# Online List Labeling: Breaking the $\log^2 n$ Barrier

#### To Appear in FOCS 2022

Michael A. Bender

Stony Brook University

Alex Conway

#### **Martín Farach-Colton**

**Rutgers University** 

Hanna Komlós

**Rutgers University** 

William Kuszmaul

Nicole Wein DIMACS

### **The List Labeling Problem**

Maintain *n* elements in sorted order in an array of size *m*,  $m \ge n(1 + \Theta(1))$ 

**Goal:** minimize # elements moved per insertion/deletion (cost)



### List Labeling — Library Analogy

A librarian wants to store books alphabetically on a shelf

No gaps between books  $\implies$  many books move after each insertion



Library with no gaps



Library with gaps

#### Formulations and (Re)Discoveries

**Priority queues** [Itai, Konheim, Rodeh, ICALP'81]

File maintenance [Willard, STOC'82, SIGMOD'86]

Order maintenance [Dietz, STOC'82] [Dietz, Sleator, STOC'87]

Balanced binary trees & scapegoat trees [Andersson, WADS'89] [Andersson, Lai, SWAT'90] [Galperin, Rivest, SODA'93]

Locality preserving dictionaries [Raman, PODS'99]

#### **Sample Applications**

#### **Cache-oblivious B-trees**

[Bender, Demaine, Farach-Colton, FOCS'00] [Bender, Demaine, Iacono, Wu, SODA'02] [Brodal, Fagerberg, Jacob, SODA'02] [Bender, Fineman, Gilbert, Kuszmaul, SPAA'16]

#### **Multi-dimensional searching**

[Nekrich, Algorithmica'07] [Nekrich, Comp Geom'09] [Mortensen, SODA'03] [Toss, Pahins, Raffin, Comba Computers & Graphics'18]

#### Packed-memory arrays and density-control arrays

[Hofri, Konheim, SICOMP'87] [Katriel, '02] [Bender, Demaine, Farach-Colton, FOCS'00] [Bender, Hu, PODS'06] [Itai, Katriel, '07] [De Leo, Boncz, ICDE'19] [Bender, Fineman, Gilbert, Kopelowitz, Montes, SODA'17] [Bender, Berry, Johnson, Kroeger, McCauley, Phillips, Simon, Singh, Zage, PODS'16]

#### **Particle simulation**

[Durand, Raffin, Faure, RIPHYS'12]

#### **Database join algorithms**

[Khayyat, Lucia, Singh, Ouzzani, Papotti, Quiané-Ruiz, Tang, Kalnis, VLDB'17]

#### Managing dynamic sparse graphs

[Wheatman, Burns, IEEE BigData'21] [Wheatman, Xu, HPEC 2018, ALENEX'21] [Pandey, Wheatman, Xu, Buluc, SIGMOD'21] [De Leo, Boncz, GRADES and NDA'19, VLDB'21]

#### **Controller problem**

[Emek, Korman, Distributed Computing'11]

#### The Classical Solution [Itai, Konheim, Rodeh '81]

Step 1: Partition the array into nested intervals



#### The Classical Solution [Itai, Konheim, Rodeh '81]

**Invariant:** The density of an interval at level *i* satisfies  $d_i \leq \tau_i$ 



**Invariant:** The density of an interval at level *i* satisfies  $d_i \leq \tau_i$ 



**Invariant:** The density of an interval at level *i* satisfies  $d_i \leq \tau_i$ 

→ When an interval is over density, go up until  $d_j \leq \tau_j$  and rebalance smoothly



**Invariant:** The density of an interval at level *i* satisfies  $d_i \leq \tau_i$ 

→ When an interval is over density, go up until  $d_j \leq \tau_j$  and rebalance smoothly



**Invariant:** The density of an interval at level *i* satisfies  $d_i \leq \tau_i$ 

→ When an interval is over density, go up until  $d_j \leq \tau_j$  and rebalance smoothly



**Invariant:** The density of an interval at level *i* satisfies  $d_i \leq \tau_i$ 

→ When an interval is over density, go up until  $d_j \leq \tau_j$  and rebalance smoothly

 $d_2 = 0.56$  $- \tau_2 = 0.59$  $d_3 = 0.63$ -1  $\tau_3 = 0.71$  $d_4 = 0.50$  $- \tau_4 = 0.84$  $d_5 = 0.50$  $- \tau_5 = 1.0$ Ε F Μ В С D G H Κ 0 Q R S T [IKR '81]

