

# 1 FAIR-COUNT-MIN: Frequency Estimation under Equal Group-wise 2 Approximation Factor

3 Anonymous Author(s)

## 4 Abstract

5 Frequency estimation in streaming data often relies on sketches  
6 like Count-Min to provide approximate answers with sublinear  
7 space. However, Count-Min sketches introduce additive errors that  
8 disproportionately impact the unpopular elements, creating fairness  
9 concerns across different groups of elements. We introduce  
10 Fair-Count-Min, a frequency estimation sketch that guarantees  
11 equal expected approximation factors across element groups, thus  
12 addressing the unfairness issue. We propose a column partitioning  
13 approach with group-aware semi-uniform hashing to eliminate  
14 collisions between elements from different groups. We provide  
15 theoretical guarantees for fairness, analyze its associated cost, and  
16 validate our findings through extensive experiments on real-world  
17 datasets in comparison with representative state-of-the-art baselines.  
18 Our experimental results demonstrate that Fair-Count-Min  
19 achieves fairness with generally small additional error while main-  
20 taining efficiency comparable to that of the Count-Min sketch.

## 25 1 Introduction

26 Estimating item frequencies in data streams is a fundamental prob-  
27 lem in computer science, where given a universe of elements, the  
28 objective is to count the appearance of each within the stream. How-  
29 ever, due to memory limitations and processing time in streaming  
30 settings, exact counting is often impractical. Consequently, approx-  
31 imating the counts using techniques such as the Count-Min (CM)  
32 sketch has become widely adopted [27]. The CM estimates the  
33 frequency of elements by randomly hashing the elements into a  
34 smaller set of buckets, while adding up the frequencies of all ele-  
35 ments in each bucket. The frequency count of each bucket is then  
36 returned as an upper bound of the frequency of elements it contains.

37 The CM sketches offer space-efficient frequency estimations  
38 with small error guarantees. Nevertheless, *the error guarantees are*  
39 *additive*. Such errors may be negligible for popular elements with  
40 a high frequency in the data stream. However, as further elabo-  
41 rated in Example 2, the unpopular groups with low frequencies in  
42 the stream are disproportionately impacted under such guarantees.  
43 That is because even small additive errors can cause a significant rel-  
44 ative increase in the count estimation of unpopular (low-frequency)  
45 elements. This imbalance gives rise to a fairness issue – an as-  
46 pect overlooked in the existing work. Let us further illustrate the  
47 implications of this issue by presenting the following example:

48 **Example 1:** *Urban traffic sensors log vehicle IDs or license plate hashes as*  
49 *cars pass intersections. Buses, taxis, and delivery trucks show up hundreds of*  
50 *times per day, while private vehicles may only appear a few times. Suppose*  
51 *on average the public transportation cars pass the intersections 100 times,*  
52 *while this number is 5 for the private vehicles. In this setting, an additive*  
53 *error of 10 is tolerable for public transportation cars, but it increases the ex-*  
54 *pected estimation for private vehicles by 200%. This distortion matters: traffic*  
55 *management systems could mistakenly flag private vehicles as suspicious,*  
56 *or misallocate toll discounts and congestion pricing.*<sup>1</sup>

To address this issue, in this paper, we introduce **Fair-Count-Min** (FCM), a frequency estimation sketch with the equal group-level approximation factor guarantee, ensuring that the relative count-estimation increase is equal for various groups, irrespective of their popularity in the data stream. To the best of our knowledge, FCM is the first frequency estimation sketch with provable multiplicative error guarantees in its estimation.

Our key idea in the design of the FCM is the partitioning of the sketch buckets strategy based on a group-aware, semi-uniform hashing scheme that prevents collisions between elements of different groups. Then, we prove that, carefully selecting the number of buckets allocated to each group, our FCM sketch guarantees an equal expected approximation factor across the groups, i.e., equal ratios of the true frequency to the estimated one. The FCM sketch applies to binary and non-binary grouping of the elements, while the grouping can be based on element popularity or any other grouping strategy, e.g., based on the demographic groups.

Since *our approach is restricted to the same memory capacity as CM*, FCM does not introduce any additional memory overhead. Furthermore, FCM and CM have the same time complexities for update and query operations. In other words, the time taken to look up the frequency estimation of an element and the time to increase the frequency of an element in the data stream are the same in CM and FCM.

In a general setting, where multiple hash functions are used in the FCM, identifying the proper number of buckets allocated to each group requires solving an equation, for which we propose efficient exact and approximation algorithms.

Furthermore, we theoretically analyze the price of fairness (PoF) of FCM. In general, PoF denotes the accuracy loss incurred when fairness constraints are imposed. It arises because accuracy-maximizing outcomes often disproportionately benefit certain groups, and enforcing fairness necessitates deviations from this optimum. Interestingly, we show that in this case achieving fairness may even reduce the total additive error for a specific case when the CM sketch uses a single random hash function, while experimentally demonstrating that the PoF is small for the general settings.

