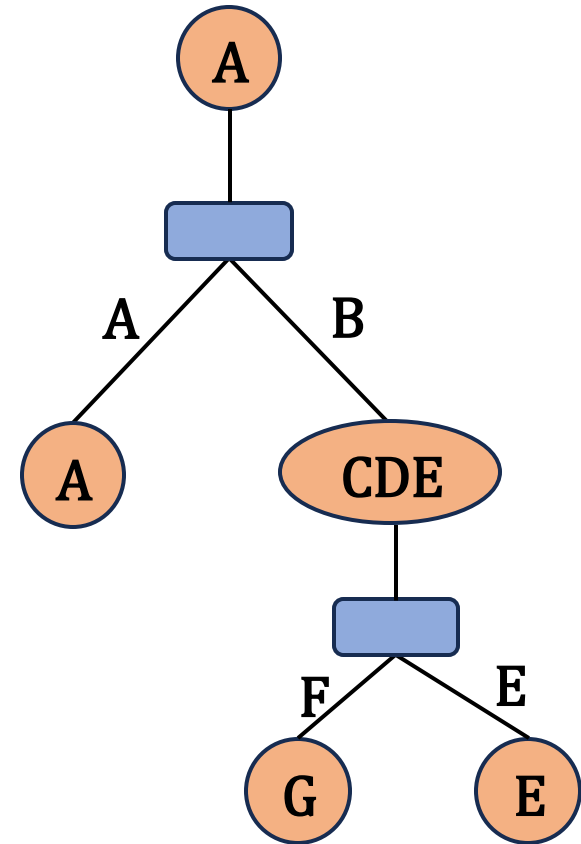
- PhaseHC



PhaseHC: Phase 1 – Marking Phase



Key Idea: Modeling the hash computation into a **dependency** graph.



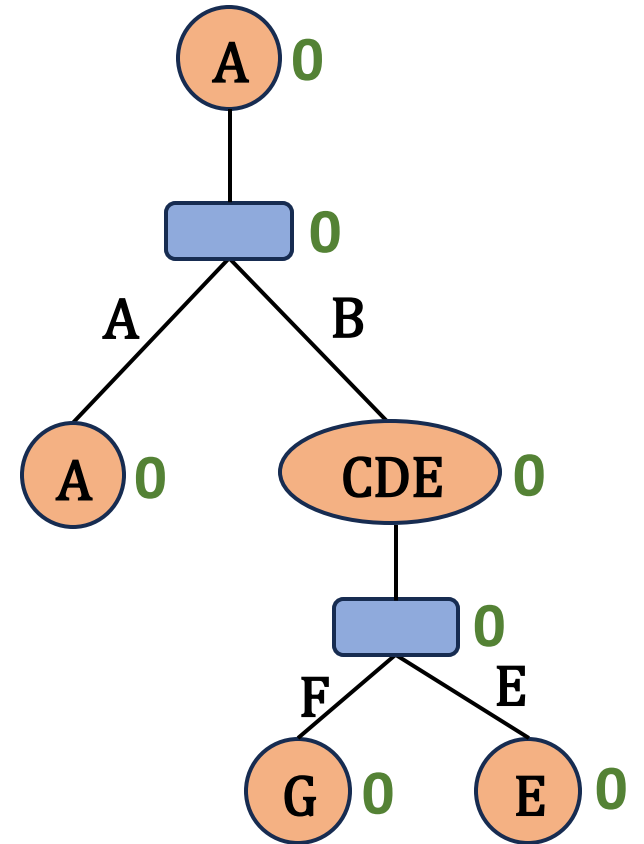


PhaseHC: Phase 1 – Marking Phase



Key Idea: Modeling the hash computation into a **dependency graph**.

Dependency counter: how many hash values it needs?





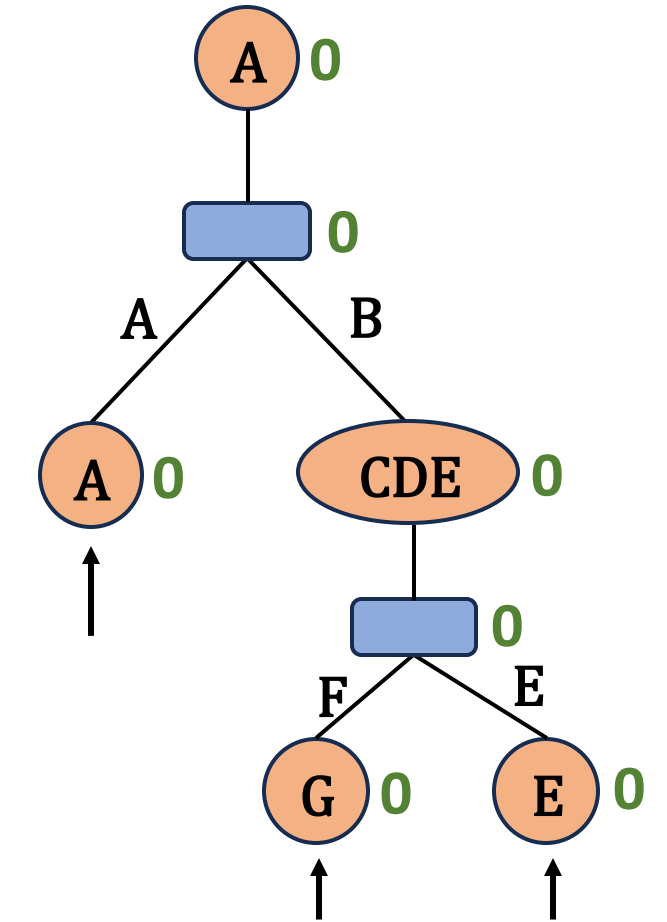
PhaseHC: Phase 1 – Marking Phase



Key Idea: Modeling the hash computation into a **dependency graph**.

Dependency counter: how many hash values it needs?

- Each thread starts from a newly inserted leaf.





