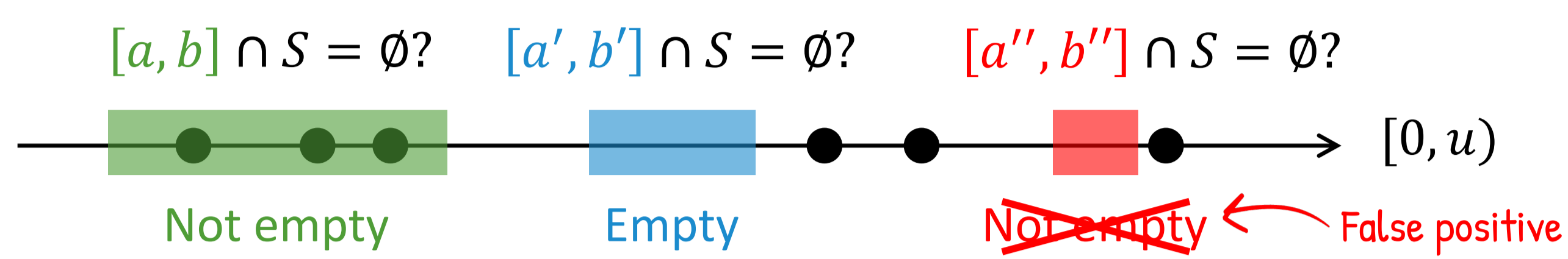


## Range Filters

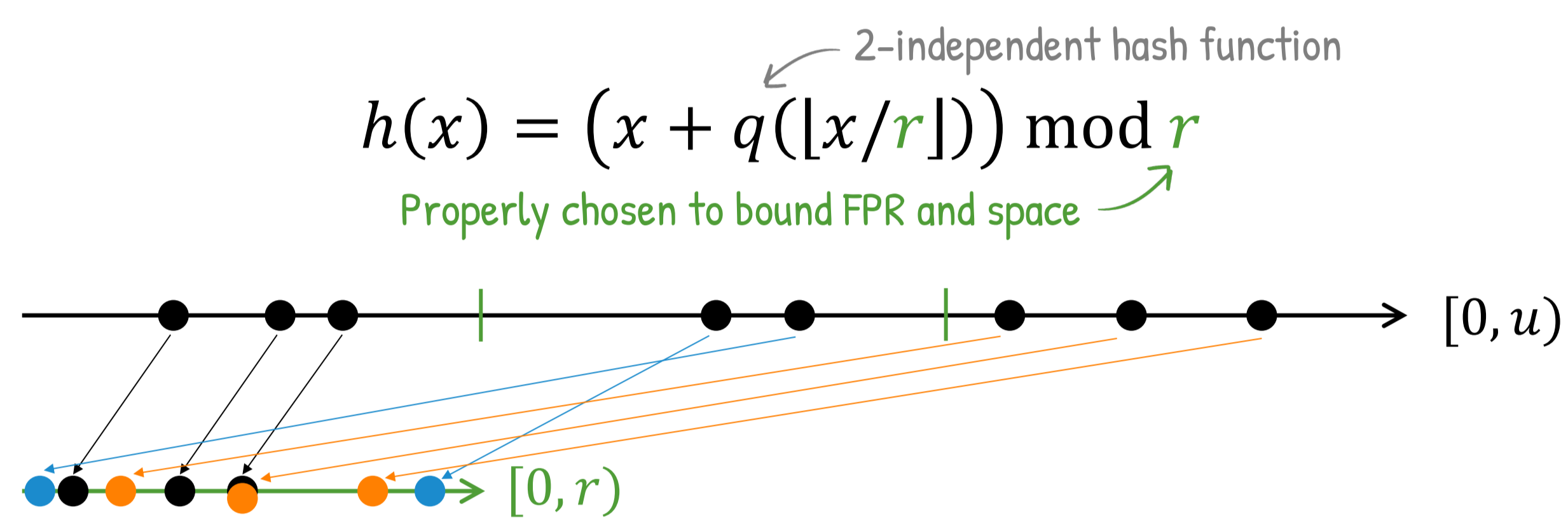
Given a set  $S$  of  $n$  keys, a range filter is a space-efficient data structure that answers range emptiness queries with a false positive probability of at most  $\varepsilon$