Beyond the theoretical analysis, we conduct extensive experiments on multiple real-world datasets, replicating our motivating examples to evaluate the practical performance of the proposed approach. In summary, our experiments verify (a) while CM is usually not fair, FCM achieves fairness in all experiments, (b) the PoF is generally small, and (c) FCM has the same memory and time efficiency as CM.

## 109 Summary of contributions:

- We propose a novel notion of *group fairness* in the frequency estimation context, defined by the requirement that the approximation factor (the *multiplicative* overestimation error) remains

<sup>1</sup>Additional motivation examples are provided in Appendix A.

| Notation                              | Description                                 |
|---------------------------------------|---|
| $\mathcal{D}$                         | data stream                                 |
| $\mathcal{U}$                         | A universe of $n$ element types             |
| $e$                                   | an element type                             |
| $f(\cdot)$                            | true frequency of an element type           |
| $\hat{f}(\cdot)$                      | estimated frequency of an element type      |
| $\mathcal{G}$                         | The set of $\ell$ groups                    |
| $g$                                   | a group with $n_g$ element types            |
| $CM$                                  | count-min 2D array                          |
| $w, d$                                | width and depth of the sketch               |
| $w_g, d_g$                            | # of columns and rows assigned to group $g$ |
| $N$                                   | length of stream                            |
| $N_g$                                 | frequency of group $g$                      |
| $h(\cdot)$                            | random hash function                        |
| $\varepsilon_A(\cdot)$                | additive error of an element type           |
| $\alpha(\cdot)$                       | approximation factor of an element type     |
| $C_i$                                 | total frequency in $i$ -th bucket           |
| $PoF$                                 | price of fairness                           |
| $\mathcal{L}_{FCM}, \mathcal{L}_{CM}$ | total additive errors of FCM and CM         |

**Table 1: Table of notations**

equal across different groups. This contrasts with traditional approaches based on *additive* error bounds, which tend to disproportionately impact unpopular groups. Using this notion of fairness, we introduce *Fair-Count-Min* (FCM), the frequency estimation sketch that guarantees group fairness.

- We introduce a *column-partitioning* technique based on semi-uniform hash functions to develop FCM. This method assigns groups to disjoint sets of hash buckets, thereby eliminating inter-group collisions and promoting fairness in estimation. Our proposed method, FCM, is agnostic to both the grouping strategy and the underlying distribution of the data elements. Furthermore, FCM does not increase the memory and time requirements of CM.
- We theoretically prove that FCM is fair, i.e., it ensures an equal expected approximation factor across groups. Furthermore, we analyze the *price of fairness* and show that it is typically small—and in some cases, can even be negative.
- We design both exact and approximate algorithms to efficiently compute the optimal allocation of buckets per group.
- We empirically evaluate FCM on several real-world datasets and compare it with representative state-of-the-art baselines for frequency estimation such as learning-based approach [48] and CM with conservative update [36]. The results show that FCM is the only approach that successfully achieves group fairness while incurring only a generally small additional additive error, all while preserving the time and space efficiency characteristics of CM. Interestingly, CM with conservative update—though effective in minimizing additive error—exhibits the poorest performance with respect to multiplicative error. The learned approach is constrained by its capacity on identifying heavy hitters from the data sample and incurs construction and query times several orders of magnitude higher than non-learned methods.

## 2 Preliminary

### 2.1 Data Model

Let  $\mathcal{D}$  be a data stream consisting of elements (or events), each belonging to a type from a finite universe  $\mathcal{U} = \{e_1, e_2, \dots, e_n\}$ .

Each element type  $e$  be associated with a group  $g \in \mathcal{G}$ , where  $\mathcal{G} = \{g_1, g_2, \dots, g_\ell\}$ . Let  $N$  be the number of elements in the data stream  $\mathcal{D}$ . We use the function  $f : \mathcal{U} \rightarrow [N]$  to refer to the frequency of each element type in  $\mathcal{D}$ . That is,  $f(e)$  is the number of times an element of type  $e$  appears in  $\mathcal{D}$ . Formally,

$$f(e) = \sum_{x \in \mathcal{D}} \mathbb{1}(x = e),$$

where  $\mathbb{1}(x = e)$  is the indicator function that equals 1 if  $x$  is of type  $e$  and 0 otherwise.

### 2.2 Count-Min Sketch

Frequency estimation is a fundamental problem in streaming data processing. Given a data stream  $\mathcal{D}$ , the goal is to efficiently estimate the frequency  $f(e)$  of each element type  $e$ . It is easy to see that finding an exact solution to the above problem requires  $\Theta(n)$  space, as it involves maintaining  $n$  counters for each element type. To overcome this limitation, frequency estimation data structures (a.k.a sketches) have been designed to provide approximate answers to such queries while using sub-linear space. In addition to being space-efficient, these sketches provide probabilistic error guarantees through tunable parameters and enable constant-time update and query operations.