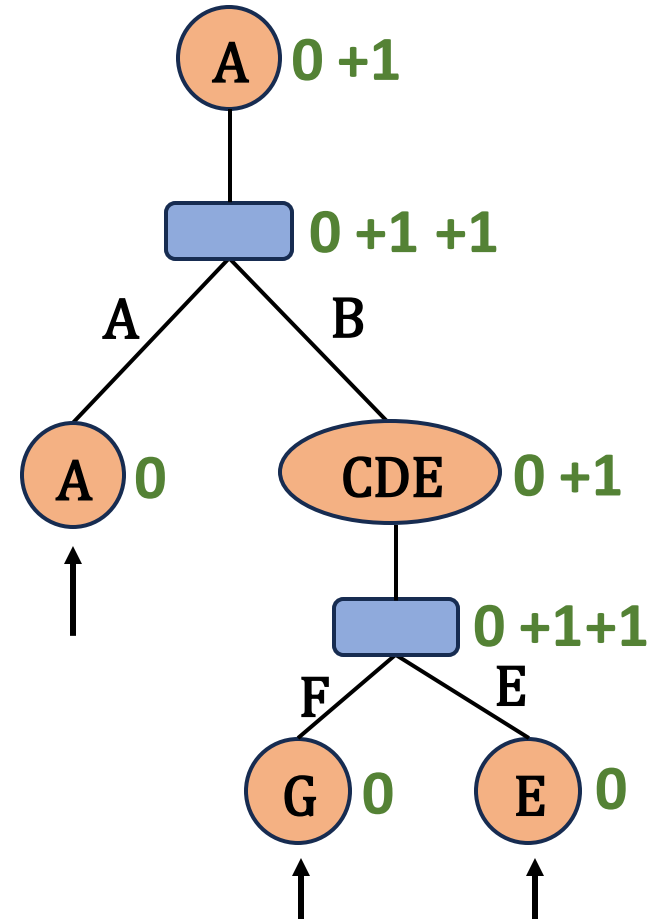
PhaseHC: Phase 1 – Marking Phase



Key Idea: Modeling the hash computation into a **dependency graph**.

Dependency counter: how many hash values it needs?

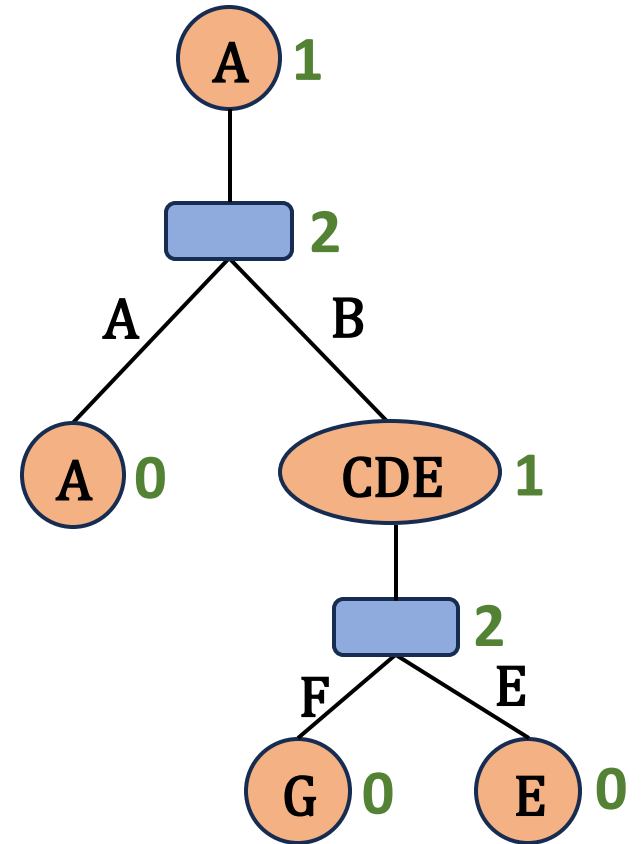
- Each thread starts from a newly inserted leaf.
- Increment counters in the path. (+1)



Key Idea: Modeling the hash computation into a **dependency graph**.

Dependency counter: how many hash values it needs?

- Each thread starts from a newly inserted leaf.
- Increment counters in the path. (+1)



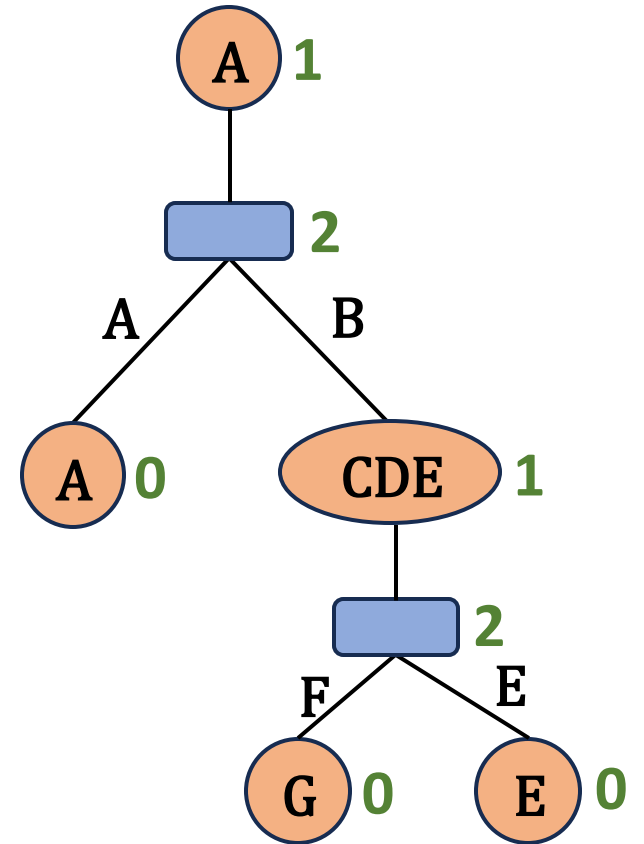


PhaseHC: Phase 2 – Computing Phase



Key Idea: Modeling the hash computation into a **dependency graph**.

The last one arrives takes the job.





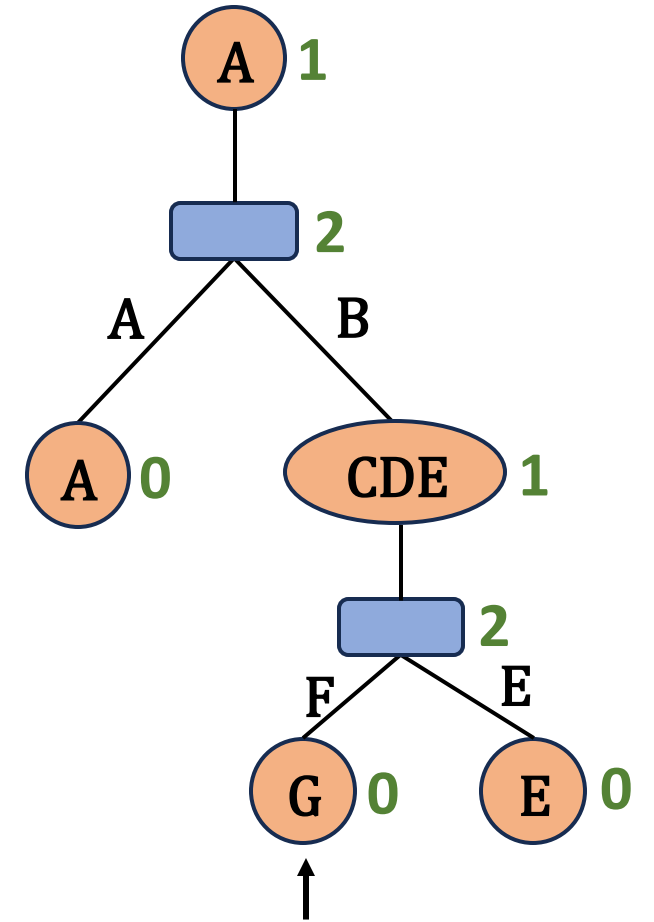
PhaseHC: Phase 2 – Computing Phase



Key Idea: Modeling the hash computation into a **dependency graph**.