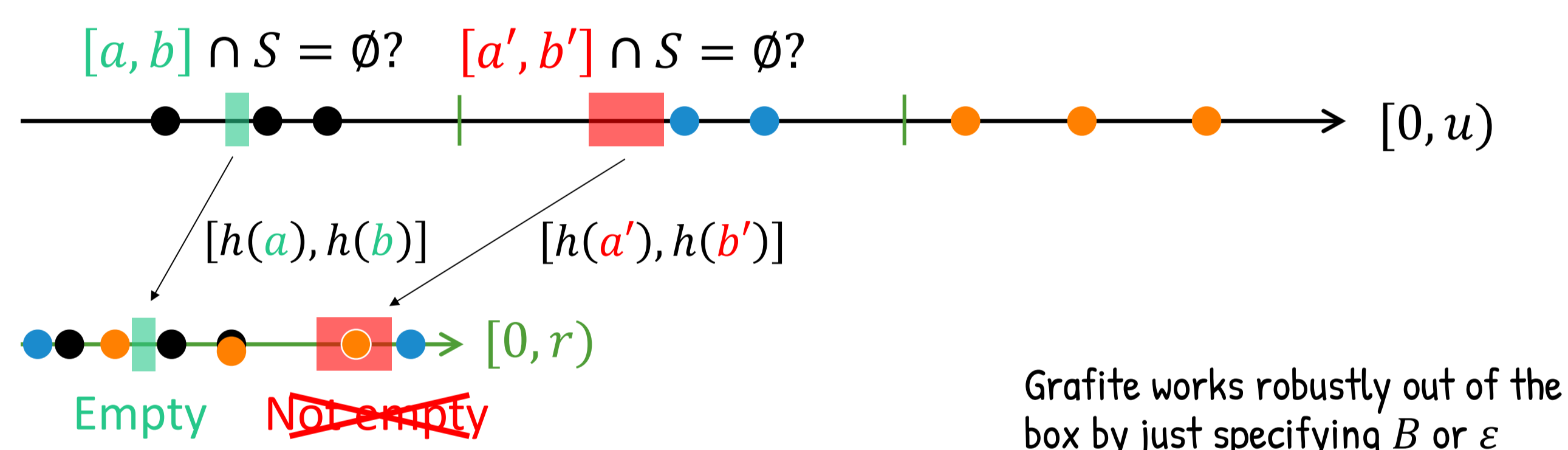
- Generalises Bloom filters from point to range queries
- Reduces I/Os of range queries in LSM-based storage engines

## Grafite: An Optimal Robust Range Filter

1. Apply a locality-preserving hash function to the keys



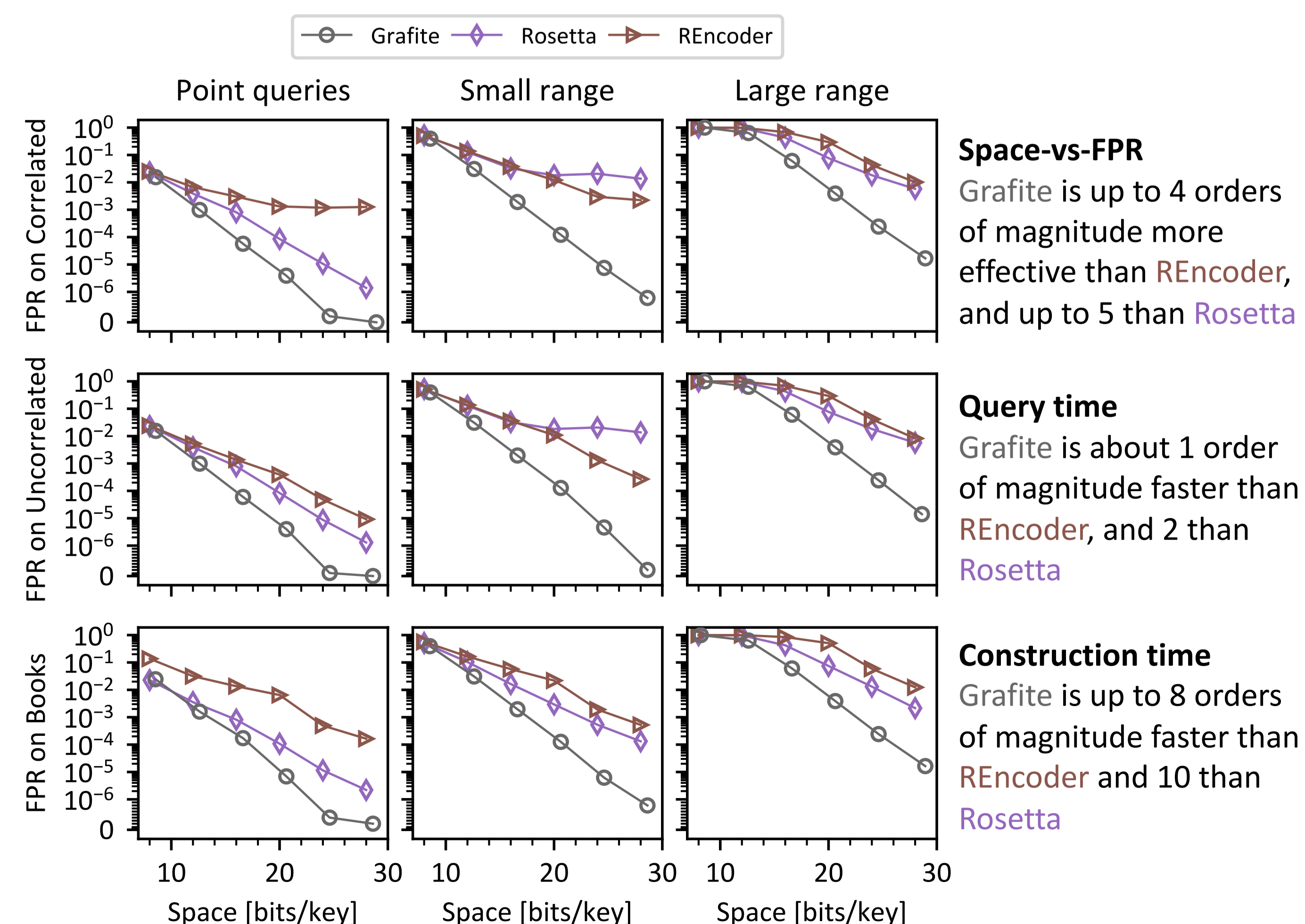
2. Compress the hash codes with the Elias-Fano code
3. Solve queries in hash space



