

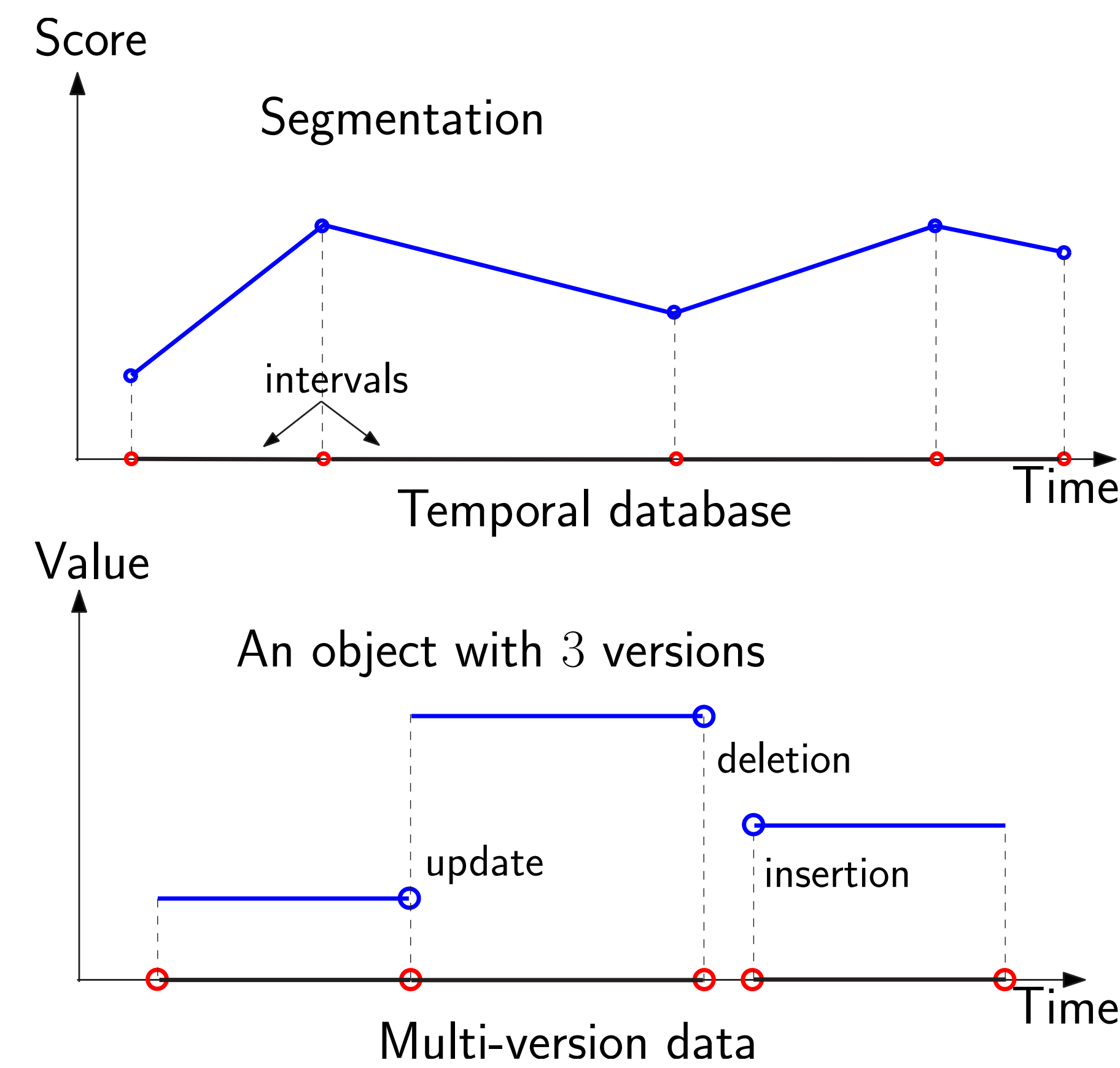
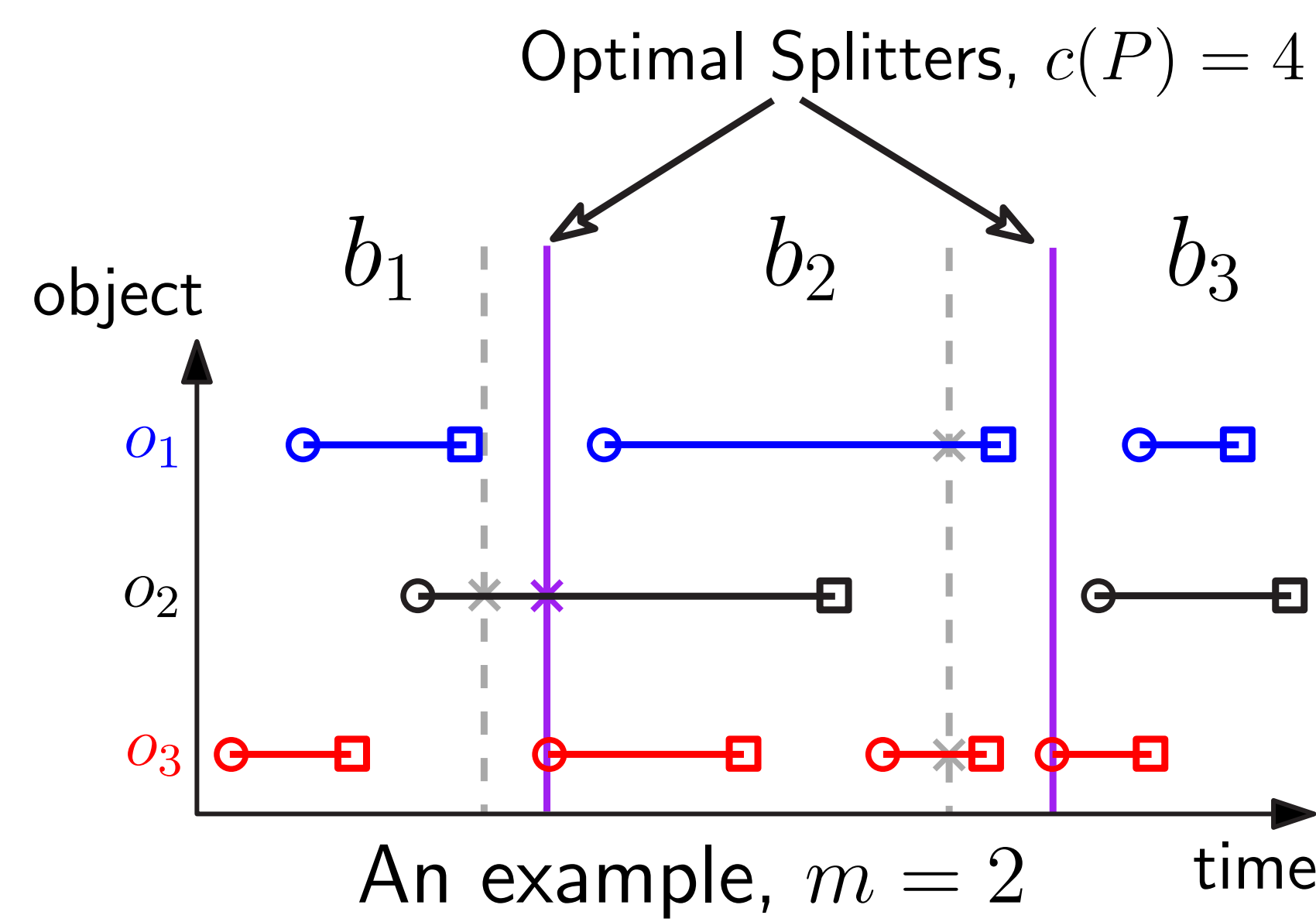
# Optimal Splitters for Temporal and Multi-version Database