 $d_1 = 0.47$ 

-1  $\tau_1 = 0.5$ 

**Invariant:** The density of an interval at level *i* satisfies  $d_i \leq \tau_i$ 

The density thresholds decrease for larger intervals

 $\tau_1 = 0.5$  $d_2 = 0.56$  $- \tau_2 = 0.59$  $d_3 = 0.63$ -1  $\tau_3 = 0.71$  $d_4 = 0.50$  $- \tau_4 = 0.84$  $d_5 = 0.50$  $--- \tau_5 = 1.0$ D Ε F Μ В С G Η Κ 0 Q R S T [IKR '81]

 $d_1 = 0.47$ 



Consider an interval at level *i* of size *s* 

- Last time it was rebalanced, it had density at most  $\tau_{i-1}$
- On next rebalance, it has density at least  $\tau_i = (1 + 1/\log n)\tau_{i-1}$
- So we've added  $\Omega(s/\log n)$  elements

Consider an interval at level *i* of size *s* 

- Last time it was rebalanced, it had density at most  $\tau_{i-1}$
- On next rebalance, it has density at least  $\tau_i = (1 + 1/\log n)\tau_{i-1}$
- So we've added  $\Omega(s/\log n)$  elements
- The amortized cost of a rebalance at a specific level is  $O(\log n)$

Consider an interval at level *i* of size *s* 

- Last time it was rebalanced, it had density at most  $\tau_{i-1}$
- On next rebalance, it has density at least  $\tau_i = (1 + 1/\log n)\tau_{i-1}$
- So we've added  $\Omega(s/\log n)$  elements
- The amortized cost of a rebalance at a specific level is  $O(\log n)$
- $\rightarrow$  Each element can be rebalanced at each of its  $\log n$  enclosing intervals

Consider an interval at level *i* of size *s* 

- Last time it was rebalanced, it had density at most  $\tau_{i-1}$
- On next rebalance, it has density at least  $\tau_i = (1 + 1/\log n)\tau_{i-1}$
- So we've added  $\Omega(s/\log n)$  elements
- The amortized cost of a rebalance at a specific level is  $O(\log n)$
- → Each element can be rebalanced at each of its  $\log n$  enclosing intervals So the total cost is:

So the total cost is:

$\log n$	amortized cost per level
$\times \log n$	levels
$= \log^2 n$	total cost

[IKR '81]

 $O(\log^2 n)$ 

#### **Upper Bound**

[Itai, Konheim, Rodeh '81]



#### **Open Problem**



#### **Open Problem**



#### **Upper Bound**

Unimproved

### Formulations and (Re)Discoveries



#### **Open Problem**









Evenly spaced



## **Our Result:**

**Theorem:** Algorithm for list labeling with  $O(\log^{3/2} n)$  expected cost per insertion/deletion

### **Our Results**





### **Our Results**





#### **Upper Bound**

Algorithm with  $O(\log^{3/2} n)$  expected cost per update

Matching lower bound for History-Independent algorithms

**Lower Bound** 

### **Our Results**



#### **Our Results**

#### **Upper Bound**

Algorithm with  $O(\log^{3/2} n)$  expected cost per update

Matching lower bound for History-Independent algorithms

#### **Lower Bound**

So what is history independence?

### History Independence — Definition

History Independence: For our problem, the distribution of occupied slots depends only on number of elements currently present





Surprising result:

•  $O(\log^2 n)$  history-independent algorithm

[Bender, Berry, Johnson, Kroeger, McCauley, Phillips, Simon, Singh, Zage, PODS'16]

Surprising result:

- O(log<sup>2</sup> n) history-independent algorithm
   [Bender, Berry, Johnson, Kroeger, McCauley, Phillips, Simon, Singh, Zage, PODS'16]
- No increase in cost/runtime to achieve history independence

Surprising result:

- O(log<sup>2</sup> n) history-independent algorithm
   [Bender, Berry, Johnson, Kroeger, McCauley, Phillips, Simon, Singh, Zage, PODS'16]
- No increase in cost/runtime to achieve history independence

- Previously viewed as a desirable property to achieve in its own right (e.g. for security guarantees)
- In contrast, we use history independence to achieve a faster algorithm

Surprising result:

- O(log<sup>2</sup> n) history-independent algorithm
   [Bender, Berry, Johnson, Kroeger, McCauley, Phillips, Simon, Singh, Zage, PODS'16]
- No increase in cost/runtime to achieve history independence

Their main idea: Use randomness to determine left/right split of elements

- Reformulate density constraint as an imbalance constraint between the left/right subintervals
- Maintains uniform distribution of left/right imbalances satisfying the constraint
- In contrast, classical algorithm always splits evenly when rebuilding

Looser balance thresholds  $\implies$  Less frequent rebalancing  $\implies$  Lower Cost

Looser balance thresholds  $\implies$  Less frequent rebalancing  $\implies$  Lower Cost