The last one arrives takes the job.

- Each thread starts from a newly inserted leaf.





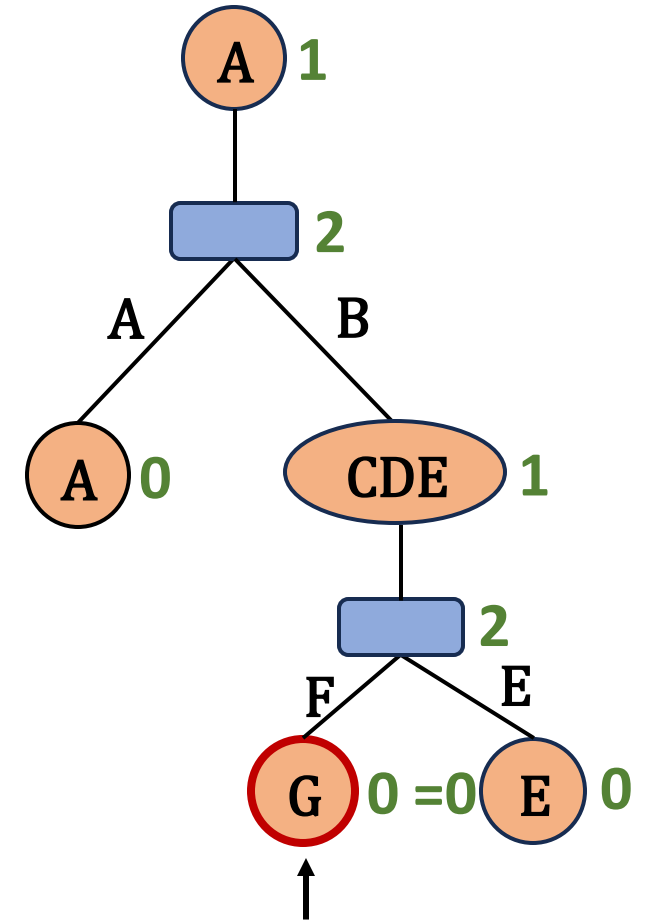
PhaseHC: Phase 2 – Computing Phase



Key Idea: Modeling the hash computation into a **dependency graph**.

The last one arrives takes the job.

- Each thread starts from a newly inserted leaf.
- Proceeds when dependencies are resolved. (=0)





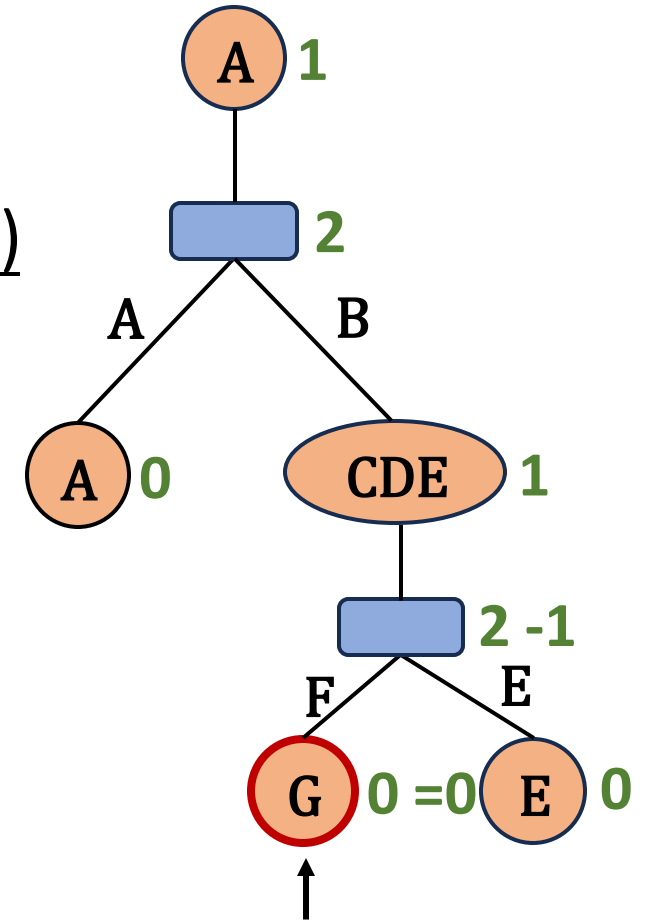
PhaseHC: Phase 2 – Computing Phase



Key Idea: Modeling the hash computation into a **dependency graph**.

The last one arrives takes the job.

- Each thread starts from a newly inserted leaf.
- Proceeds when dependencies are resolved. (=0)
- Compute hash value and decrement the counter (-1)