CM is among the most widely used sketches for frequency estimation. It represents the frequency estimates using a two-dimensional array  $CM$  of  $d$  rows and  $w$  columns. The sketch employs  $d$  independent hash functions  $h_1, h_2, \dots, h_d$ , where each  $h_i : \mathcal{U} \rightarrow [w]$  maps an element type to a column index (called *bucket* or *bin* in the rest of the paper) in row  $i$ . These hash functions distribute element types across different buckets to reduce the impact of hash collisions. When a new element of type  $e$  arrives in the stream, the sketch updates its counts as follows:

$$CM[i, h_i(e)] \leftarrow CM[i, h_i(e)] + 1, \quad \forall i \in [d].$$

Each row records a count for  $e$ , though potential overestimations may occur due to hash collisions. To estimate the frequency of an element type<sup>2</sup>  $e$ , the sketch returns:

$$\hat{f}(e) = \min_{i=1}^d CM[i, h_i(e)].$$

Taking the minimum across multiple rows mitigates over-counting errors. As previously mentioned, the estimation error in CM occurs due to hash collisions. For each element type  $e$  the estimated frequency  $\hat{f}(e)$  is always an upper bound on the true frequency  $f(e)$ :

$$\hat{f}(e) = f(e) + \varepsilon_A(e), \tag{1}$$

where the additive error  $\varepsilon_A(e)$  corresponds to the combined frequency of all other elements that hash into the same bucket as  $e$ . In expectation, the additive error  $\varepsilon_A(e)$  is bounded by  $\frac{N}{w}$ . Formally:

$$\mathbb{E}[\hat{f}(e)] = f(e) + \frac{\sum_{\forall e' \neq e, h(e')=h(e)} f(e')}{w} \leq f(e) + \frac{N}{w} \tag{2}$$

Several variations of the CM exist, including the CM with conservative updates [36], the Count-Mean-Min sketch [30], and other related frequency estimation structures such as the Count sketch [21].

<sup>2</sup>In the rest of the paper, we use “element  $e$ ” and “element type  $e$ ” interchangeability to refer to elements of type  $e$ .

and Spectral Bloom Filters [25]. While using CM as the foundation for the development of our sketch, and deriving our theoretical analysis, we compare our work against other baselines in experiments.

### 2.3 Fairness: Equal Expected Approximation Factor

Existing CM sketches provide an additive upper-bound error  $\epsilon_A$  on the frequency estimation of the elements. However, while  $\epsilon_A$  may be a negligible error for the popular frequent elements, it can be intolerable for the unpopular ones which appear infrequently in the sequence. In other words, the significance of the frequency estimation error is relative to the frequency value. To further clarify this using a toy example, let us consider Example 2.

**Example 2:** Consider a universe of events of types  $\mathcal{U} = \{e_1, e_2, \dots, e_n\}$ , where  $\{e_1, \dots, e_{\frac{n}{2}}\}$  belong to the group  $g_1$ , while  $\{e_{\frac{n}{2}+1}, \dots, e_n\}$  belong to  $g_2$ . Now, consider a CM sketch with a hash function  $h$  that maps the events to  $w = 200$  bins. Let  $N = 10000$  be the number of elements in the data stream. Hence, the additive estimation error is bounded by  $\epsilon_A \leq 50$ . Suppose the expected frequencies for group  $g_1$  and  $g_2$  are 1000 and 10, respectively. As a result, the expected-frequency estimation is bounded by 1050 for  $g_1$ , resulting in a small error of only 5% for group  $g_1$ . On the other hand, this error is as high as 500% for  $g_2$ .

From Example 2, it is clear that the additive error does not provide a strong guarantee independent of the element frequencies. Specifically, providing large estimation errors relative to small frequency values, it is *unfair* for unpopular groups. This motivates us to instead promote **approximation factor** as a stronger frequency estimation guarantee.

**DEFINITION 1 (APPROXIMATION FACTOR).** Let  $\hat{f}(e)$  be the frequency estimation of the element type  $e$  by a count-min sketch CM. The approximation factor of CM for  $e$  is  $\alpha(e) = \frac{\hat{f}(e)}{f(e)}$ .

Using Definition 1, we define the notion of *group fairness* that equalizes the expected approximation factor across various groups. Specifically, a CM is called fair if it satisfies Definition 2.