Wangchao Le, Feifei Li, Yufei Tao, Robert Christensen

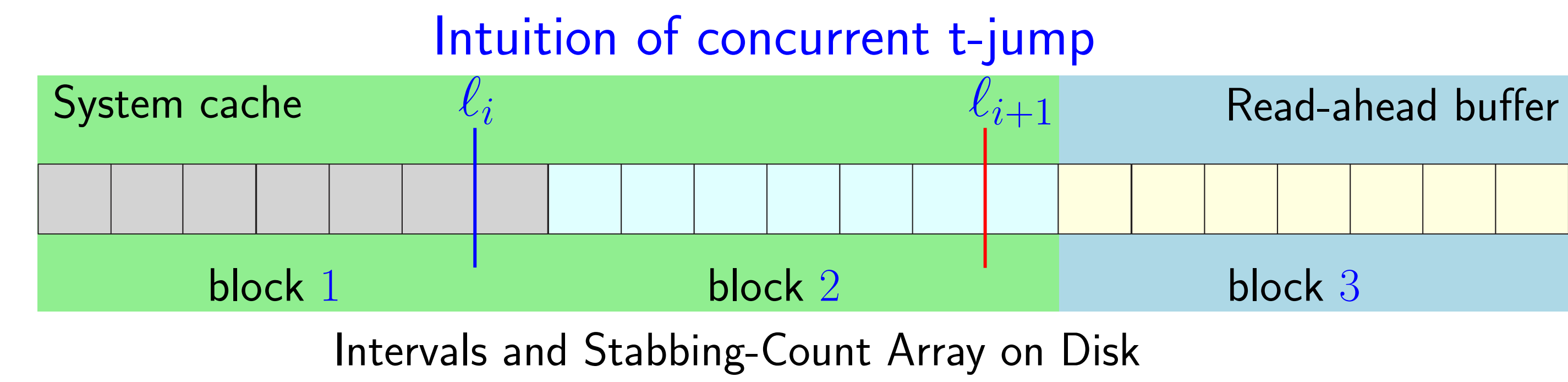
## Motivation and Problem Formulation

Temporal and multi-version data used extensively in:

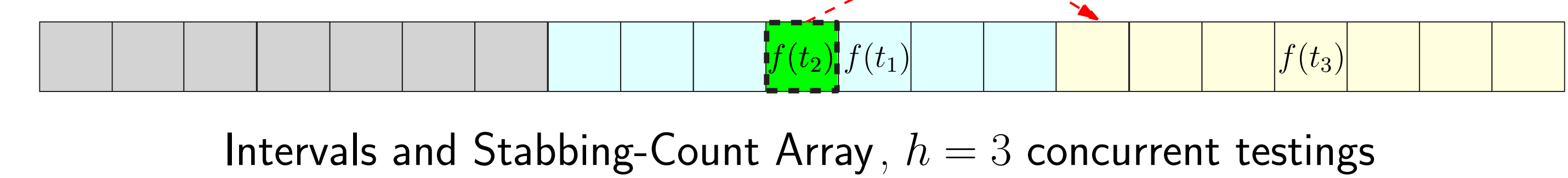
- financial market
- scientific applications
- data warehousing



## External Memory Method



- t-jump scans **forwardly**, next block to be read is **uniquely defined**.
- One execution requires  $O(1)$  space.



- initialize  $h$  threads of cost- $t$  tests  $1 \leq t_1 < t_2 < \dots < t_h \leq N$
- $f(t_i)$  the frontier of cost- $t_i$  testing
- at any time activate the thread with  $\min(f(t_i))$

## Cost Analysis

sort disk-resident array of  $O(N/B)$  blocks  $SORT(N) = \frac{(N/B) \log_{M/B}(N/B)}{M/B}$

construct the stabbing count array  $= O(SORT(N))$  I/Os

One round of Concurrent Cost- $t$  testings  $= O(N/B)$  I/Os

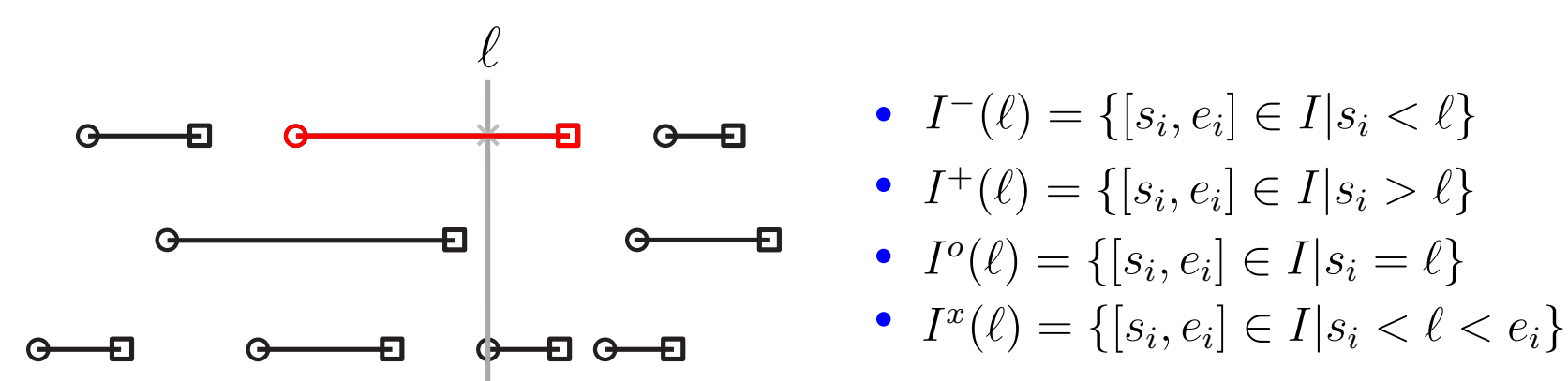
# rounds of Concurrent Cost- $t$  testings  $= O(\log_M N) \leq O(\log_{M/B} N/B)$

