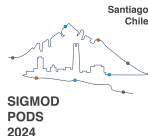


FairHash: A Fair and Memory/Time-efficient Hashmap

Nima Shahbazi, Stavros Sintos, Abolfazl Asudeh

Department of Computer Science
University of Illinois Chicago

ACM SIGMOD, Santiago Chile
June 2024



Motivation

Little attention to **fairness-aware data structures**

Hashmaps are the founding block of many applications

- Bloom filter, Count sketches, Min-wise hashing, etc.

This paper:

- Revisits hashmaps through the lens of group fairness

Review: Traditional Hash Functions

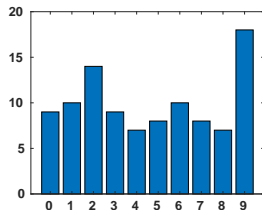
Traditional k-wise independent hashing [Sie89]

- Randomly map a key to a random value in a specific output range
- Unlikely that independent random value assignment distribute points uniformly in the buckets
- (Related topic: The [Occupancy Problem](#) [MR95])

Example

100 iid integers in range $[0, 9]$

- Not uniformly distributed within the buckets
- Number of collisions minimized when uniform distribution is satisfied



Review: Data-informed Hashmaps

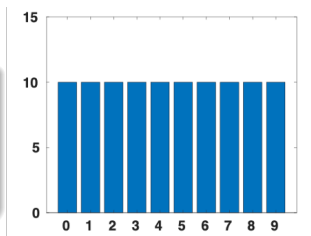
Learn a hash function that **uniformly distributes** the data across different buckets [KBC⁺18]¹

- CDF of data is constructed
- Range of values are partitioned into equi-size buckets

Example

Data-informed Hashmap learned over 100 integers in range $[0, 9]$

- Uniformly distributed within the buckets



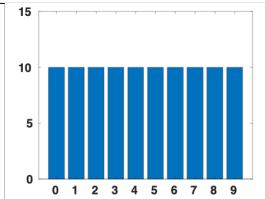
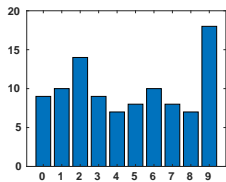
¹We refer to [KBC⁺18] as CDF-based hashmap.

Motivation at a Glance

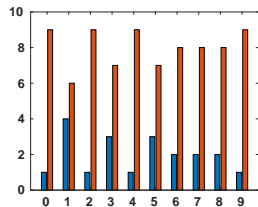
Traditional hashmaps

Data-informed hashmaps

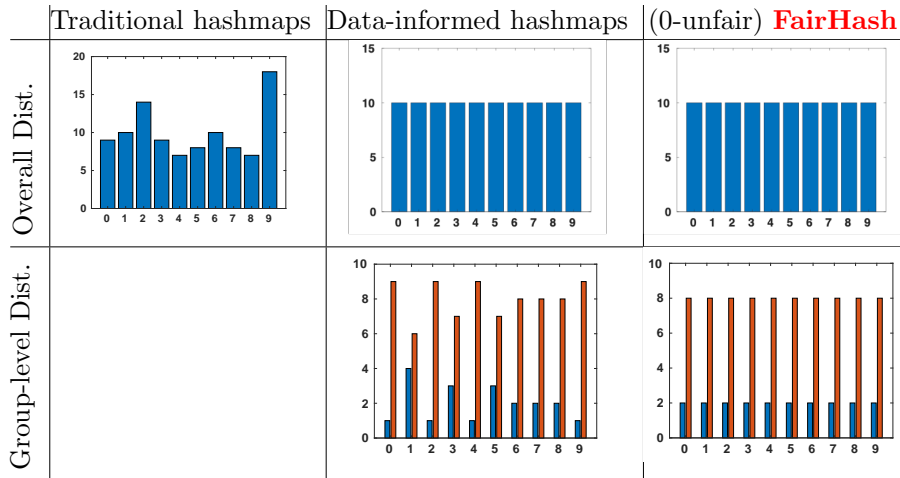
Overall Dist.



Group-level Dist.



Motivation at a Glance



Fairness Definitions

Given

- 1 A set P of n points in \mathbb{R}^d
 - ▶ Each point belong to one of the k demographic groups $\mathcal{G} = \{\mathbf{g}_1, \dots, \mathbf{g}_k\}$
- 2 A hashmap \mathcal{H} with
 - ▶ m buckets, b_1, \dots, b_m
 - ▶ a hash function $h : \mathbb{R}^d \rightarrow [1, m]$

that maps each point $p \in P$ to one of the m buckets.

Fairness Definitions

Collision Probability

- \forall random pairs $p \in P, q \in P$:
 - ▶ $Pr[h(p) = h(q)] = \frac{1}{m}$

Single fairness

- \forall random points $p_i \in \mathbf{g}_i$:
 - ▶ $Pr[h(p_i) = h(x)] = \dots = Pr[h(p_k) = h(x)] = \frac{1}{m}$

Pairwise fairness

- \forall random pairs $p_i \in \mathbf{g}_i$ and $q_i \in \mathbf{g}_i$:
 - ▶ $Pr[h(p_i) = h(q_i)] = \dots = Pr[h(p_k) = h(q_k)] = \frac{1}{m}$
- The **strongest notion of fairness**: if satisfied, the other two are also satisfied.

(ε, α) -hashmap

ε -unfairness

A hashmap is ε -**unfair**, if and only if

$$\frac{\max_{\mathbf{g} \in \mathcal{G}}(Pr_{\mathbf{g}})}{1/m} \leq (1 + \varepsilon) \Rightarrow \max_{\mathbf{g} \in \mathcal{G}}(Pr_{\mathbf{g}}) \leq \frac{1}{m}(1 + \varepsilon) \quad (1)$$

α -memory

We say a hashmap with m buckets satisfies α -**memory**, if and only if it stores at most $\alpha(m - 1)$ boundary points.