**DEFINITION 2 (GROUP-FAIR COUNT-MIN (FCM)).** A count-min sketch is group-fair iff:

$$\forall g_i, g_j \in \mathcal{G}, \quad \mathbb{E}_{e \in g_i} [\alpha(e)] = \mathbb{E}_{e' \in g_j} [\alpha(e')]$$

Without loss of generality and to ease the explanations, we use the binary groups  $g_1$  and  $g_2$  for explaining our sketches and drawing the analyses. Our results readily hold for non-binary and *arbitrary grouping* of element types, as we shall further discuss in Section 3.4.

### 2.4 Problem Formulation

With the necessary terms and notations defined, we now formally define our problem of interest:

**DEFINITION 3.** Given a data stream of elements  $\mathcal{D}$  drawn from a universe of element types  $\mathcal{U}$  where each element type belongs to a group  $g \in \mathcal{G}$ , design a group-fair count-min sketch, whose overall additive error increase compared to CM (price of fairness) is small.

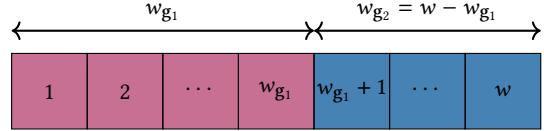


Figure 1: Illustration of a column-partitioning based group-fair min-count with one row, i.e., one hash function  $h(\cdot)$ .

### 2.5 Solution Overview

A main issue that causes the unbounded multiplicative error in CM is that high-frequency elements can collide with low-frequency ones, adding a major overestimation to their counts. Similarly, a sketch in our case might not be group-fair if elements in  $g_1$  collide with elements in  $g_2$ .

One way to reduce the likelihood of such collisions is to substantially increase the number of bins  $w$  and/or the number of independent hash functions  $d$  used in the estimation, i.e., the number of rows in the sketch. However, this approach has a fundamental drawback: to ensure that an element  $e \in g_1$  is unlikely to collide with any element in  $g_2$  across at least one of the  $d$  rows, the sketch size must scale linearly with  $n$ , the number of distinct element types—defeating the purpose of sketching.

Instead, to design group-fair CM, our goal is to *make it impossible for elements in different groups to collide*. Specifically, we ensure our goal by designing group-aware “semi-uniform” hash schemes that isolate elements of  $g_1$  and  $g_2$  by strategically partitioning and assigning the columns of the sketch to  $g_1$  and  $g_2$  (Section 3). Proving that our sketch satisfies Definition 2, we study its price of fairness (Section 5).

## 3 Fair-Count-Min using Group-Aware Semi-Uniform Hashing

Our idea for generating group-fair CM (FCM) is to utilize “semi-uniform” hash functions that separate element types of various groups. Specifically, reserving  $w_{g_1}$  bins for the group  $g_1$  and  $w_{g_2} = w - w_{g_1}$  for the group  $g_2$ , the semi-uniform hash function  $h(\cdot)$  assigns each element in  $g_1$  to the first  $w_{g_1}$  bins and the others to the remaining  $w_{g_2}$  bins (Figure 1):

$$h : \begin{cases} g_1 \rightarrow [w_{g_1}] \\ g_2 \rightarrow [w_{g_2}] \end{cases}$$

The key question we shall answer in the rest of this section is whether there exists a value  $w_{g_1}$  (hence  $w_{g_2} = w - w_{g_1}$ ) for which a CM based on the semi-uniform hash scheme becomes group-fair.

In the following, first, we focus on the case where  $d = 1$ , i.e., the sketch is based on only one hash function  $h$ . Next, we extend our analysis to the general values of  $d$ .

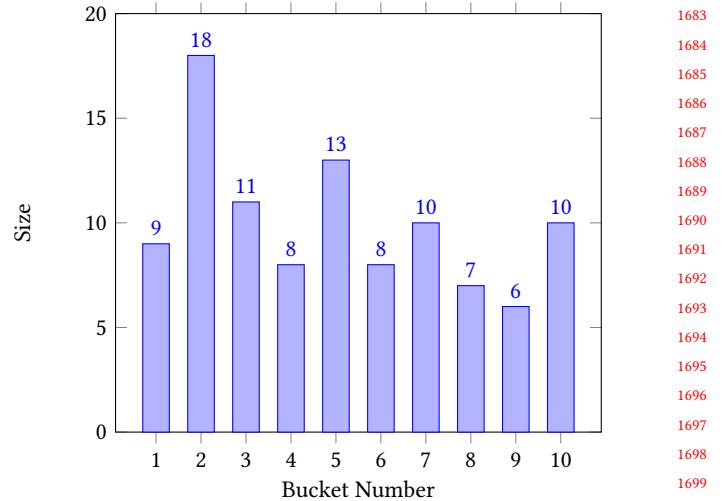
### 3.1 $d = 1$

First, let us derive the expected approximation factor for an element  $e$  in group  $g_1$ . For the case of  $d = 1$ , we employ the stronger guarantee of the average approximation factor instead of the expected one. For clarity, we use the expectation notation  $\mathbb{E}$  to also denote the average. Let  $n_{g_1}$  be the number of element types in  $g_1$ .