Cost to find  $t^*$   $= SORT(N)$  at most

Retrieve the optimal splitter  $= O(\min(k, N/B))$  I/Os

## Baseline Method – Dynamic Programming

Given a splitter  $\ell$  and a set of intervals  $I$  stored in an array.



$$c(P^*(I, k)) = \min_{\ell_k \in S(I)} \{ \max_{\ell_k} \{ c(P^*(I^-(\ell_k), k-1), \lambda) \} \}$$

$\lambda = |I^o(\ell) + I^r(\ell) + I^+(\ell)|$

A sub-problem:  $c(P^*(I^-(\ell_k), k-1))$

• How many ways to place  $\ell_k$ ?  $\ell_k \in S(I)$

Algorithm cost:  $O(kN^2)$

## Internal Memory Method

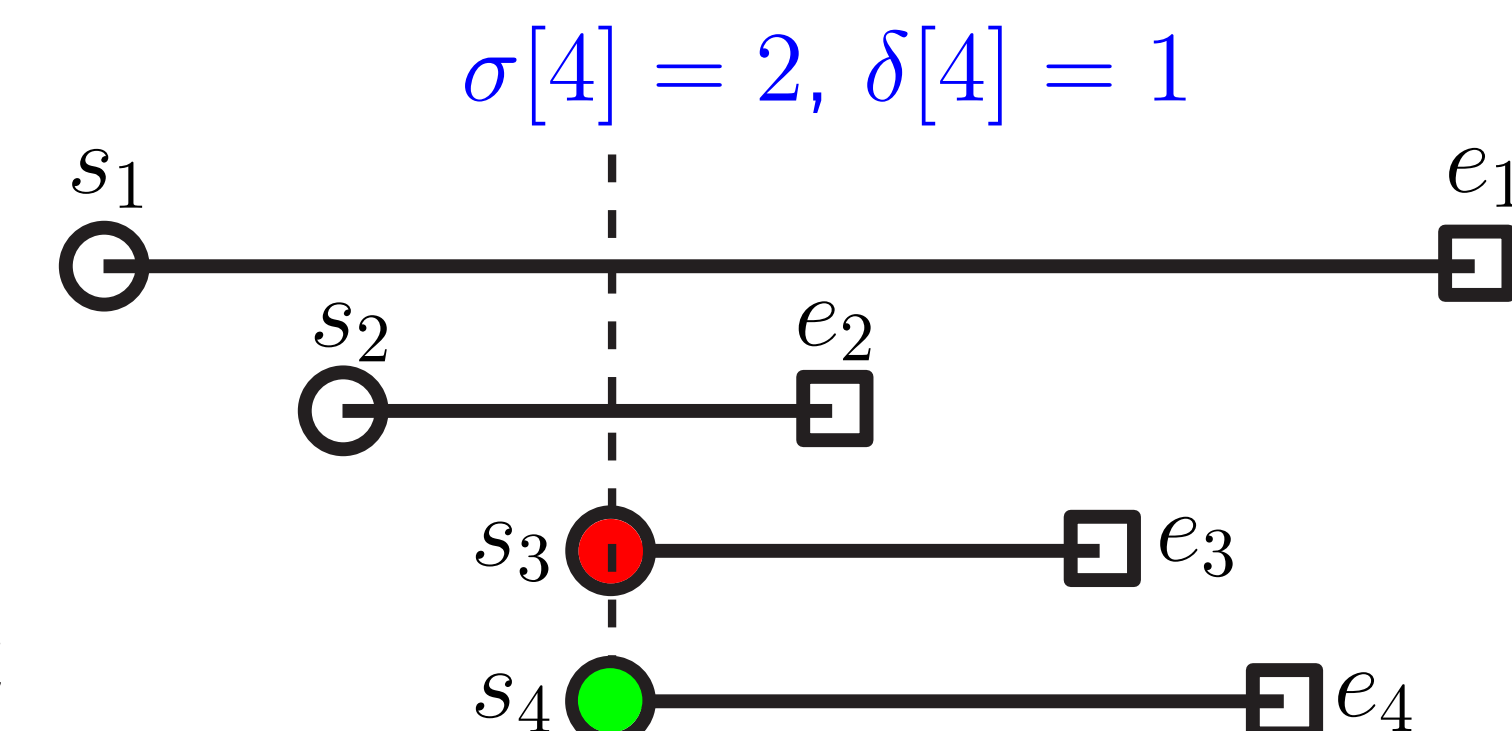
Sketch of Cost- $t$  algorithm

1. Optimal cost  $t^*$  is in range of  $\mathbf{R} = [1, N]$
2. Binary Search on  $\mathbf{R}$ 
  - if  $t$  is infeasible, then any  $t' < t$  is also infeasible
  - Solve  $O(\log N)$  instances of Cost- $t$  splitters problem
3. Report  $t^*$  as result when  $t^*$  is **feasible** but  $t^* - 1$  is **infeasible**