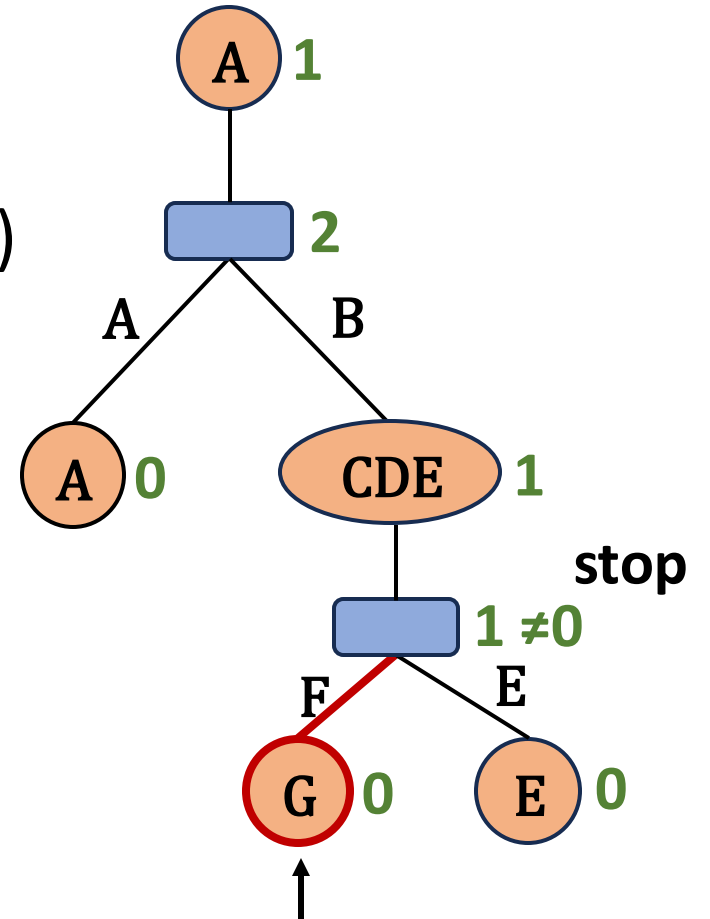
PhaseHC: Phase 2 – Computing Phase



Key Idea: Modeling the hash computation into a **dependency graph**.

The last one arrives takes the job.

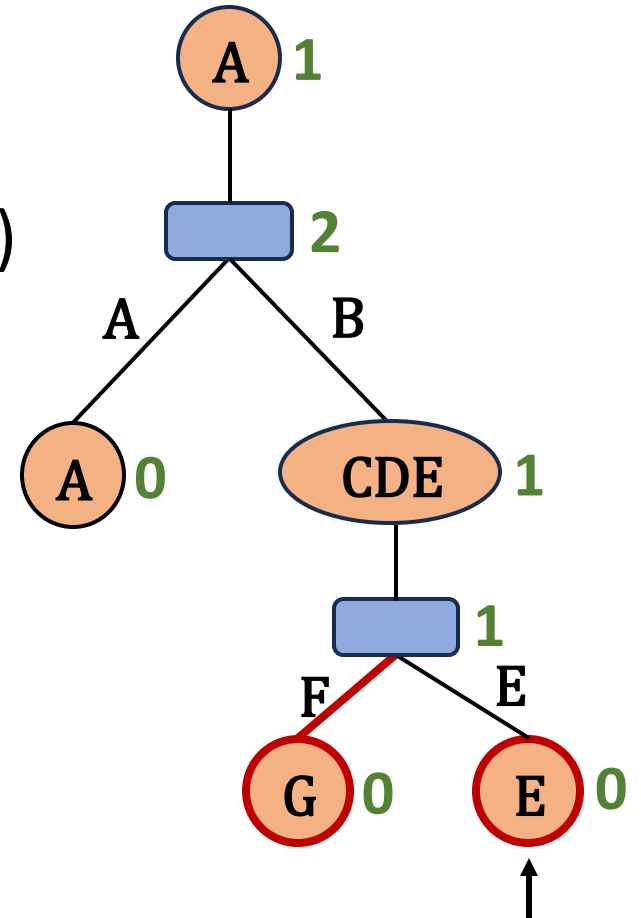
- Each thread starts from a newly inserted leaf.
- Proceeds when dependencies are resolved. (=0)
- Compute hash value and decrement the counter (-1)



Key Idea: Modeling the hash computation into a **dependency graph**.

The last one arrives takes the job.

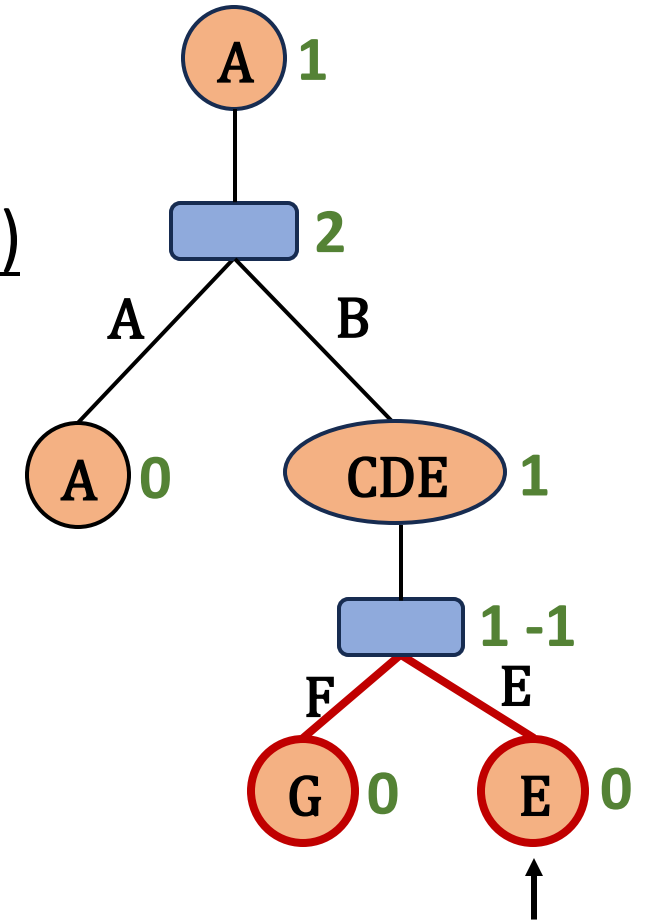
- Each thread starts from a newly inserted leaf.
- Proceeds when dependencies are resolved. (=0)
- Compute hash value and decrement the counter (-1)



Key Idea: Modeling the hash computation into a **dependency graph**.

The last one arrives takes the job.

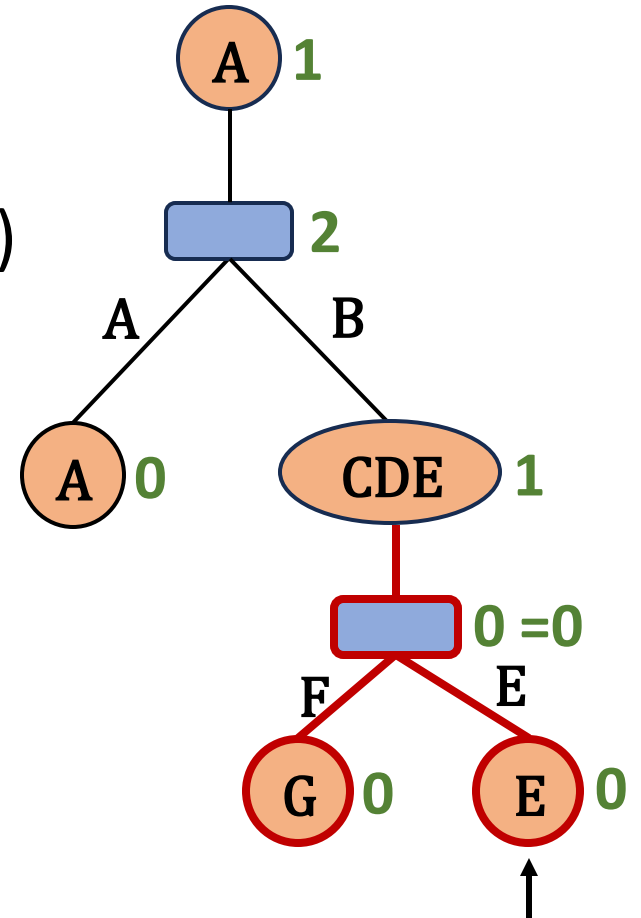
- Each thread starts from a newly inserted leaf.
- Proceeds when dependencies are resolved. (=0)
- Compute hash value and decrement the counter (-1)



Key Idea: Modeling the hash computation into a **dependency graph**.