## 1625 APPENDIX

## 1626 A Additional Motivation Examples

1628 **Example 3:** On platforms like Twitter, hashtags appear as a rapid, high-  
 1629 volume data stream. Automated or troll-bot hashtags can generate millions  
 1630 of occurrences per hour, while rare but potentially critical hashtags may  
 1631 surface only a few thousand times. In a CM, a small additive error can  
 1632 severely distort counts for the rarer tags. Although this may seem like a  
 1633 “free boost”, the resulting collision-induced overcounts undermine reliability,  
 1634 making it hard to determine whether a rare hashtag is truly trending or  
 1635 merely inflated by overlap with high-volume tags. As a result, trend detection,  
 1636 resource allocation, and moderation decisions such as filtering become  
 1637 noisier—false signals may emerge while genuinely urgent topics are delayed  
 1638 or misclassified.



1683  
 1684  
 1685  
 1686  
 1687  
 1688  
 1689  
 1690  
 1691  
 1692  
 1693  
 1694  
 1695  
 1696  
 1697  
 1698  
 1699  
 1700  
**Figure 33: Illustration of a random distribution of  $n = 100$  elements to  $w = 10$  buckets.**

1640  
 1641  
 1642  
 1643  
 1644  
**Example 4:** In network monitoring, consider two groups of IPs: malicious  
 1645 and benign IPs. While benign IP addresses (e.g., Google DNS) generate massive  
 1646 traffic, the malicious IPs might only appear sporadically. A CM with a small  
 1647 additive error has a negligible impact on the counts of the benign group. This  
 1648 error, on the other hand, can artificially inflate the counts for the malicious  
 1649 IPs, hiding meaningful anomalies. This makes it harder to detect subtle but  
 1650 critical security threats.  
 1651

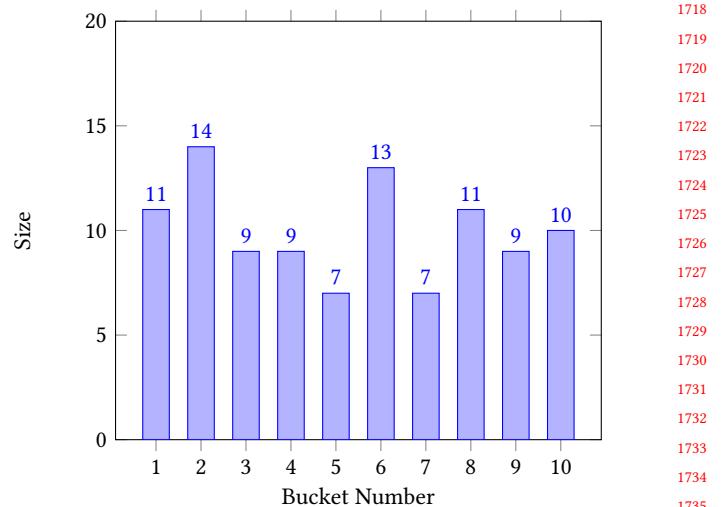
1652  
 1653  
 1654  
 1655  
 1656  
 1657  
**Example 5:** In online advertising, platforms track how often each ad or  
 1658 user click event occurs. Popular ads may generate millions of impressions,  
 1659 while niche ads may only get a few hundred. With CM, small additive errors  
 1660 hardly affect big advertisers but can completely distort the performance  
 1661 metrics for small advertisers—for example, inflating a niche ad’s 50 clicks  
 1662 into 500 due to collisions. This creates unfair reporting, potentially wasting  
 1663 budgets or hiding click fraud.  
 1664

## 1665 B Price of Fairness: Randomness vs. Uniformity

1666 In Section 5, we observed a counterintuitive result for  $d = 1$ , proving  
 1667 a negative price of fairness (PoF) for FCM, compared to a regular  
 1668 CM sketch that uses a random hash function. This result can be  
 1669 explained by the fact that a random hash function is unlikely to  
 1670 uniformly distribute the elements to the buckets (an interesting  
 1671 related topic is the *occupant problem* [59]). As a result, some of the  
 1672 buckets will have more than  $\frac{n}{w}$  elements in them. For example,  
 1673 Figure 33 illustrates a random hashing of  $n = 100$  element types  
 1674 into  $w = 10$  buckets. While a uniform distribution would allocate  
 1675 10 elements to each bucket, the random hashing allocated up to 18  
 1676 elements in one bucket.  
 1677

1701 Elements in large buckets will suffer from a large additive error,  
 1702 and there are a large number of them in such buckets. Hence, such  
 1703 buckets significantly increase the total additive error of the CM.