(ε, α) -hashmap

A hashmap that is ε -**unfair** and satisfies α -**memory**.

Comparisons

Hashmap	Architecture	Query time	Collision probability	Single fairness	Pairwise fairness
traditional	data-independent	$O(1)$	✗	✗	✗
CDF-based	data-dependent	$O(\log m)$	✓	✓	✗
FAIRHASH	data-dependent	$O(\log m)$	✓	✓	✓

Summary of algorithmic results

Algorithm	Assumptions		(ϵ, α) -hashmap	Performance ^a	
	No. Attributes	No. Groups		Query time	Pre-processing time
RANKING	$d \geq 2$	$k \geq 2$	$(\epsilon_R, 1)$	$O(\log m)$	$O(n^d \log n)$
SWEEP&CUT	$d \geq 1$	$k \geq 2$	$(0, \frac{n}{m})$	$O(\log n)$	$O(n \log n)$
NECKLACE _{2g}	$d \geq 1$	2	$(0, 2)$	$O(\log m)$	$O(n \log n)$
NECKLACE _{kg}	$d \geq 1$	$k > 2$	$(0, k(4 + \log n))$	$O(\log(km \log n))$	$O(mk^3 \log n + knm(n + m))$

^aThe approximate collision results are provided in the paper [SSA24].

Observation: Only the *ordering* between the tuples specify the buckets in the CDF-based hashmap.

Idea:

- Combine the attribute of a point $p \in P$ into a single score $f(p)$, using a (ranking) function $f : \mathbb{R}^d \rightarrow \mathbb{R}$
- Construct the hashmap on $f(p)$.
- **Objective:** Find the function f , according to which the unfairness is minimized.

Algorithm Overview. Use **computational geometry** concepts and Linear functions $f(p) = w^\top p$

- Consider points as **hyperplanes in the dual space**
- Design a **Ray-sweeping** algorithm to efficiently find all possible orderings \rightarrow return a function that **minimizes unfairness**.

Cut-based Algorithms

Observation: Buckets do not necessarily need to be continuous!

Idea:

- Partition the values into **more than m “bins”**.
- Many-to-one mapping: Several bins are assigned to each bucket.

Theorem

Independent of how the points are distributed and their orders, there always exists a cut-based hashmap that is **0-unfair**.

Algorithm Overview. Make two sorted passes over P

- 1 First pass: (knowing the number of tuples each bucket should contain from each group) For every tuple record the bucket it should belong
- 2 Second pass: add a cut between each pair of points that belong to different buckets.

Necklace Splitting Problem [AG21]

Divide a necklace of T beads of n' types between k' agents, such that

- ① all agents receive the same amount of beads from each type.
- ② the number of splits to the necklace is minimized.

Reduction: **points** \rightarrow **beads**; **group** \rightarrow **bead type**; **buckets** \rightarrow **agents**

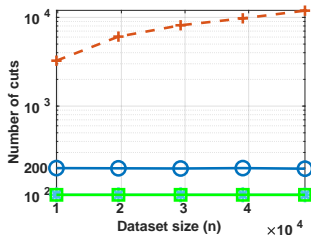
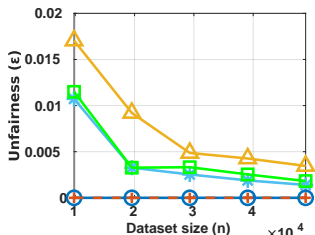
Algorithm Overview (2-groups): Iterative Algorithm

- Consider sorted P as a **circle** (p_n comes before p_1)
- **Key idea:** The circle *always* has at least one consecutive window of size $\frac{n}{m}$ that contains $\frac{|g_1|}{m}$ points from g_1 (and hence $\frac{|g_2|}{m}$ points from g_2).
- At each iteration, *efficiently find* such a window; **carve it out** of the circle; **connect the two ends** to form the circle for the next iteration

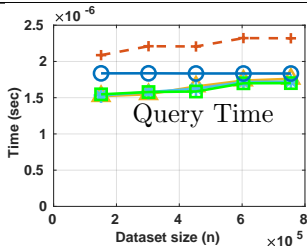
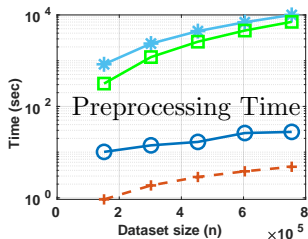
Highlighted Experiment Results

▲ Fairness-agnostic ★ Ranking-1000 ■ Ranking-100 ○ Necklace Splitting + Sweep & Cut

ADULT



CHICAGOPOP








Thank you!



Figure: Github Repository

References

-  Noga Alon and Andrei Graur, *Efficient splitting of necklaces*, ICALP 2021, 2021.
-  Tim Kraska, Alex Beutel, Ed H Chi, Jeffrey Dean, and Neoklis Polyzotis, *The case for learned index structures*, SIGMOD, 2018, pp. 489–504.
-  Rajeev Motwani and Prabhakar Raghavan, *Randomized algorithms*, Cambridge university press, 1995.
-  Alan Siegel, *On universal classes of fast high performance hash functions, their time-space tradeoff, and their applications*, FOCS, 1989, pp. 20–25.
-  Nima Shahbazi, Stavros Sintos, and Abolfazl Asudeh, *Fairhash: A fair and memory/time-efficient hashmap*, SIGMOD **2** (2024), no. 3, 1–29.