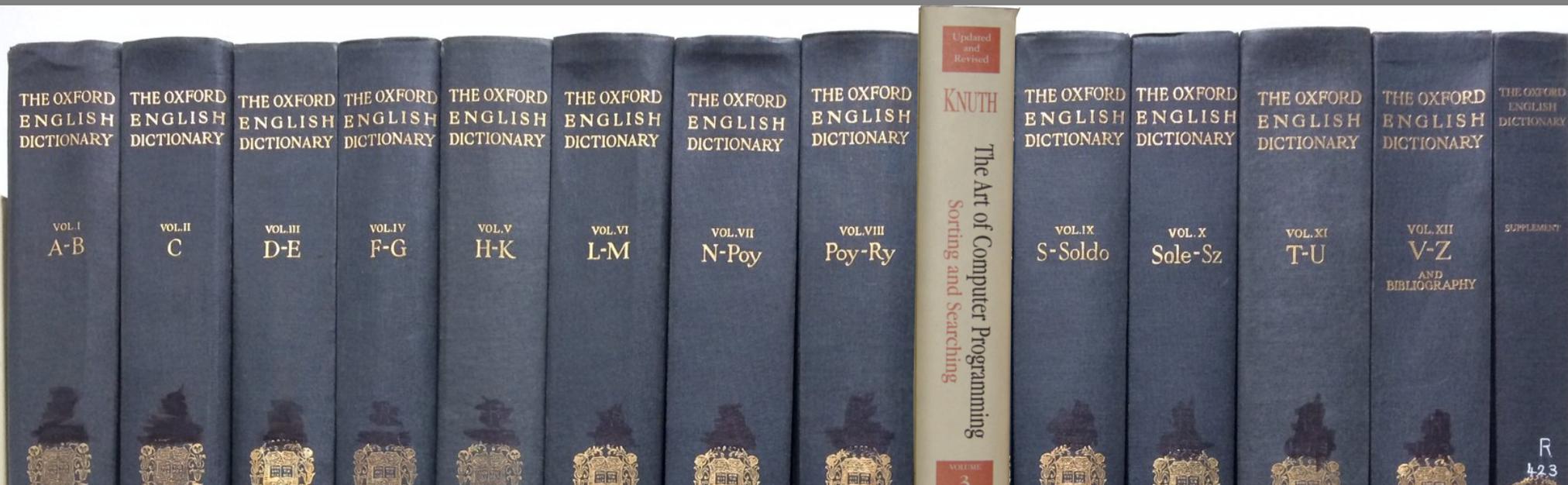


Dynamic Space Efficient Hashing

Presentation · 05. September 2017

Tobias Maier and Peter Sanders

INSTITUTE OF THEORETICAL INFORMATICS · ALGORITHMICS GROUP

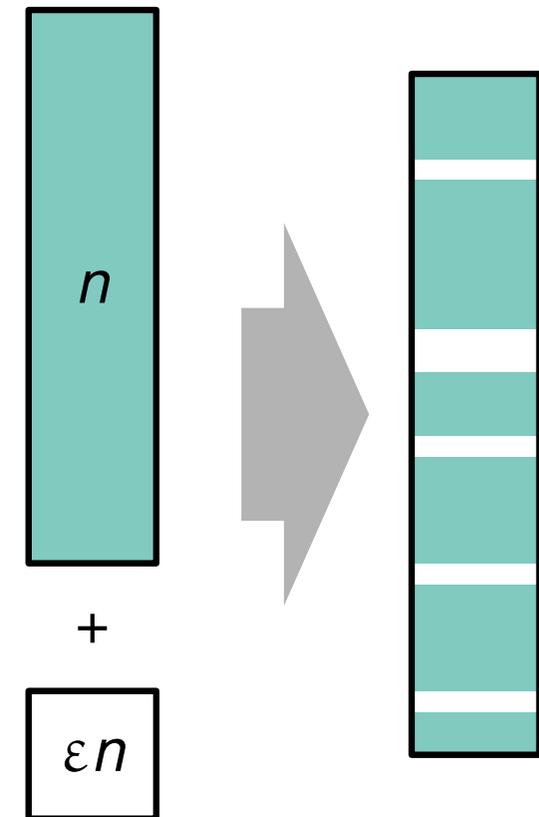


Basics – Hash Tables

- ubiquitous **dictionary** datastructure
 - ▶ insert
 - ▶ find (preferably $O(1)$)
 - ▶ erase (preferably $O(1)$)
- we do not consider chaining
 - ▶ not space efficient for small elements

“Classic” Space Efficient Hashing

- densely filled table
- lots of collisions
 - ▶ needs good collision handling
- static size (post-initialization)
 - ▶ fixed number of elements