Algorithm cost:  
 $O(N \log N)$

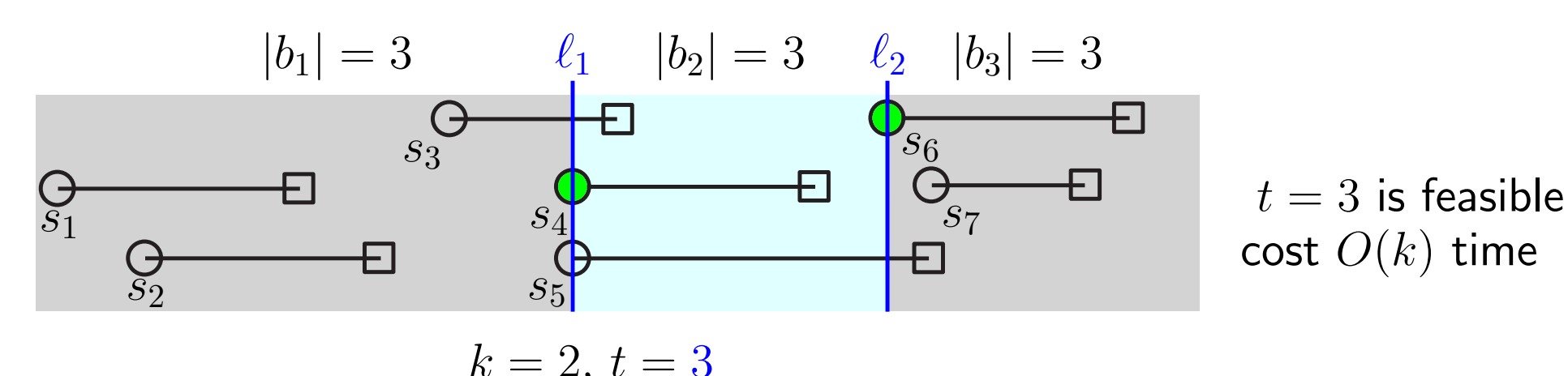
## Stabbing Count Array

- Sort  $I = \{[s_1, e_1] \dots [s_N, e_N]\}$
- The stabbing-count array for  $I$ 
  - $\forall s_i \in I$ , maintain two counts  $\sigma, \delta$ 
    - \*  $\sigma[i] = |I^x(s_i)|$ , # intervals strongly covering  $s_i$
    - \*  $\delta[i] = |I^o(s_i)|$ , # intervals in  $I^o(s_i)$  with ids less than  $i$



## t-jump method

1. place splitters in ascending order
2.  $l_{i+1}$  is pushed as far as possible from  $l_i$ , let each new  $b_i$  have size  $t$
3. if not achievable, move  $l_{i+1}$  backward just enough to form the new  $b_i$



## Experimental Results

	Internal	External
Data set	a subset of <i>Meme</i>	a subset of <i>Temp</i>
size	$\sim 21$ MB	$\sim 4.1$ GB
$N$	$\sim 1$ million	$\sim 200$ million
$k$	40	5000
$h$	not applicable	5