Roughly speaking, densities differ by a  $\log^{-1/2} n$  factor instead of  $\log^{-1} n$ 

Looser balance thresholds  $\implies$  Less frequent rebalancing  $\implies$  Lower Cost

Roughly speaking, densities differ by a  $\log^{-1/2} n$  factor instead of  $\log^{-1} n$ 

**Pro:** Cost of maintaining our new invariant:  $O(\log n) O(\log^{1/2} n)$  at each level

 $\implies$  Final cost of  $O(\log^{3/2} n)$ 

Looser balance thresholds  $\implies$  Less frequent rebalancing  $\implies$  Lower Cost

Roughly speaking, densities differ by a  $\log^{-1/2} n$  factor instead of  $\log^{-1} n$ 

**Pro:** Cost of maintaining our new invariant:  $O(\log n) O(\log^{1/2} n)$  at each level

 $\implies$  Final cost of  $O(\log^{3/2} n)$ 

Con: Some intervals may overflow



#### **Problem: Overflows Can Happen**

Con: Some intervals may overflow

• Density can increase by a factor of  $1 + \Theta(\log^{-1/2} n)$  at each level:



#### **Problem: Overflows Can Happen**

**Con:** Some intervals may overflow

• Density can increase by a factor of  $1 + \Theta(\log^{-1/2} n)$  at each level:

→ Density at bottom level: can be  $d(1 + \Theta(\log^{-1/2} n))^{\log n} = \omega(1)$ 

Intuition: Revert to classical solution in regions that get "too" dense

Intuition: Revert to classical solution in regions that get "too" dense

• Higher density  $\implies$  More stringent balance condition  $\implies$  Prevents overflow

Intuition: Revert to classical solution in regions that get "too" dense

- Higher density 
   —> More stringent balance condition 
   —> Prevents overflow
- Lower density  $\implies$  Less stringent balance condition  $\implies$  Less work

Intuition: Revert to classical solution in regions that get "too" dense

- Higher density  $\implies$  More stringent balance condition  $\implies$  Prevents overflow
  - $O(\log^2 n) \cos t$
- Lower density  $\implies$  Less stringent balance condition  $\implies$  Less work
  - →  $O(\log^{3/2} n) \cos t$

Intuition: Revert to classical solution in regions that get "too" dense

- Higher density  $\implies$  More stringent balance condition  $\implies$  Prevents overflow
  - $O(\log^2 n) \cos t$
- Lower density  $\Longrightarrow$  Less stringent balance condition  $\Longrightarrow$  Less work

→  $O(\log^{3/2} n) \cos t$ 

We show: Vast majority of the array is sparse  $\implies O(\log^{3/2} n)$  expected total cost

### Last Idea: Eliminate Potential Vulnerabilities

Hide the dense regions from the adversary

#### Last Idea: Eliminate Potential Vulnerabilities

Hide the dense regions from the adversary

- Our strategy may create some dense regions in the array
  - An adversary could target and exploit them to drive up cost

#### Last Idea: Eliminate Potential Vulnerabilities

Hide the dense regions from the adversary

- Our strategy may create some dense regions in the array
  - An adversary could target and exploit them to drive up cost

• **Solution:** History independence!

- The adversary can neither find nor create dense regions
- Additional idea: Apply a random shift See the paper

Prior  $O(\log^2 n)$  History-Independent Algorithm

[BBJKMPSSZ '16]

Prior  $O(\log^2 n)$  History-Independent Algorithm

[BBJKMPSSZ '16]

Loosen rebalancing rule

Speeds up the algorithm

Prior  $O(\log^2 n)$  History-Independent Algorithm

[BBJKMPSSZ '16]

Loosen rebalancing rule

Speeds up the algorithm

Stricter rebalancing rule in dense regions — Prevents overflow

Prior  $O(\log^2 n)$  History-Independent Algorithm

[BBJKMPSSZ '16]

Loosen rebalancing rule

Speeds up the algorithm

Stricter rebalancing rule in dense regions — Prevents overflow

Smoothly vary between rebalance rules – Reduces cost from  $\tilde{O}(\log^{3/2} n)$  to  $O(\log^{3/2} n)$ 

#### Summary of Results & New Open Problem



# Thank you!

- Michael A. Bender, Stony Brook University <u>bender@cs.stonybrook.edu</u>
- Alex Conway, VMWare Research <u>aconway@vmware.com</u>
- Martín Farach-Colton, Rutgers University <u>martin@farach-colton.com</u>
- Hanna Komlós, Rutgers University <u>hkomlos@gmail.com</u>
- William Kuszmaul, MIT <u>william.kuszmaul@gmail.com</u>
- Nicole Wein, DIMACS <u>nicole.wein@rutgers.edu</u>