The last one arrives takes the job.

- Each thread starts from a newly inserted leaf.
- Proceeds when dependencies are resolved. (=0)
- Compute hash value and decrement the counter (-1)





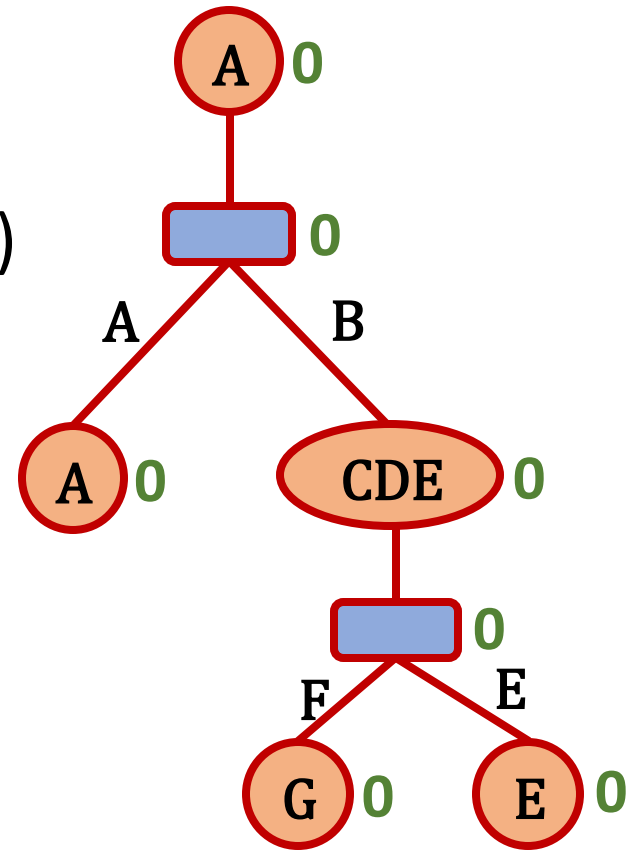
PhaseHC: Phase 2 – Computing Phase



Key Idea: Modeling the hash computation into a **dependency graph**.

The last one arrives takes the job.

- Each thread starts from a newly inserted leaf.
- Proceeds when dependencies are resolved. (=0)
- Compute hash value and decrement the counter (-1)





Overview



- Motivation
- Challenges
- Solutions
- **Results**
- Takeaways



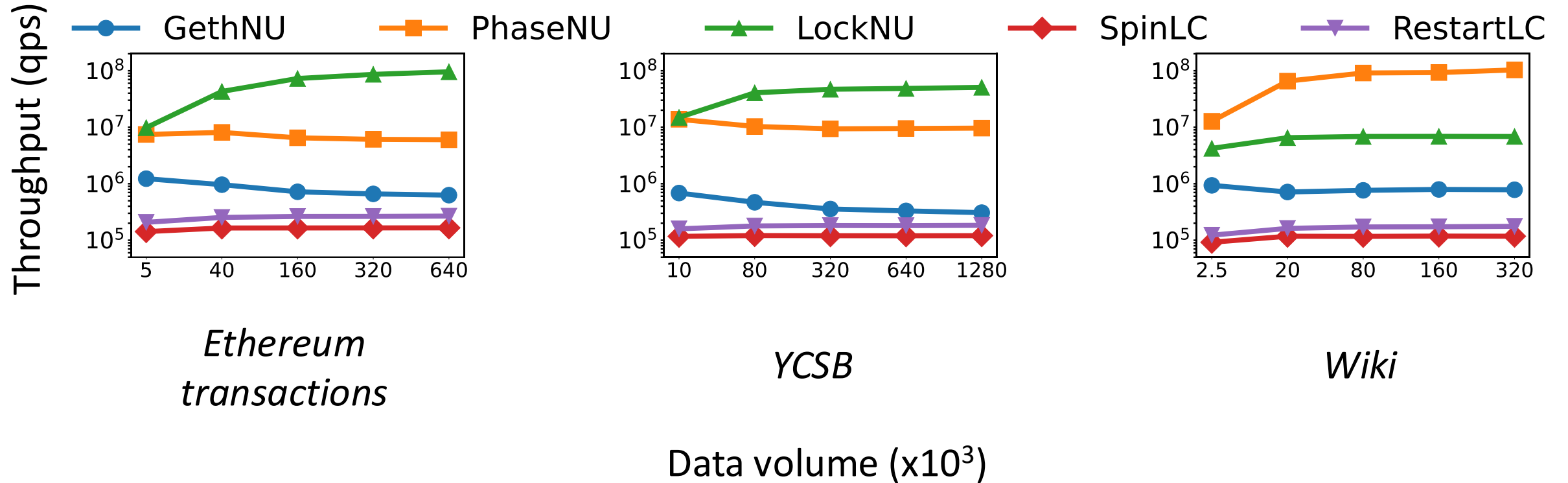
Experiment Setup



- **GPU:** NVIDIA V100
- **CPU:** Intel Xeon Gold 6230R (2 sockets x 26 cores x 2 threads)
- **Baseline system:** Go-Ethereum (Geth)