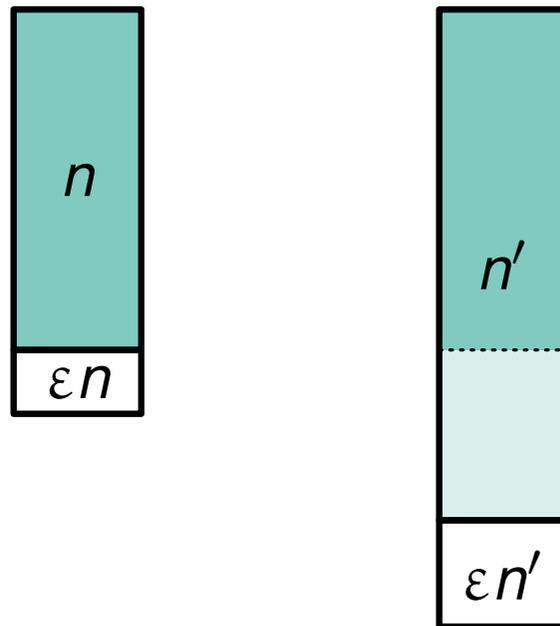
Final Size Not Known A Priori

■ conservative estimate

$$n \leq n'$$

▶ strict bound might not be reasonable

▶ less space efficient



Final Size Not Known A Priori

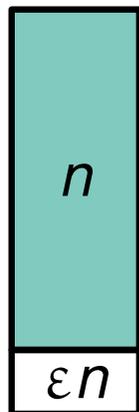
- conservative estimate

- optimistic estimate

- ▶ might overflow

- ▶ needs growing strategy

$$n \approx n'$$



slow



needs growing

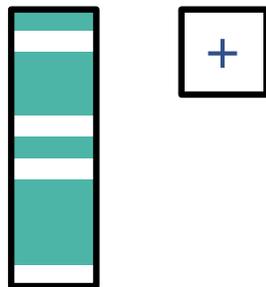
Final Size Not Known A Priori

- conservative estimate
- optimistic estimate
- number of elements changes over time
 - ▶ cannot be initialized with max size

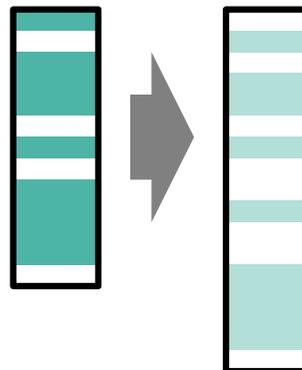
Resizing

- growing has to be in **small steps**
- basic approaches

additional table

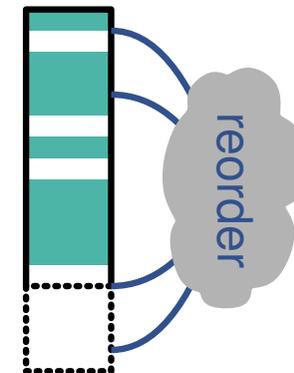


full migration



most common
in libraries

inplace+reorder

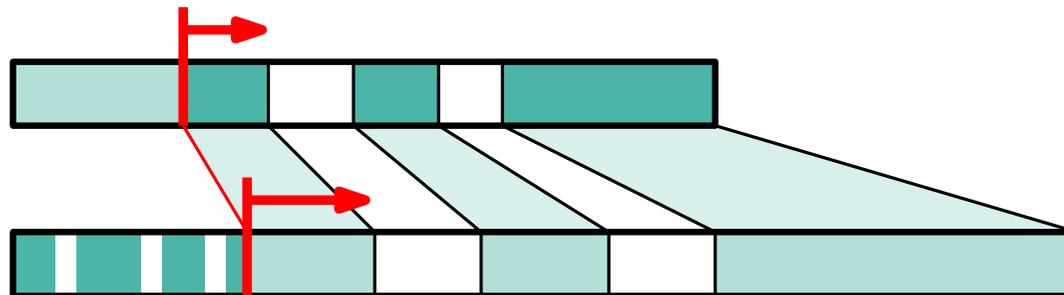


Secondary Contribution – Efficient Growing

- addressing the table (no powers of two)
 - ▶ conventional wisdom: modulo table size
 - ▶ faster: use hash value as scaling factor

$$idx(k) = h(k) \cdot \frac{size}{maxHash + 1}$$

- very fast migration due to cache efficiency

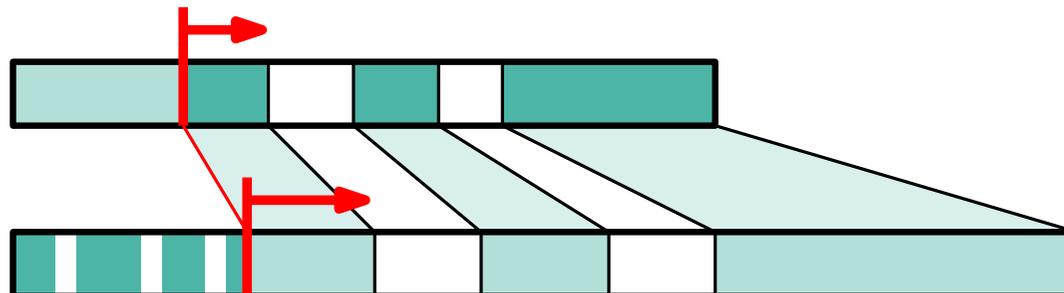


Secondary Contribution – Efficient Growing

- addressing the table (no powers of two)
 - ▶ ~~conventional wisdom: modulo table size~~
 - ▶ faster: use hash value as **scaling** factor