1704 FCM increases the size-uniformity of the buckets by reserving  
 1705  $w_g = \frac{n_g}{n} w$  for each group  $g$ . For example, Figure 34 illustrates  
 1706 the allocation of the 100 elements to the 10 buckets, by dividing  
 1707 (arbitrarily) the elements into two groups of size  $n_{g_1} = 50$ ,  $n_{g_2} = 50$ .  
 1708 As a result, each group is allocated 5 buckets, and the distribution  
 1709 of the elements is more uniform, reducing the maximum size of the  
 1710 buckets (in this example) to 14. Increasing the uniformity (reducing  
 1711 the change of large-size buckets) helps to reduce the total additive  
 1712 error of FCM versus CM.  
 1713



1714  
 1715  
 1716  
 1717  
 1718  
 1719  
 1720  
 1721  
 1722  
 1723  
 1724  
 1725  
 1726  
 1727  
 1728  
 1729  
 1730  
 1731  
 1732  
 1733  
 1734  
 1735  
 1736  
**Figure 34: Illustration of a random distribution of  $n = 100$  elements with two groups of sizes  $n_{g_1} = 50$  and  $n_{g_2} = 50$  to  $w = 10$  buckets, based on FCM.**

## 1741 B.1 PoF for Uniform hashing ( $d = 1$ )

1742 Using a hashing scheme (such as data-informed hashmaps [52])  
 1743 that ensures uniformity by allocating exactly  $\frac{n}{w}$  elements to each  
 1744 bucket can resolve the non-uniformity issue in the CM sketches.  
 1745 In the following, we compute the PoF of FCM for  $d = 1$  when a  
 1746 uniform hashing scheme is used.

1747 Since the hashing is uniform, the size of each bucket  $i$  in the CM  
 1748 sketch is

$$1749 C_i = \frac{n}{w}$$

1750 Therefore, each element  $e_j$  collides with exactly  $\frac{n}{w} - 1$  other ele-  
 1751 ments. In other words, the frequency of each element contributes  
 1752 to the additive error of exactly  $\frac{n}{w} - 1$  other elements.

1753 As a result, the total additive error can be computed as

$$\begin{aligned} 1754 \mathcal{L}_{CM} &= \sum_{j=1}^n \varepsilon_A(e_j) \\ 1755 &= \sum_{j=1}^n f(e_j) \left( \frac{n}{w} - 1 \right) = \left( \frac{n}{w} - 1 \right) \sum_{j=1}^n f(e_j) \\ 1756 &= N \left( \frac{n}{w} - 1 \right) \end{aligned} \quad (16)$$

1757 Now, let us compute  $\mathcal{L}_{FCM}$ , the total additive error for the regu-  
 1758 lar fair-count-min sketch.

$$1759 \mathcal{L}_{FCM} = \sum_{j=1}^n \varepsilon_A(e_j) = \sum_{e_j \in g_1} \varepsilon_A(e_j) + \dots + \sum_{e_j \in g_\ell} \varepsilon_A(e_j)$$

1760 Following the same calculation as in Equation 16, for every group  
 1761  $g \in \mathcal{G}$ ,

$$1762 \sum_{e_j \in g} \varepsilon_A(e_j) = N_g \left( \frac{n_g}{w_g} - 1 \right)$$

1763 Hence,

$$1764 \mathcal{L}_{FCM} = \sum_{e_j \in g_1} \varepsilon_A(e_j) + \dots + \sum_{e_j \in g_\ell} \varepsilon_A(e_j) = \sum_{g \in \mathcal{G}} N_g \left( \frac{n_g}{w_g} - 1 \right)$$

1765 From Theorem 1, we know,  $w_g = \frac{n_g}{n} w$ ,  $\forall g \in \mathcal{G}$ , while  $\sum_{g \in \mathcal{G}} N_g =$   
 1766  $N$ . Therefore,

$$\begin{aligned} 1767 \mathcal{L}_{FCM} &= \sum_{g \in \mathcal{G}} N_g \left( \frac{n_g}{w_g} - 1 \right) \\ 1768 &= \sum_{g \in \mathcal{G}} N_g \left( \frac{n_g}{\frac{n_g}{n} w} \right) - \sum_{g \in \mathcal{G}} N_g = \sum_{g \in \mathcal{G}} N_g \left( \frac{n_g}{w} \right) - N \\ 1769 &= \left( \frac{n}{w} \right) \sum_{g \in \mathcal{G}} N_g - N \quad // \text{since } \frac{n_g}{w_g} = \frac{n}{w}, \forall g \in \mathcal{G} \\ 1770 &= N \left( \frac{n}{w} \right) - N = N \left( \frac{n}{w} - 1 \right) \end{aligned} \quad (17)$$

