

Today's Agenda

Recap

T-Tree

Versioned Latch Coupling

Latch-Free Bw-Tree

Conclusion

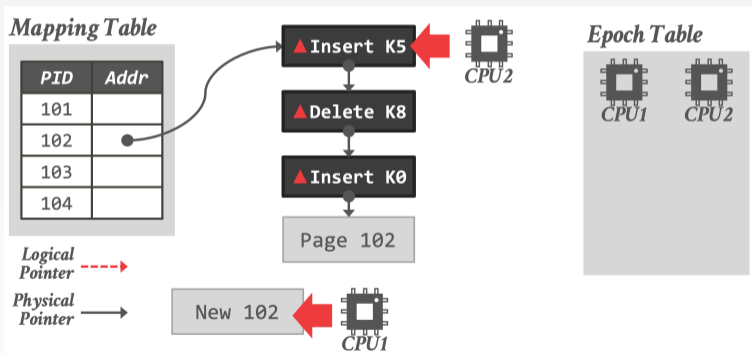
Epoch-based Garbage Collection

- Maintain a **global epoch counter** that is periodically updated (*e.g.*, every 10 ms).
 - ▶ Keep track of what threads enter the index during an epoch and when they leave.
- Mark the current epoch of a node when it is marked for deletion.
 - ▶ The node can be reclaimed once all threads have left that epoch (and all preceding epochs).
- *a.k.a.*, **Read-Copy-Update (RCU)** in Linux.

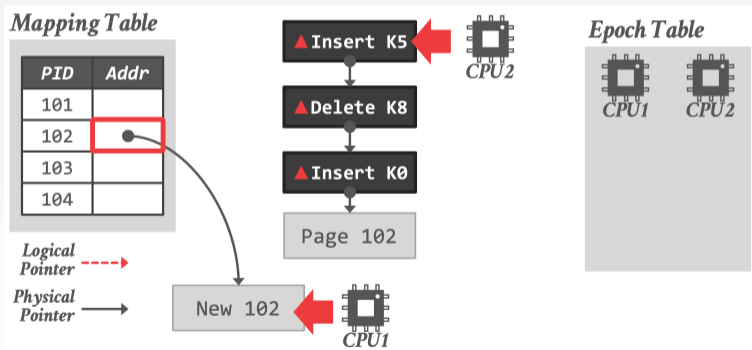
Bw-Tree: Epoch-based Garbage Collection

- Operations are tagged with an epoch number
- Each epoch tracks the threads that are part of it and the objects that can be reclaimed.
- Thread joins an epoch prior to each operation
- Garbage for an epoch reclaimed only when all threads have exited the epoch.