$$\text{idx}(k) = h(k) \cdot \frac{\text{size}}{\text{maxHash} + 1}$$

- very fast migration due to cache efficiency

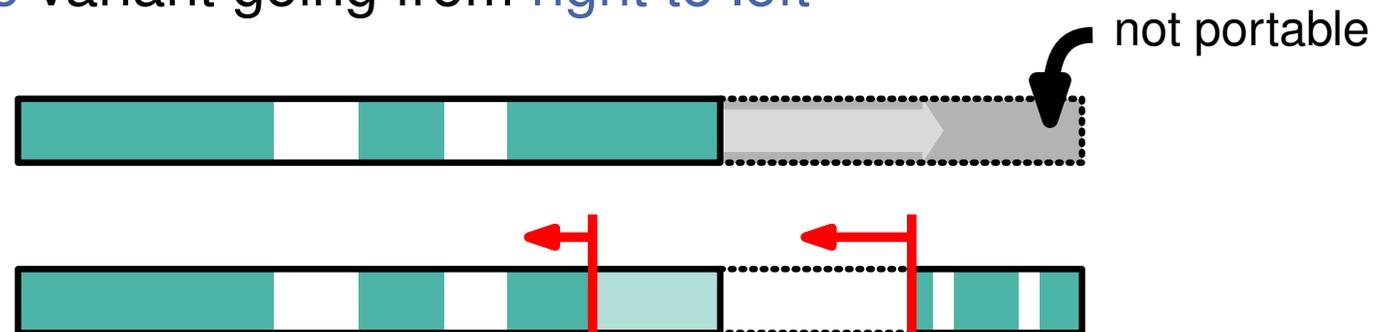


Secondary Contribution – Efficient Growing

- addressing the table (no powers of two)
 - ▶ ~~conventional wisdom: modulo table size~~
 - ▶ faster: use hash value as **scaling** factor

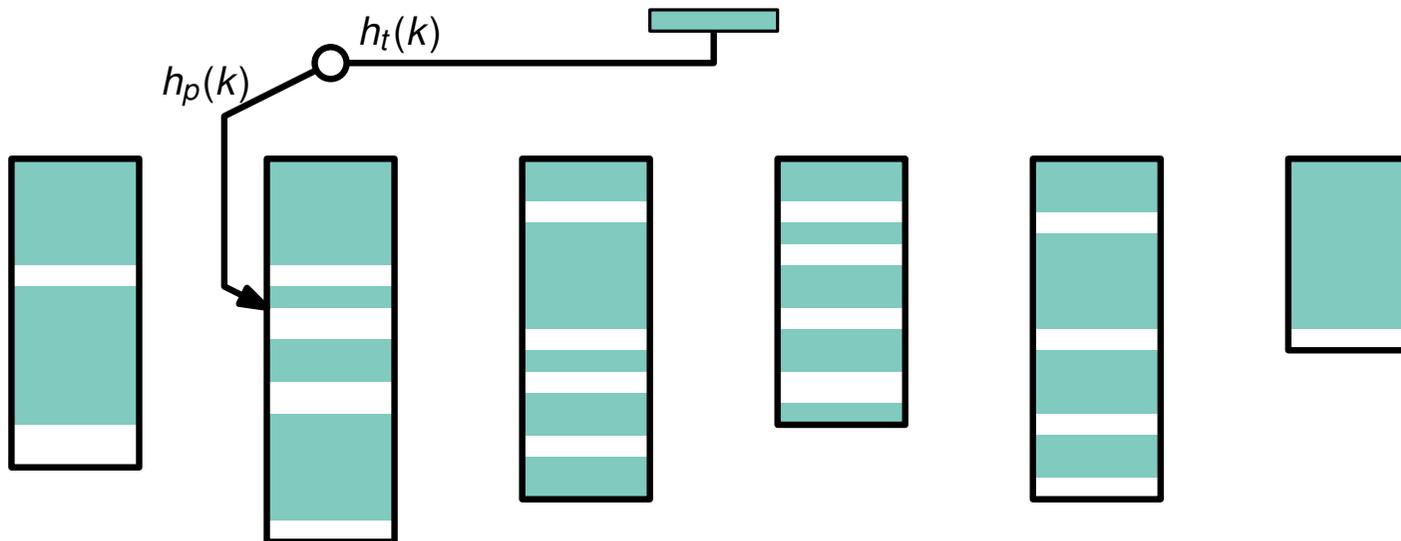
$$\text{idx}(k) = h(k) \cdot \frac{\text{size}}{\text{maxHash} + 1}$$

- very fast migration due to cache efficiency
- **inplace** variant going from **right to left**



Multi Table Approach

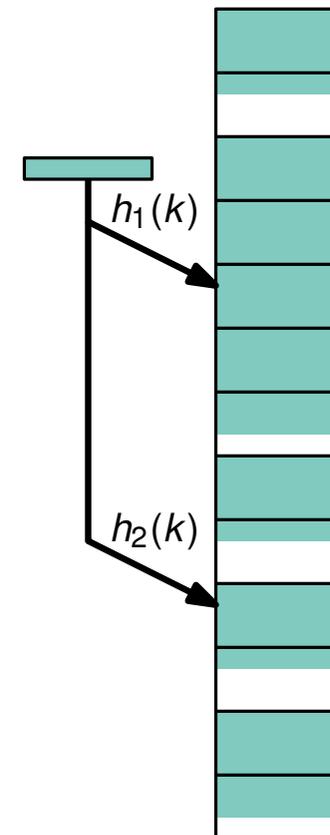
- $T = 2^c$ subtables with expected equal count
 - ▶ reduces memory during subtable migration
- $h(k) \Rightarrow h_t(k)$ for the subtable $h_p(k)$ within the table



H -ary B -Bucket Cuckoo Hashing

[Pagh, Dietzfelbinger, Mehlhorn, Mitzenmacher, ...]