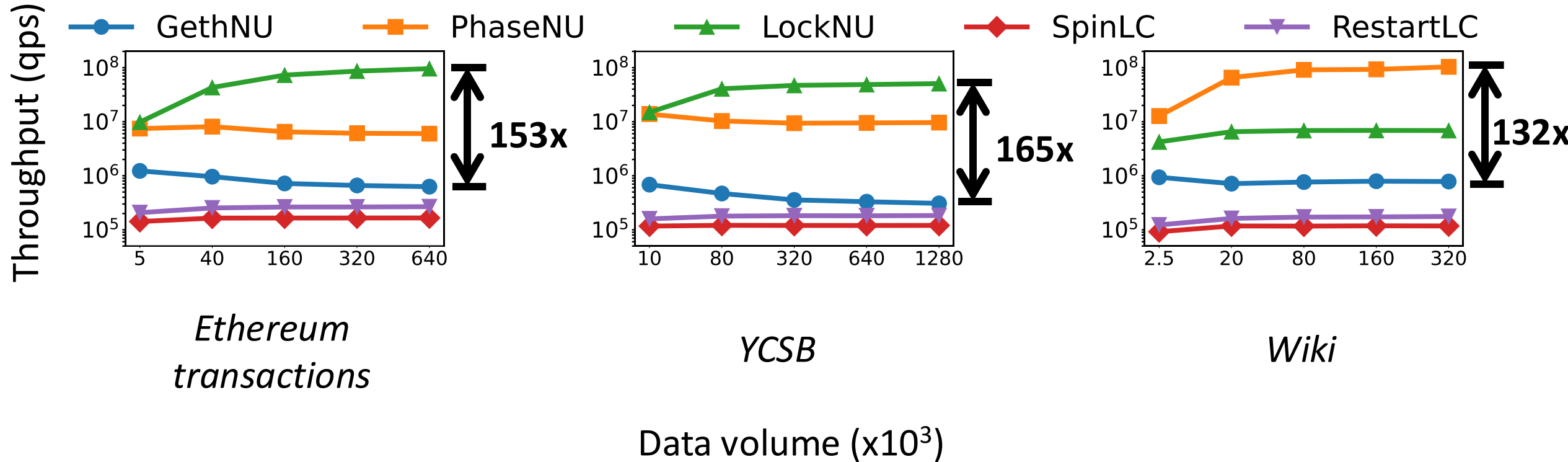
Index Benchmark: Node Update





Index Benchmark: Node Update

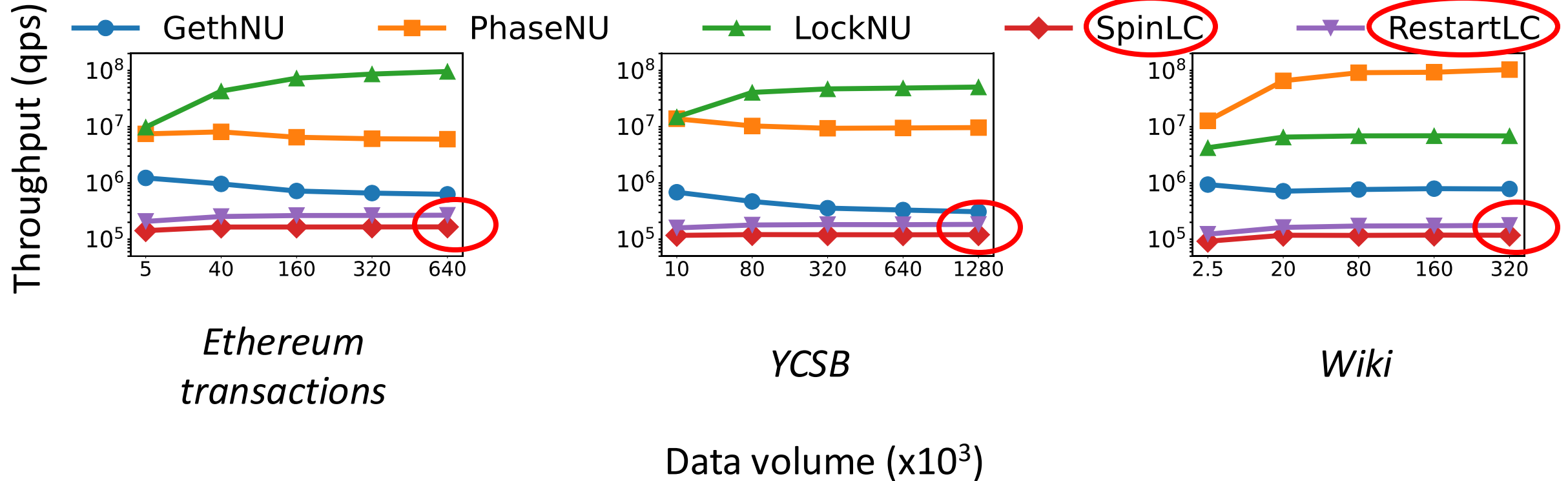
- Up to 165x speed up.
- LockNU scales better with long keys.



LockNU or PhaseNU? See our paper for the data-driven decision model!

Index Benchmark: Node Update

- Traditional lock coupling methods do not scale in GPU.