1771 Finally, combining Equations 12, 16, and 17, we get,

$$1772 PoF = \mathcal{L}_{FCM} - \mathcal{L}_{CM} = N \left( \frac{n}{w} - 1 \right) - N \left( \frac{n}{w} - 1 \right) = 0 \quad (18)$$

1773 That is, the PoF is zero for  $d = 1$ , when a uniform hashing scheme  
 1774 is used.

## 1799 B.2 PoF $d = 1$ vs. $d > 1$ on Synthetic

### 1800 C Additional Experiment Results

1801 In this section, we present comprehensive experimental results  
 1802 across five datasets: 1) GOOGLE NGRAM, 2) CENSUS, 3) NYC BUS,  
 1803 4) DDOS, 5) REDDIT-TWITTER. We also provide the absolute values of the ap-  
 1804 proximation factors and additive errors for each group to facilitate a  
 1805 clearer understanding of the trends in unfairness and the associated  
 1806 price of fairness under each setting.

| 1857 | Total Additive Error |             |  | $n_{g_i}$ ( $n = 10,000$ )                   |  |  |  |  |  |  |  | 1915       |            |
|------|----------------------|-------------|--|--|--|--|--|--|--|--|--|------------|------------|
|      |                      |             |  | 9000   | 8000   | 7000   | 6000   | 5000   | 4000   | 3000   | 2000   |            |            |
| 1858 | $d = 1$              | theoretical | CM   | 19,047,019                                   | 28,054,816                                   | 37,081,001                                   | 46,092,289                                   | 55,121,050                                   | 64,147,002                                   | 73,186,882                                   | 82,323,243                                   | 91,265,426 |            |
| 1859 |                      |             | FCM  | 18,985,224                                   | 28,028,987                                   | 37,003,399                                   | 46,050,674                                   | 55,115,538                                   | 64,089,586                                   | 73,165,549                                   | 82,249,909                                   | 91,228,281 |            |
| 1860 |                      | empirical   | CM   | 19,039,343                                   | 28,085,806                                   | 37,053,427                                   | 46,096,489                                   | 55,113,194                                   | 64,234,176                                   | 73,227,052                                   | 82,283,374                                   | 91,328,074 |            |
| 1861 |                      |             | FCM  | 18,991,509                                   | 27,982,573                                   | 37,002,246                                   | 46,071,394                                   | 55,086,890                                   | 64,133,815                                   | 73,168,189                                   | 82,230,649                                   | 91,276,271 |            |
| 1862 |                      | PoF         | <span style="color: green;">↓</span> 47,834  | <span style="color: green;">↓</span> 103,233 | <span style="color: green;">↓</span> 51,181  | <span style="color: green;">↓</span> 25,095  | <span style="color: green;">↓</span> 26,304  | <span style="color: green;">↓</span> 100,361 | <span style="color: green;">↓</span> 58,863  | <span style="color: green;">↓</span> 52,725  | <span style="color: green;">↓</span> 51,803  | 1920       |            |
| 1863 |                      | $d = 5$     | empirical                                    | CM   | 7,964,348                                    | 11,572,548                                   | 16,458,026                                   | 22,023,181                                   | 28,305,699                                   | 34,442,308                                   | 40,959,703                                   | 47,230,290 | 54,257,770 |
| 1864 |                      |             |  | FCM  | 11,695,556                                   | 17,114,933                                   | 22,666,115                                   | 28,053,739                                   | 33,856,154                                   | 39,448,942                                   | 44,913,211                                   | 50,089,698 | 55,893,637 |
| 1865 |                      | PoF         | <span style="color: red;">↑</span> 3,731,208 | <span style="color: red;">↑</span> 5,542,385 | <span style="color: red;">↑</span> 6,208,089 | <span style="color: red;">↑</span> 6,030,558 | <span style="color: red;">↑</span> 5,550,455 | <span style="color: red;">↑</span> 5,006,634 | <span style="color: red;">↑</span> 3,953,508 | <span style="color: red;">↑</span> 2,859,408 | <span style="color: red;">↑</span> 1,635,867 | 1923       |            |

Table 2: comparison of CM and FCM w.r.t additive errors and the price of fairness across different  $n_{g_i}$  values for  $d = 1$  and  $d = 5$ .