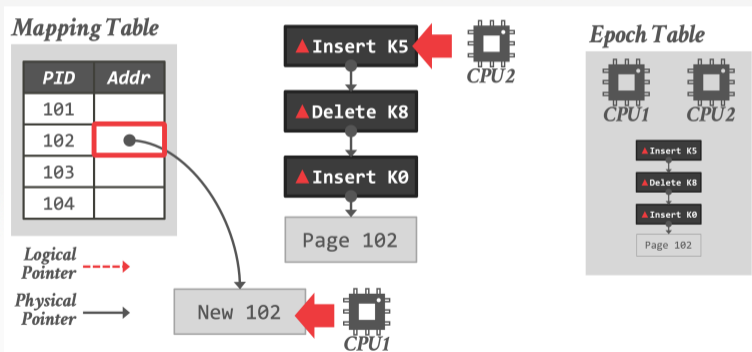
Bw-Tree: Epoch-based Garbage Collection



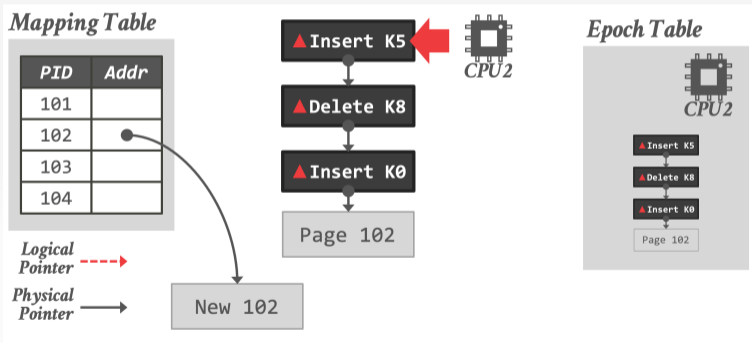
Bw-Tree: Epoch-based Garbage Collection



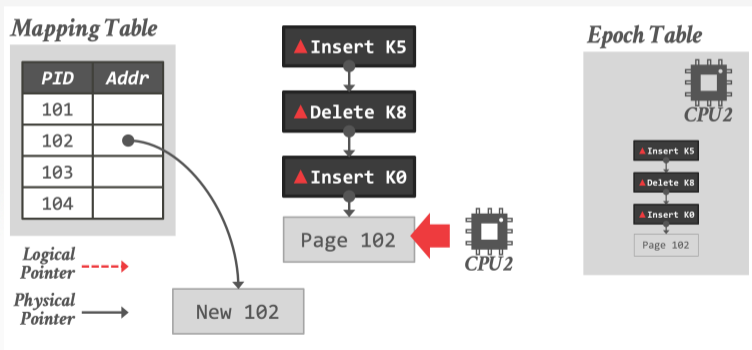
Bw-Tree: Epoch-based Garbage Collection



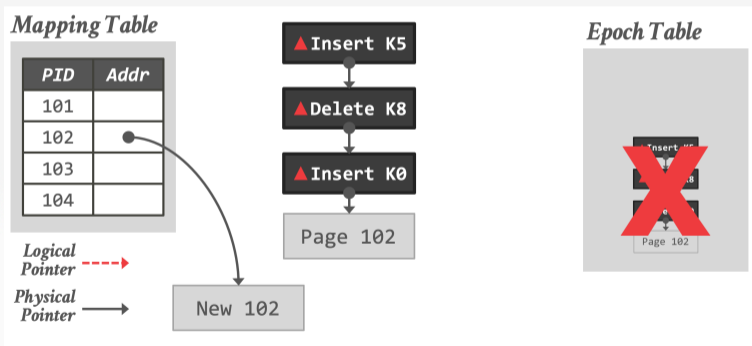
Bw-Tree: Epoch-based Garbage Collection



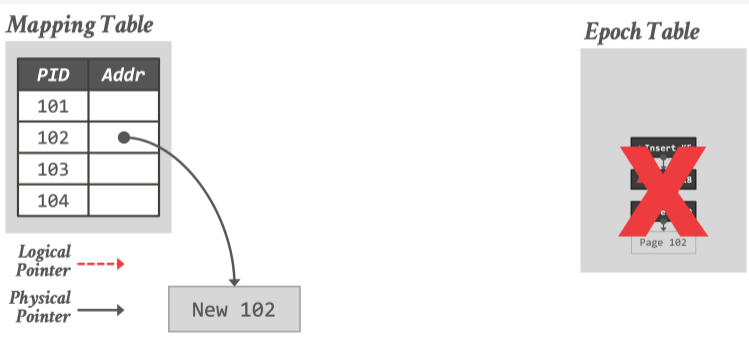
Bw-Tree: Epoch-based Garbage Collection



Bw-Tree: Epoch-based Garbage Collection



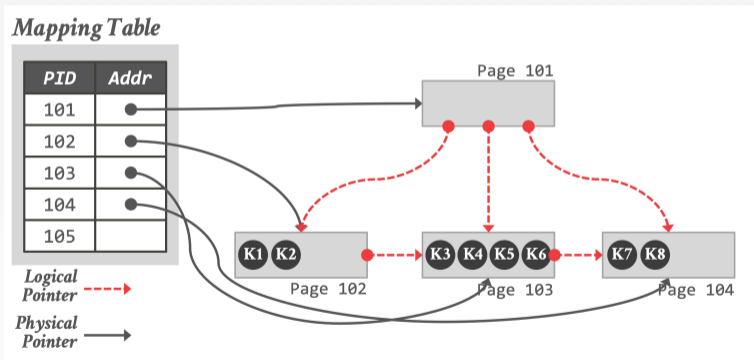
Bw-Tree: Epoch-based Garbage Collection



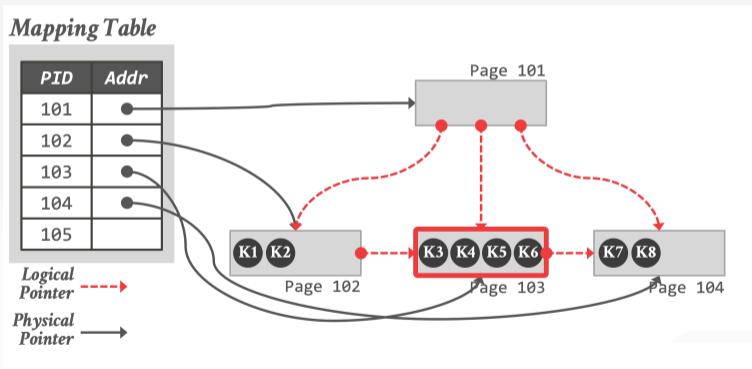
Bw-Tree: Structure Modification Operations

- **Split Delta Record**
 - ▶ Mark that a subset of the base page's key range is now located at another page.
 - ▶ Use a logical pointer to the new page.
- **Separator Delta Record**
 - ▶ Provide a shortcut in the modified page's parent on what ranges to find the new page.