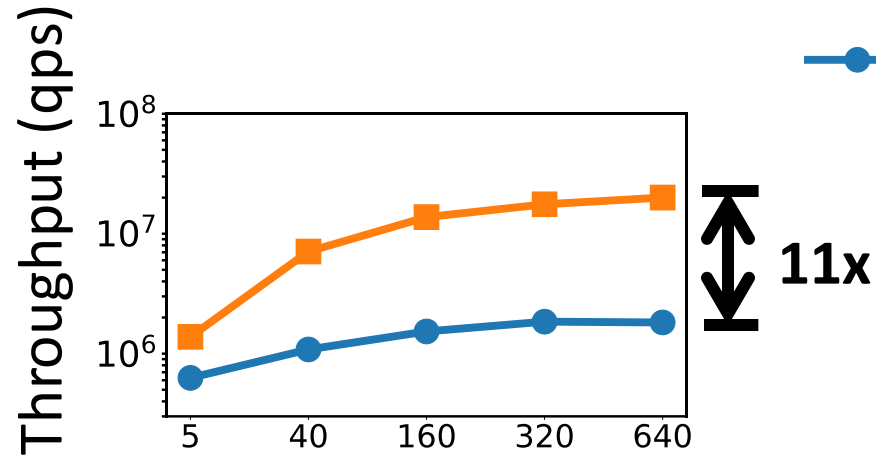




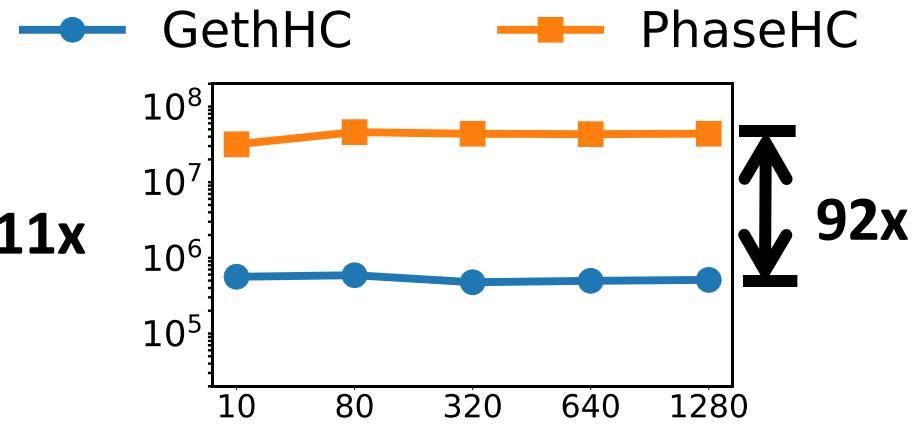
Index Benchmark: Hash Compute



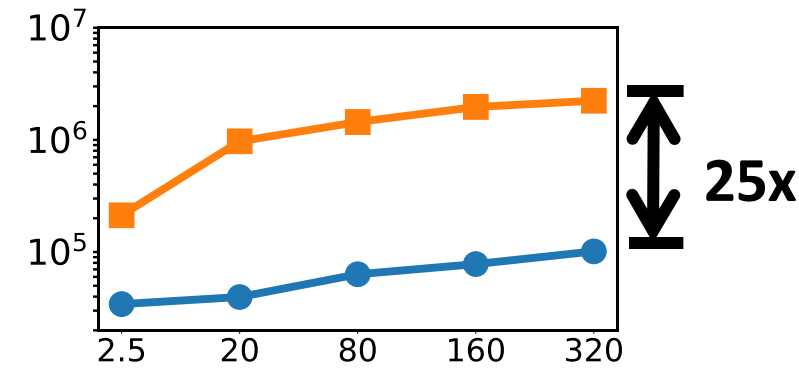
- Up to 92x speedup.



Ethereum transactions



YCSB



Wiki

Data volume (x10³)