Grafite works robustly out of the box by just specifying  $B$  or  $\varepsilon$

**Theorem.** Given a space budget of  $B$  bits/key, the query time of Grafite is  $O(1)$  and the false positive probability  $\varepsilon$  is no more than  $\ell/2^{B-2}$ , where  $\ell$  is the query range size

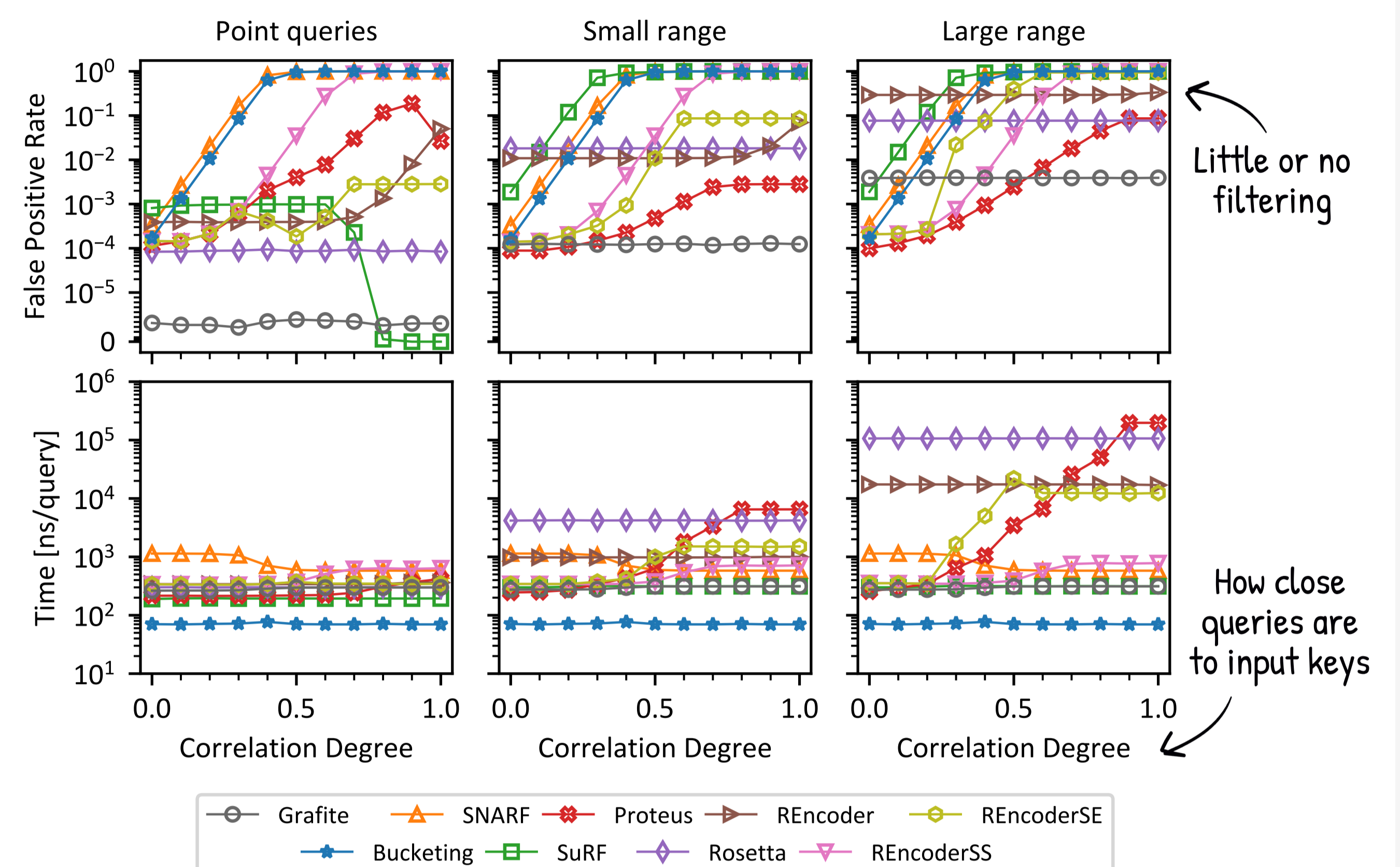
## Experiments with Robust Range Filters