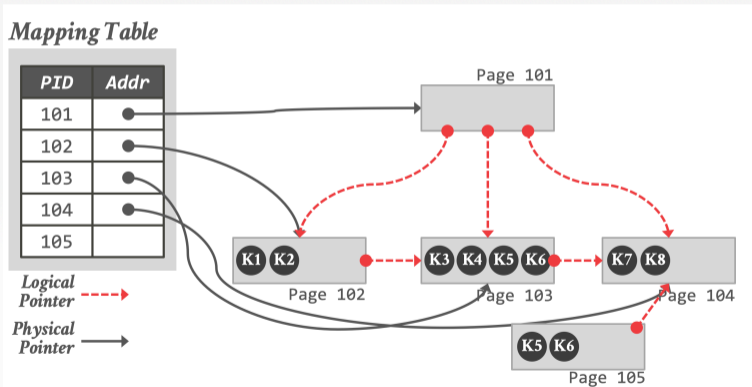
Bw-Tree: Structure Modification Operations



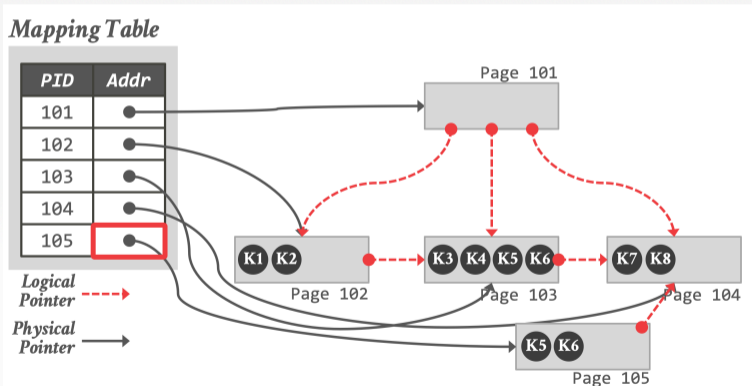
Bw-Tree: Structure Modification Operations



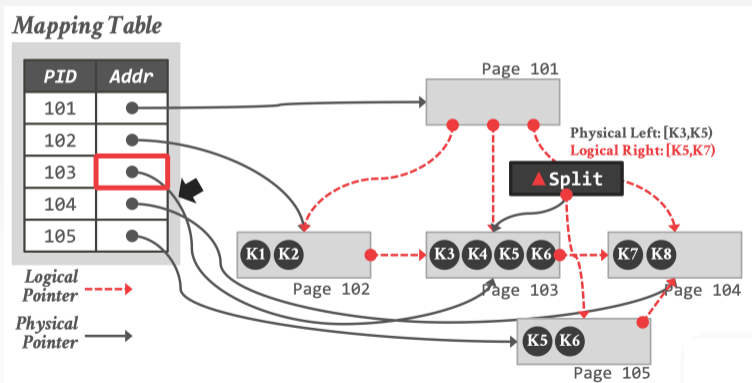
Bw-Tree: Structure Modification Operations



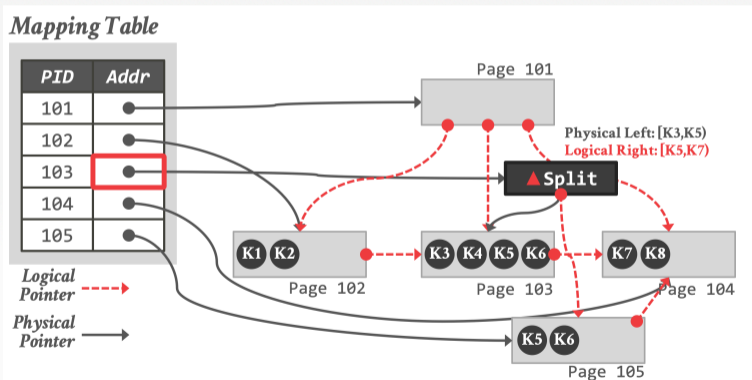
Bw-Tree: Structure Modification Operations



Bw-Tree: Structure Modification Operations



Bw-Tree: Structure Modification Operations



Conclusion

- Managing a concurrent index looks a lot like managing a database.
- Versioning and garbage collection are widely used mechanisms for increasing concurrency.
- BwTree illustrates how to design complex, latch-free data structures with only CaS instruction.