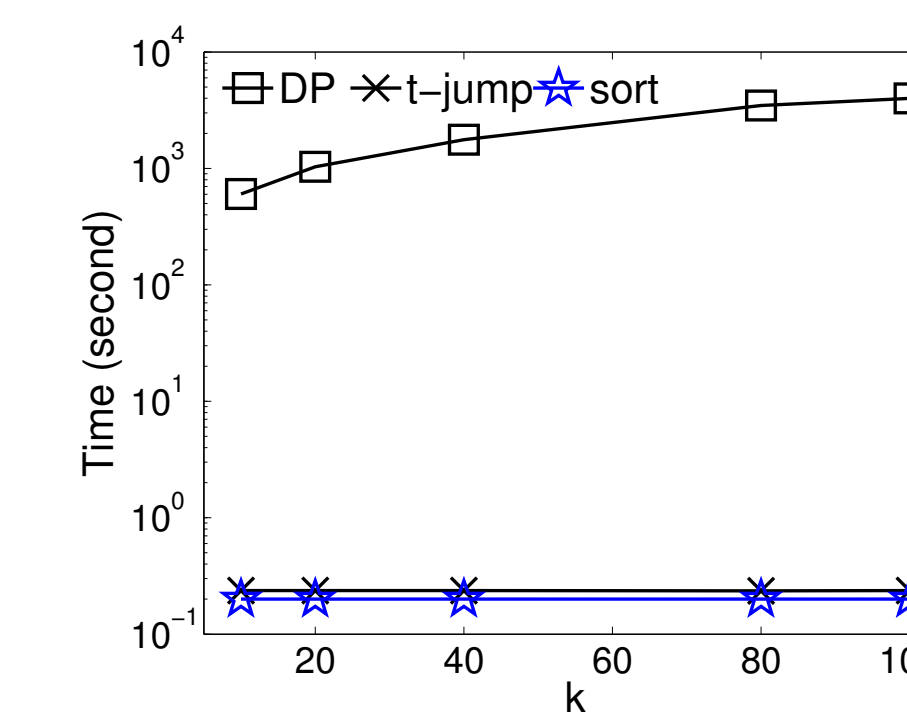


Figure 1: Running time of internal memory methods vary  $k$

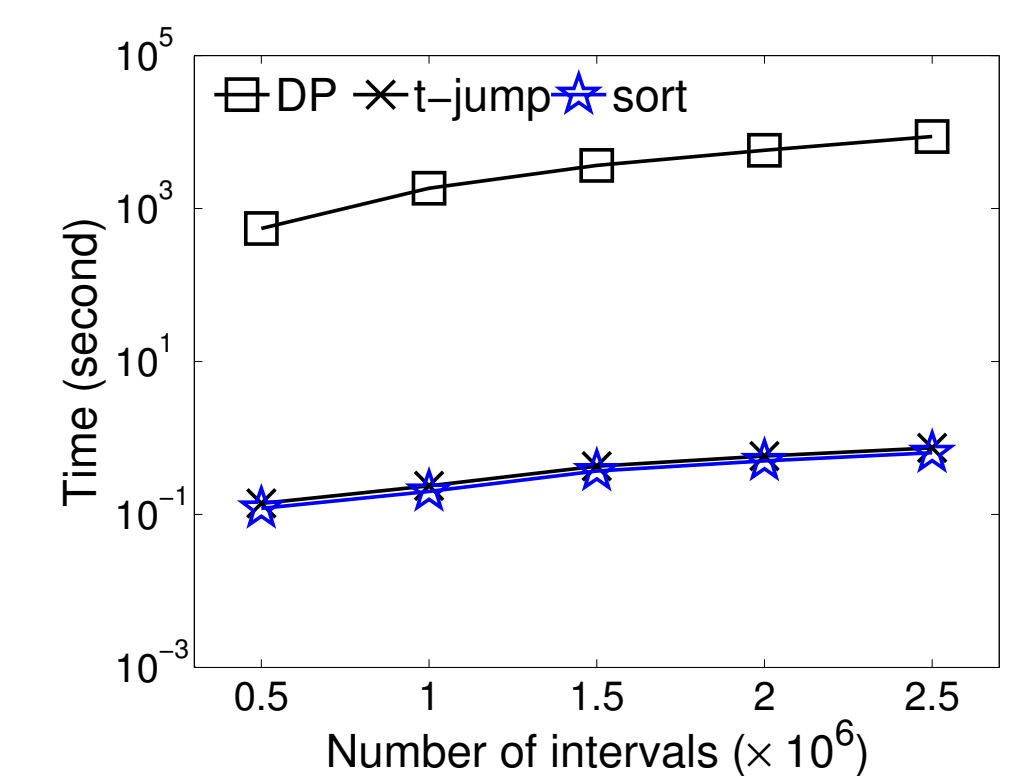


Figure 2: Running time of internal memory methods vary  $N$

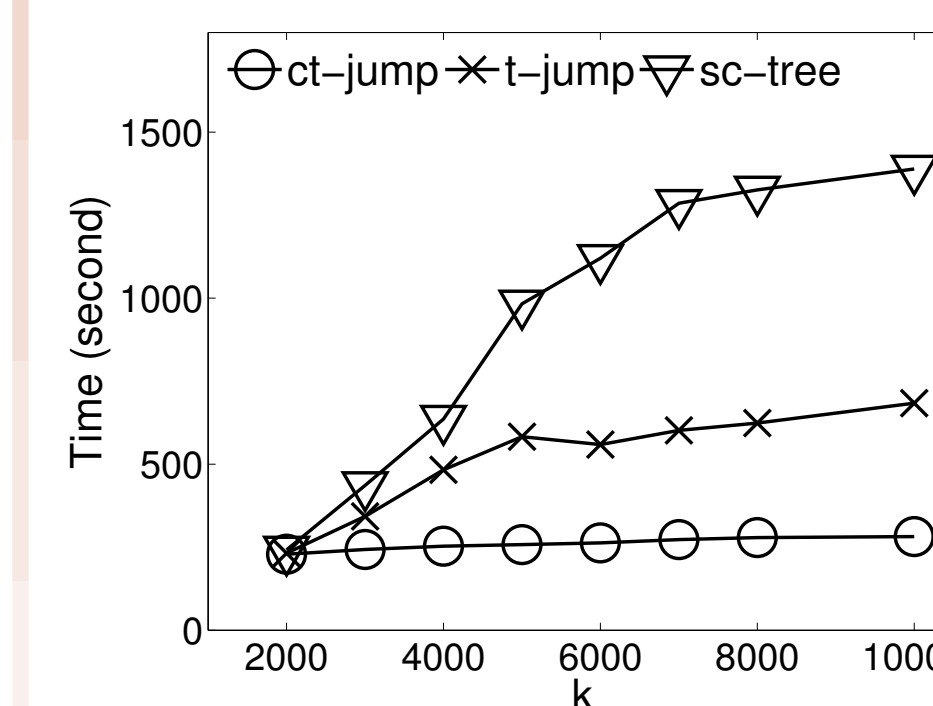


Figure 3: Query time of Queryable splitters vary  $k$