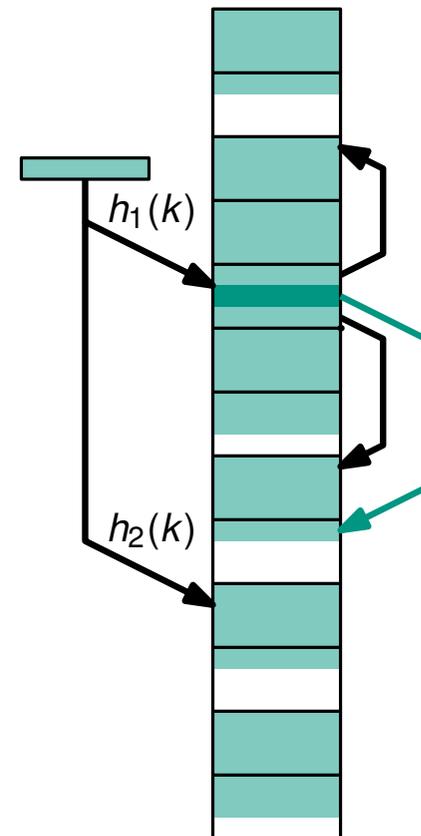
- buckets of B cells
- H alternative buckets per element
 $h_1(k), \dots, h_H(k)$
- if buckets are full, move existing elements
 - breadth-first-search



H -ary B -Bucket Cuckoo Hashing

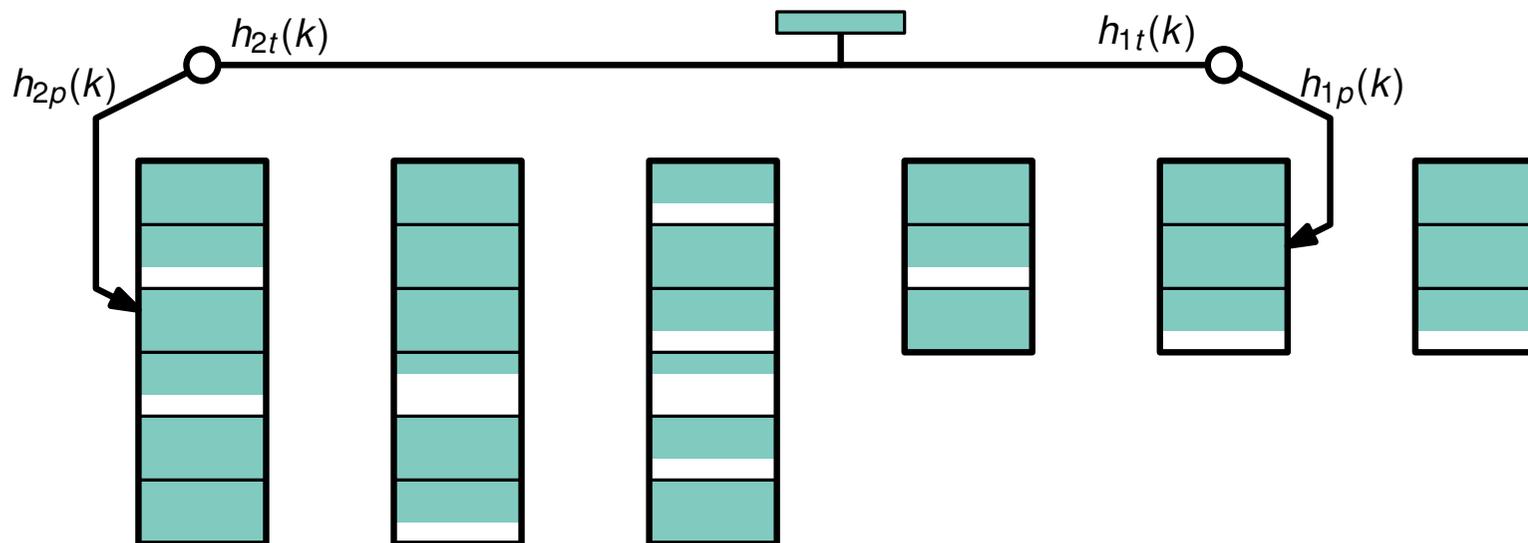
[Pagh, Dietzfelbinger, Mehlhorn, Mitzenmacher, ...]

- buckets of B cells
- H alternative buckets per element
 $h_1(k), \dots, h_H(k)$
- if buckets are full, move existing elements
 - ▶ breadth-first-search



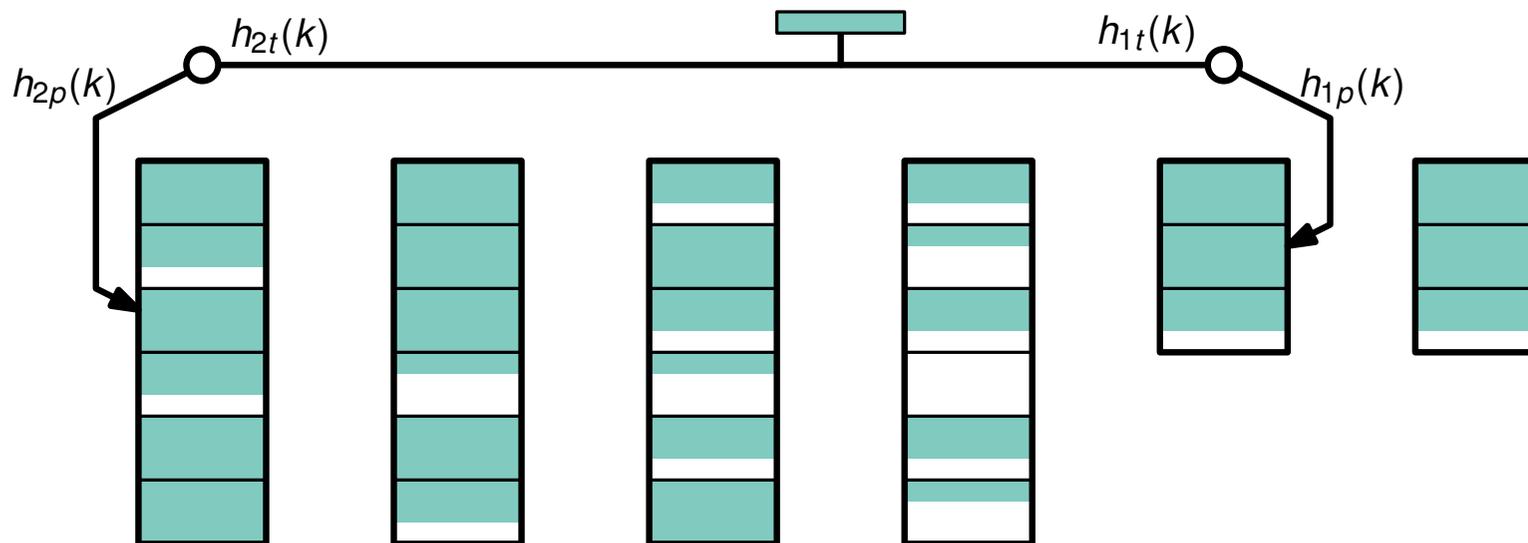
Contribution – Dynamic Space Efficient Cuckoo Table