**Take-home message.** If robustness guarantees are needed regardless of input data and future queries, Grafite is the range filter of choice

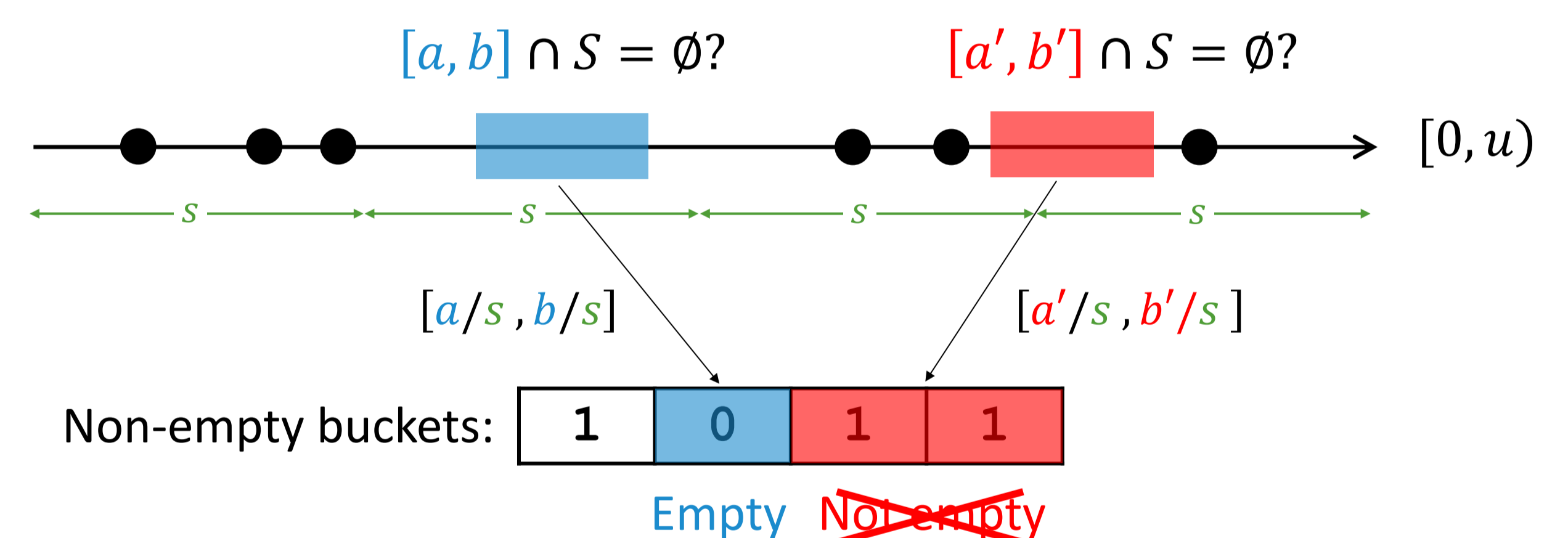
## Issues with Current Range Filters

- Highly complex**  
Sophisticated designs, hard to evaluate and deploy
- Fragile**  
Inconsistent FPR and query times across different datasets
- Adversarial queries**  
Easy to issue queries that result in false positives (thus I/Os)