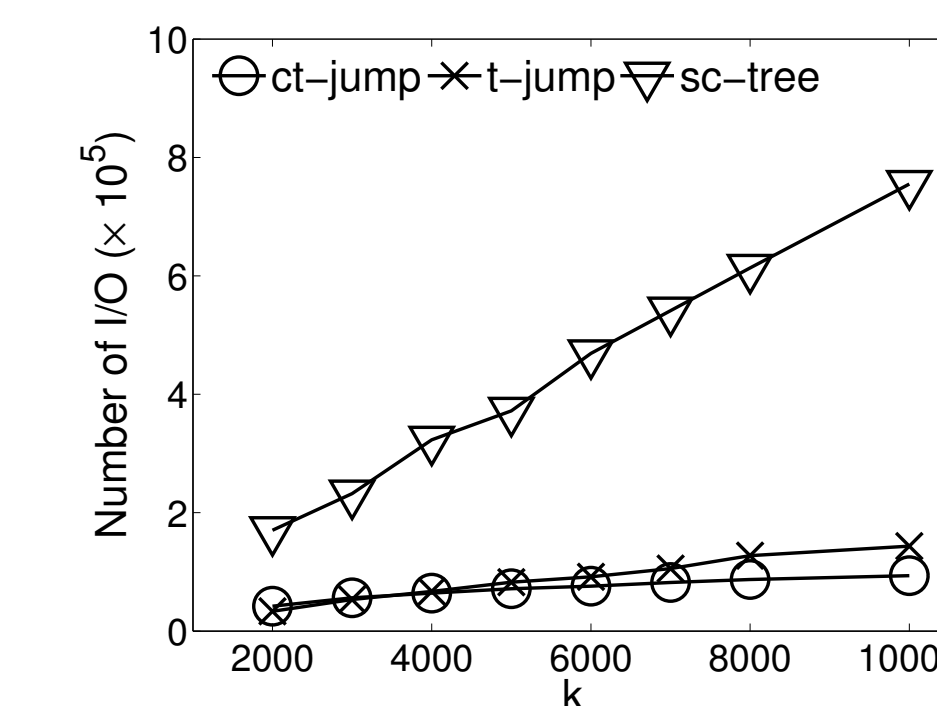


Figure 4: Query IO of Queryable splitters vary  $k$

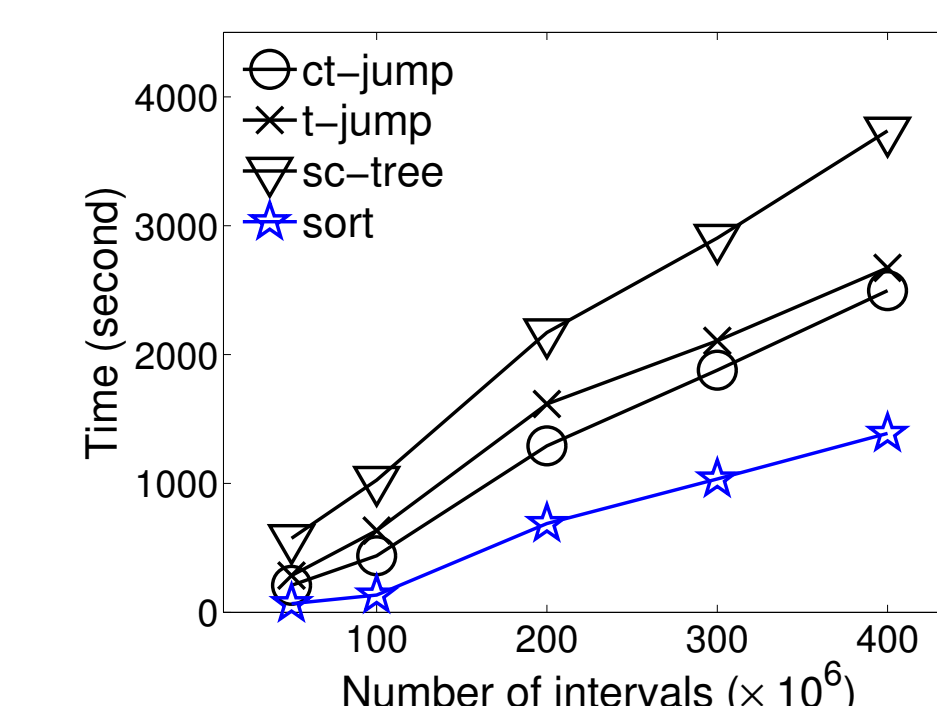


Figure 5: Query time of static splitter vary  $N$

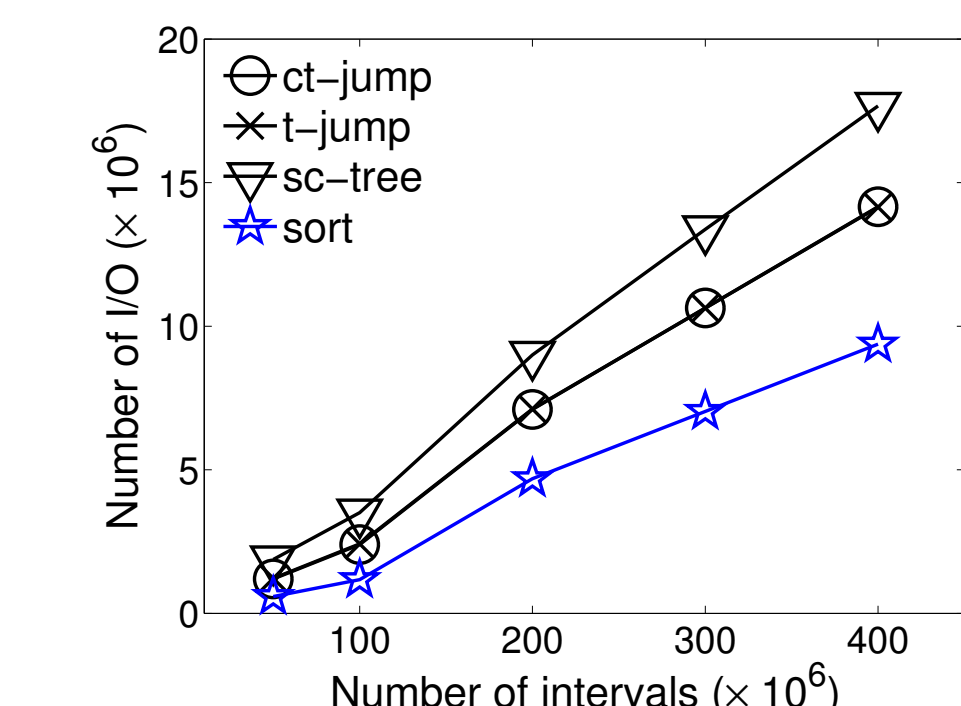


Figure 6: Total IO of static splitter vary  $N$

## Conclusion

Temporal and multi-version databases often generate massive amounts of data. Therefore, it becomes increasingly important to store and process this data in a distributed and parallel fashion. This makes an important contribution in solving the optimal splitters problem, which is essential in enabling efficient distributed and parallel processing of such data.