- use **subtables** of unequal size (use **powers of 2**)
 - ▶ $h_i(k) \Rightarrow h_{it}(k)$ table and $h_{ip}(k)$ position in table
 - ▶ doubling one subtable \Leftrightarrow small overall factor
- use displacements to equalize **load imbalance**



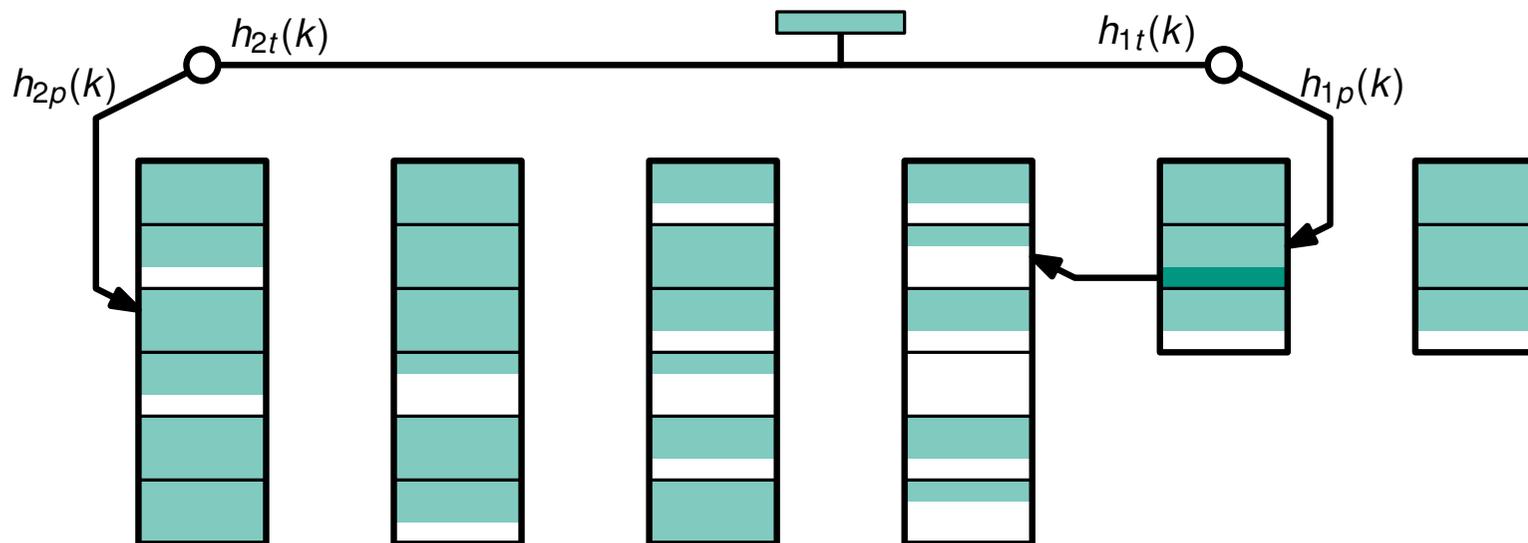
Contribution – Dynamic Space Efficient Cuckoo Table

- use **subtables** of unequal size (use **powers of 2**)
 - ▶ $h_i(k) \Rightarrow h_{it}(k)$ table and $h_{ip}(k)$ position in table
 - ▶ doubling one subtable \Leftrightarrow small overall factor
- use displacements to equalize **load imbalance**