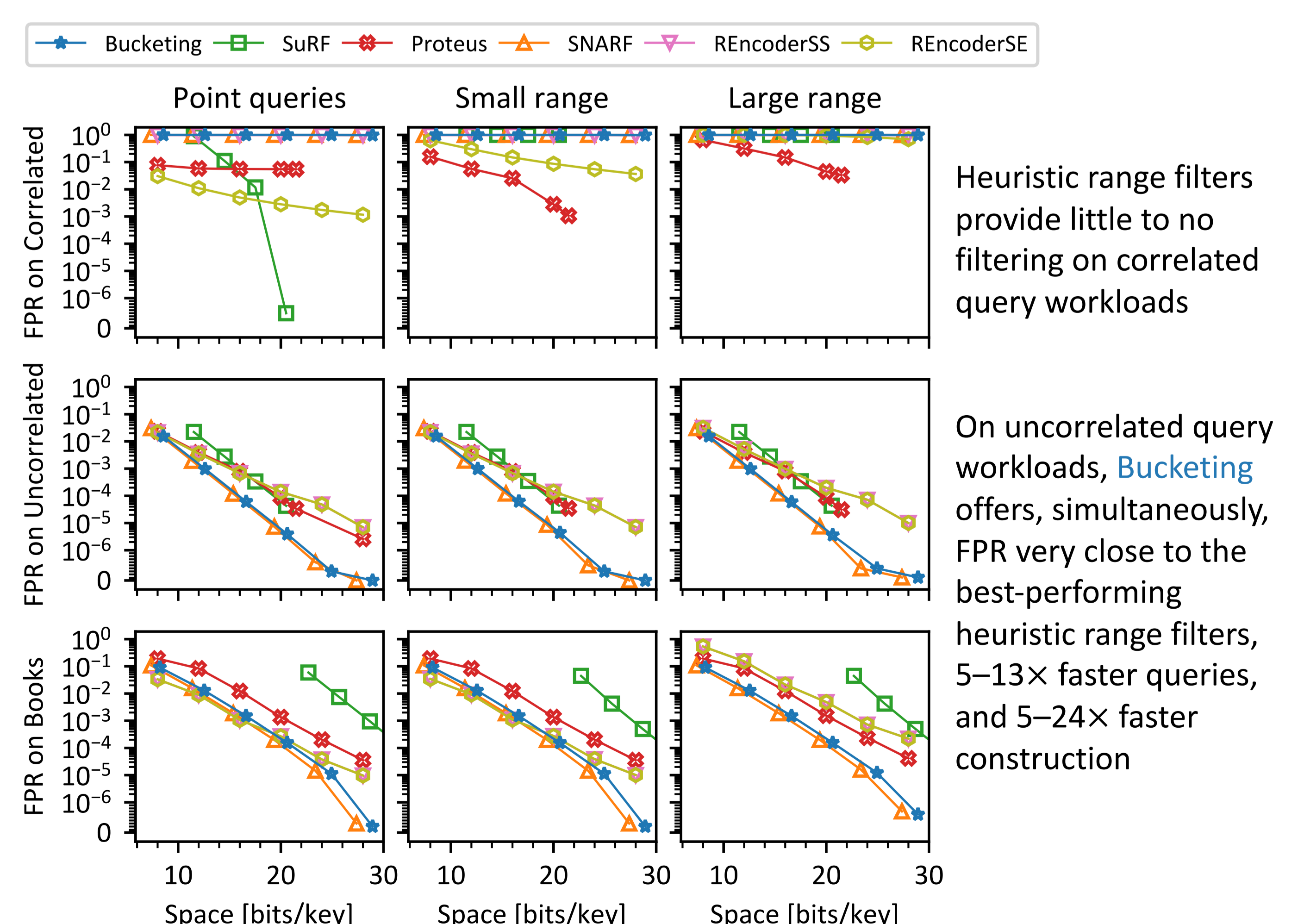


## Bucketing: A Simple Heuristic Range Filter

1. Divide universe into buckets of equal size  $s$  ← Impacts FPR and space
2. Mark non-empty buckets with a compressed bit-vector
3. Solve queries by mapping ranges to bit-vector positions



## Experiments with Heuristic Range Filters



**Take-home message.** Heuristic range filters sacrifice robustness to work well on some inputs only, but simpler solutions like Bucketing exist