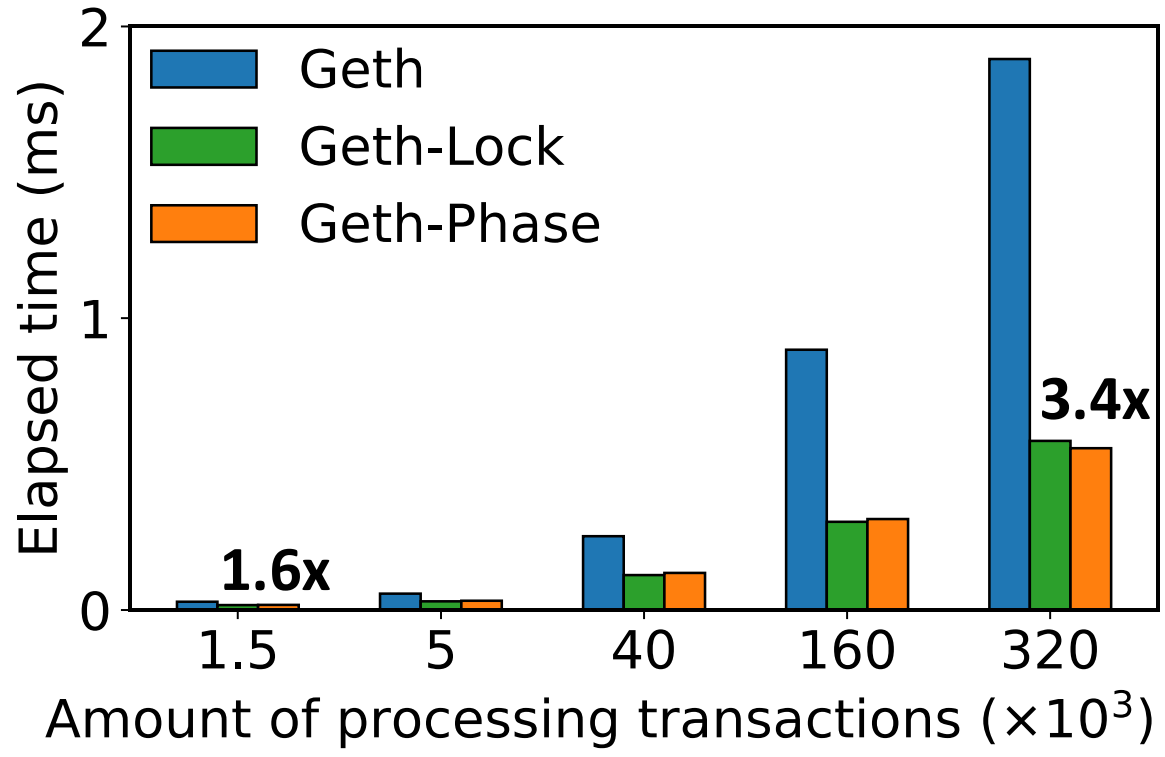


Integrating Into Real Systems



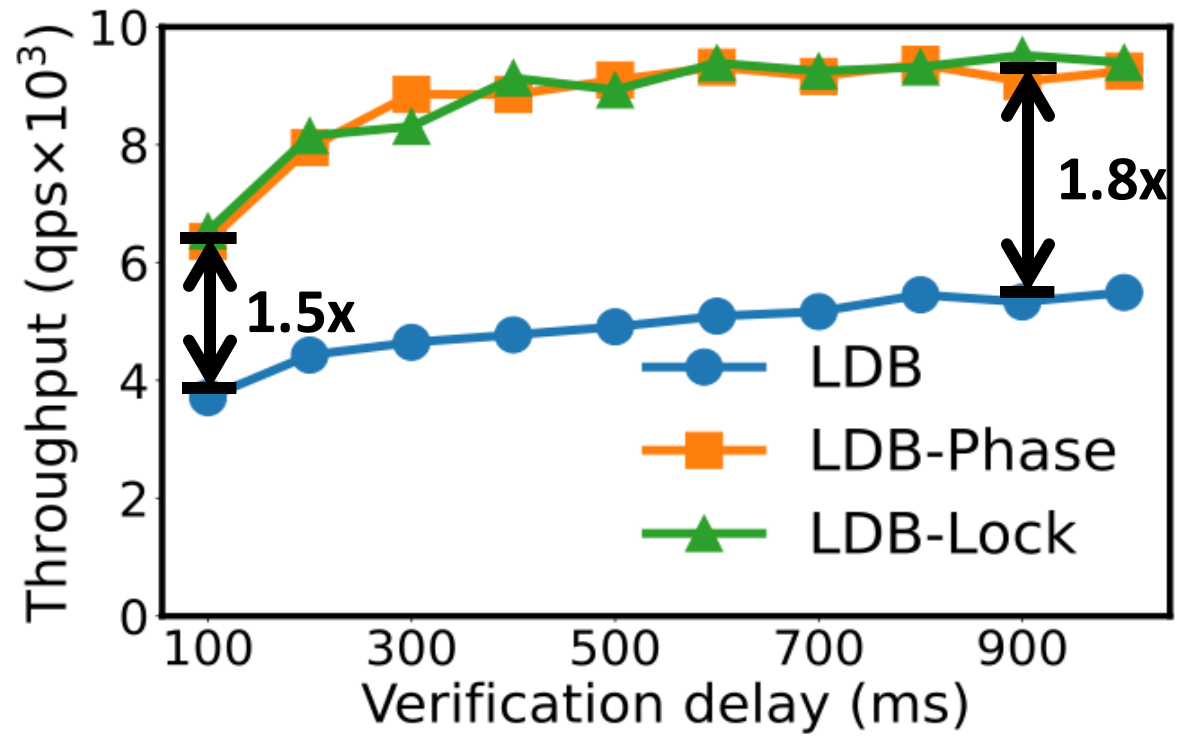
Go-Ethereum

Time to generate a block



LedgerDB

End-to-end throughput

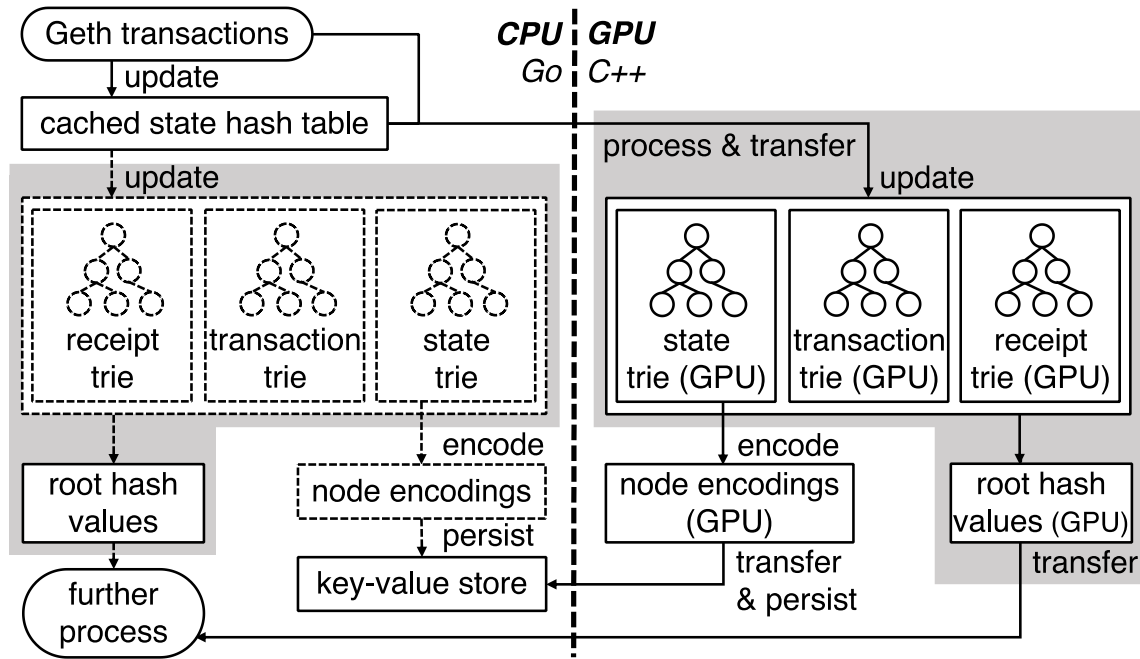




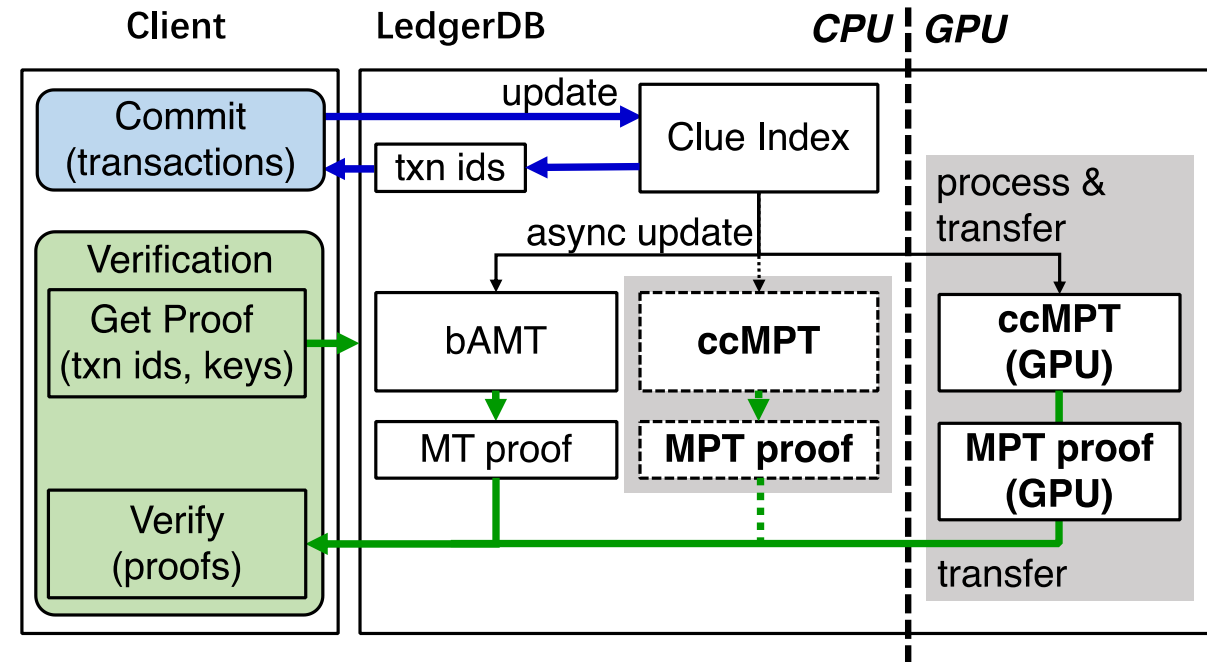
Integrating Into Real Systems



Go-Ethereum



LedgerDB





Overview



- Motivation
- Challenges
- Solution
- Results
- **Takeaways**



Takeaways

- Updating MPT on the GPU is **100+** times faster than on the CPU.
- Offloading MPT to GPU can improve the system's throughput.
- Lock-base and lock-free methods are suitable for different input distributions.



Thank you

Email: dbggroup@sustech.edu.cn

Code: <https://github.com/DBGGroup-SUSTech/GPUDB-Prefetch>