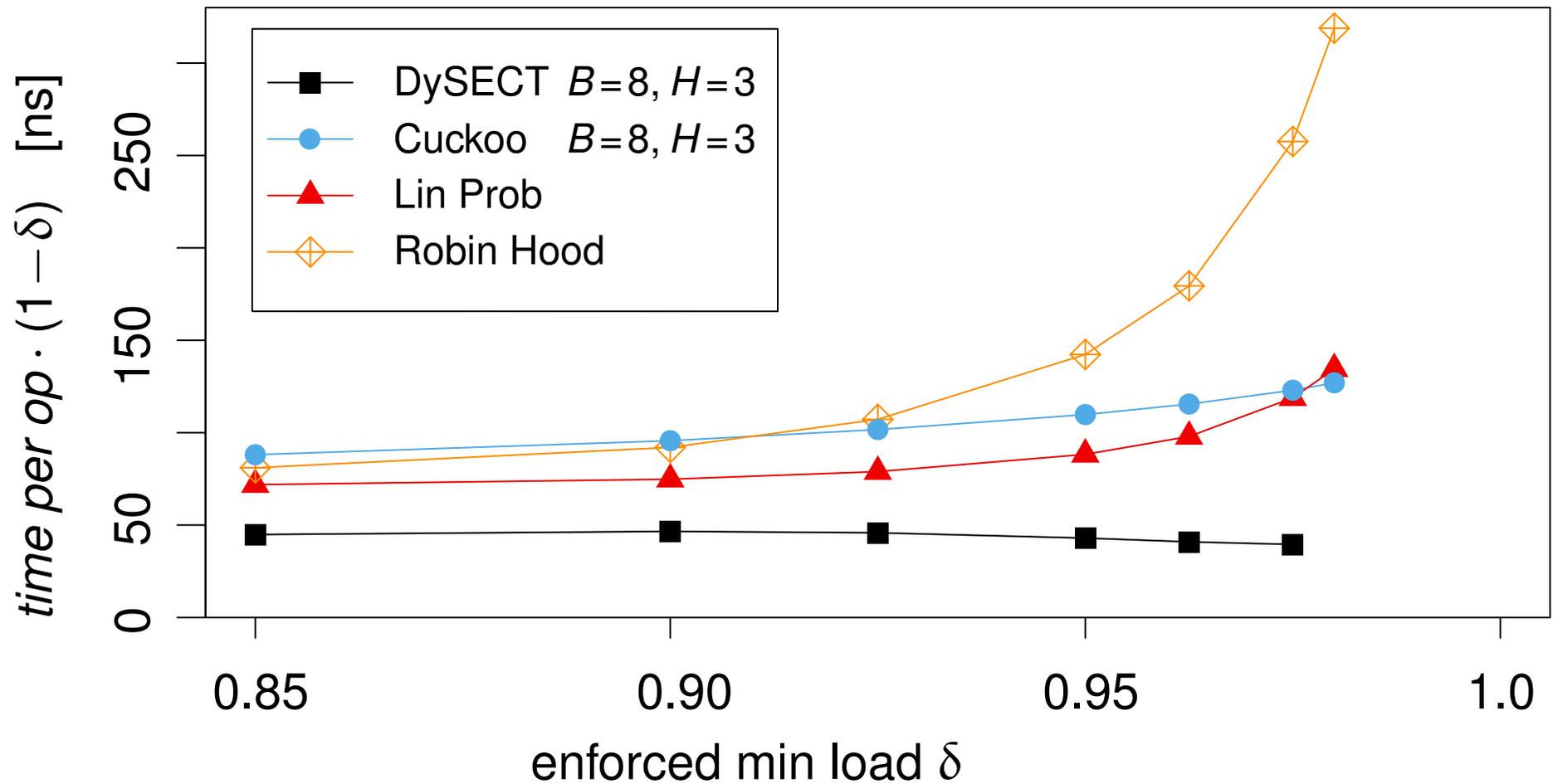


Contribution – Dynamic Space Efficient Cuckoo Table

- use **subtables** of unequal size (use **powers of 2**)
 - ▶ $h_i(k) \Rightarrow h_{it}(k)$ table and $h_{ip}(k)$ position in table
 - ▶ doubling one subtable \Leftrightarrow small overall factor
- use displacements to equalize **load imbalance**

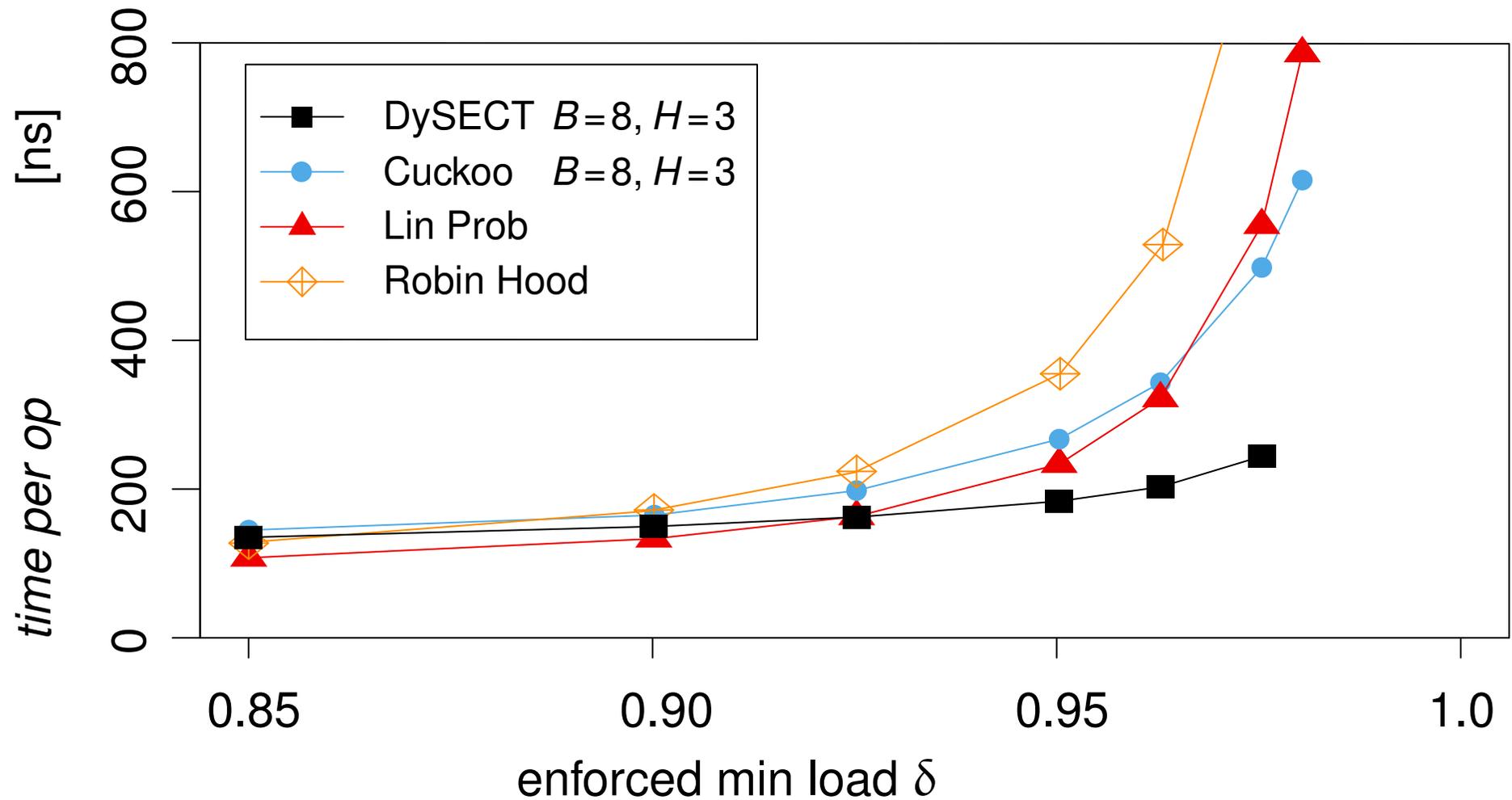


Result – Insertion into Growing Table



■ $\frac{1}{1-\delta}$ “expected time” per insertion

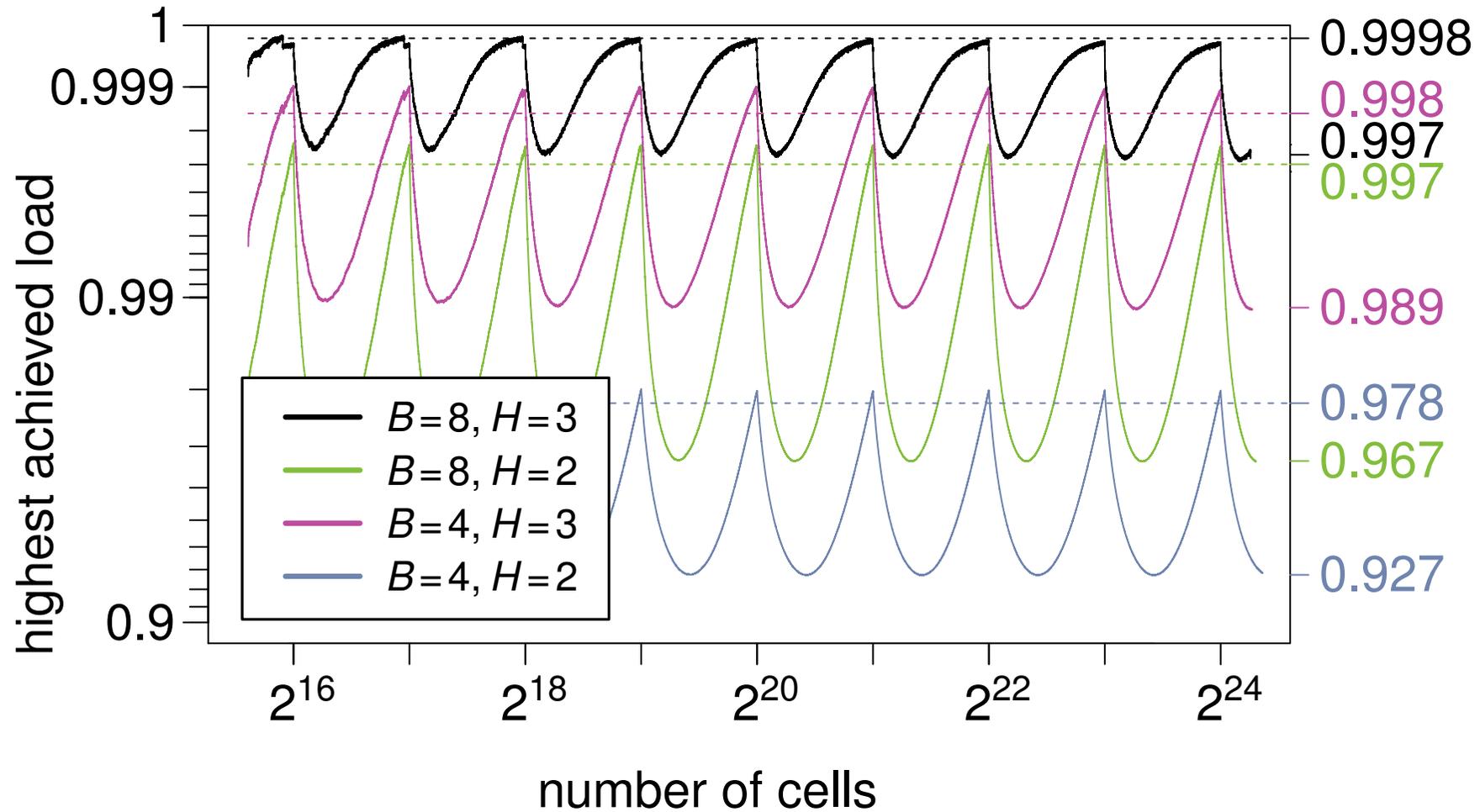
Result – Word Count Benchmark



■ CommonCrawl (avg. 12 \times)

■ not normalized

Result – Load Bound



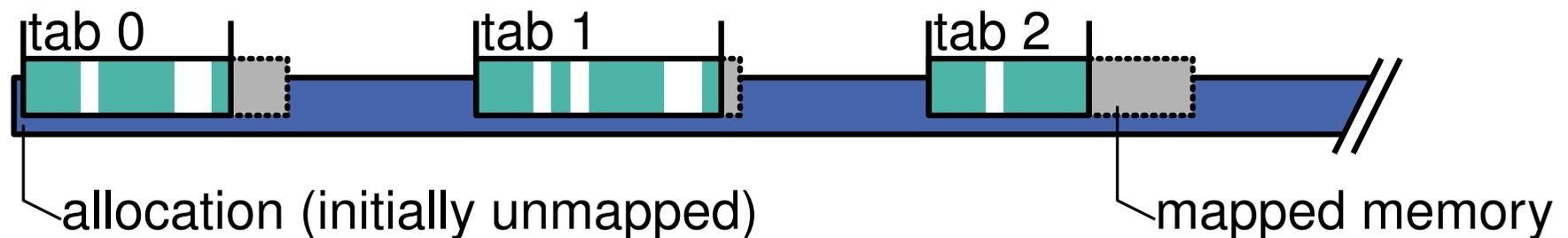
■ we are in cooperation to prove bounds

Conclusion

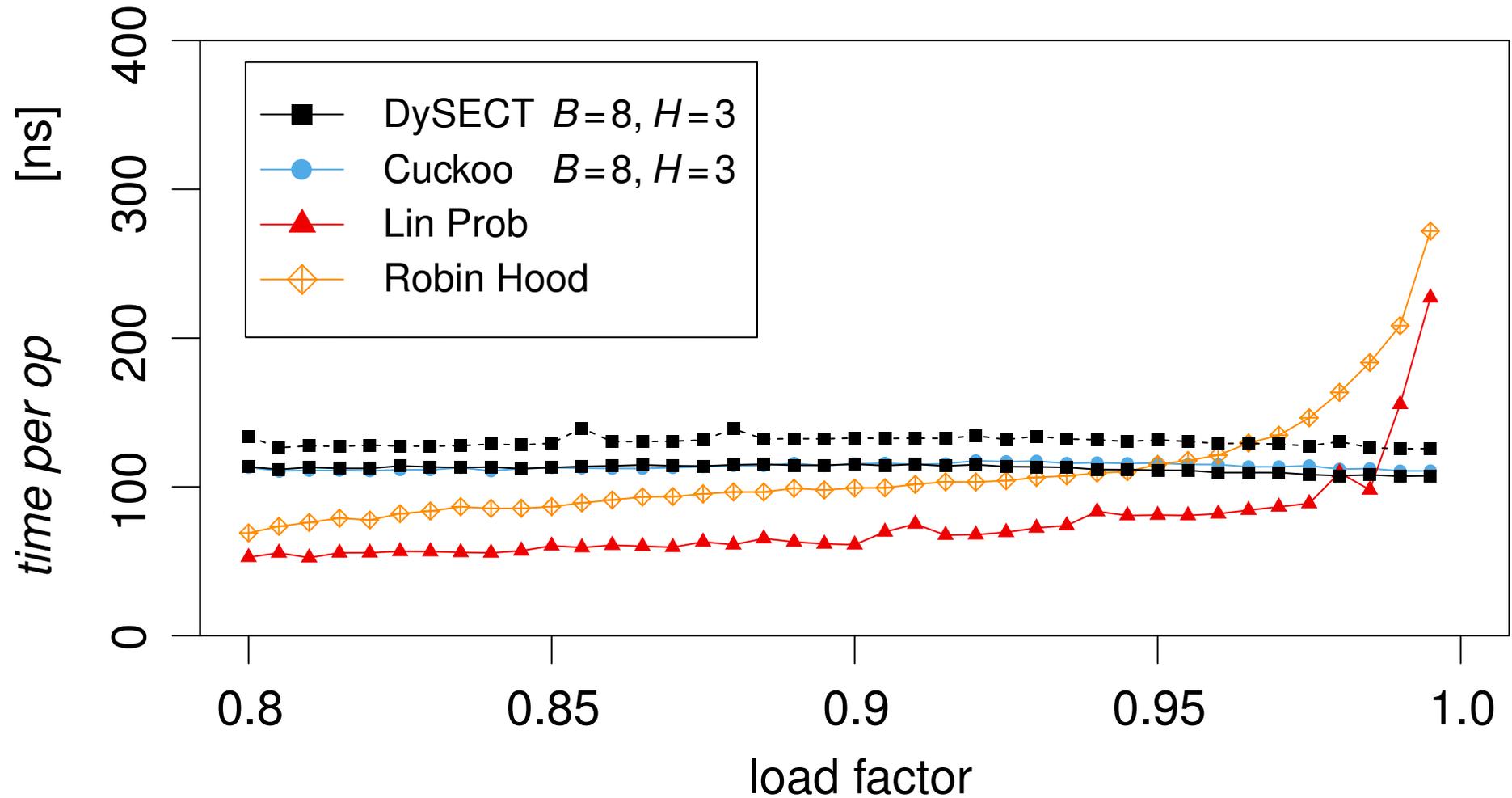
- only **dynamic** tables offer true **space efficiency**
- **lack** of published **work on dynamic** hash tables
 - even simple techniques are largely unpublished
- **DySECT**
 - no overallocation
 - constant lookup
 - addressing uses bit operations
- **cuckoo displacement** offers more untapped potential
- code available: <https://github.com/TooBiased/DySECT>

(Ab)using Overallocation

- subtables are islands of **physical memory** in a **virtual allocation**
- writing to virtual memory \approx increasing local allocation
 - ⊕ **inplace** growing
 - ⊕ no explicit indirection
 - ⊖ limited portability



Result - Successful Find



Result - Unsuccessful Find