|      |      |
|------|------|
| 1866 | 1916 |
| 1867 | 1917 |
| 1868 | 1918 |
| 1869 | 1919 |
| 1870 | 1920 |
| 1871 | 1921 |
| 1872 | 1922 |
| 1873 | 1923 |
| 1874 | 1924 |
| 1875 | 1925 |
| 1876 | 1926 |
| 1877 | 1927 |
| 1878 | 1928 |
| 1879 | 1929 |
| 1880 | 1930 |
| 1881 | 1931 |
| 1882 | 1932 |
| 1883 | 1933 |
| 1884 | 1934 |
| 1885 | 1935 |
| 1886 | 1936 |
| 1887 | 1937 |
| 1888 | 1938 |
| 1889 | 1939 |
| 1890 | 1940 |
| 1891 | 1941 |
| 1892 | 1942 |
| 1893 | 1943 |
| 1894 | 1944 |
| 1895 | 1945 |
| 1896 | 1946 |
| 1897 | 1947 |
| 1898 | 1948 |
| 1899 | 1949 |
| 1900 | 1950 |
| 1901 | 1951 |
| 1902 | 1952 |
| 1903 | 1953 |
| 1904 | 1954 |
| 1905 | 1955 |
| 1906 | 1956 |
| 1907 | 1957 |
| 1908 | 1958 |
| 1909 | 1959 |
| 1910 | 1960 |
| 1911 | 1961 |
| 1912 | 1962 |
| 1913 | 1963 |
| 1914 | 1964 |
|      | 1965 |
|      | 1966 |
|      | 1967 |
|      | 1968 |
|      | 1969 |
|      | 1970 |
|      | 1971 |
|      | 1972 |

1973

1974

1975

1976

1977

1978

1979

1980

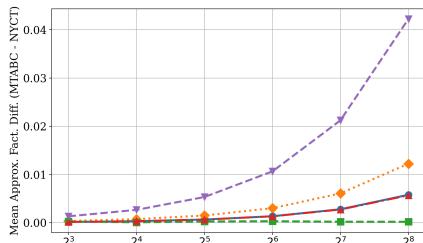
1981

1982

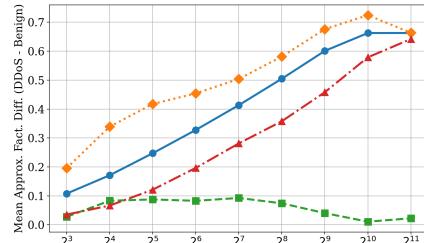
1983

1984

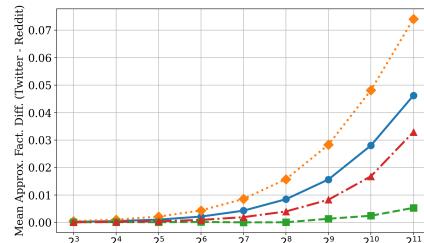
1985



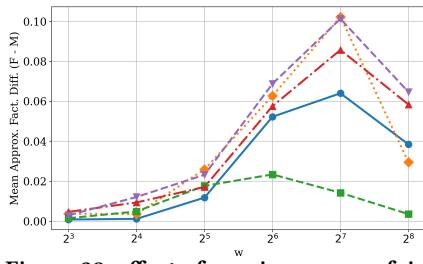
**Figure 35: effect of varying  $w$  on unfairness, NYC BUS,  $d = 5$ .**



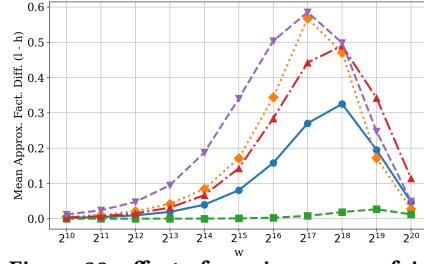
**Figure 36: effect of varying  $w$  on unfairness, DDoS,  $d = 5$ .**



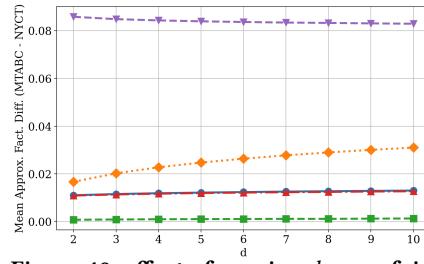
**Figure 37: effect of varying  $w$  on unfairness, REDDIT-TWITTER,  $d = 5$ .**



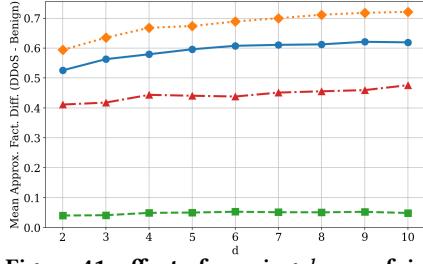
**Figure 38: effect of varying  $w$  on unfairness, CENSUS,  $d = 5$ .**



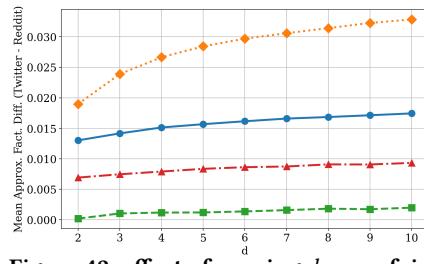
**Figure 39: effect of varying  $w$  on unfairness, GOOGLE NGRAM,  $d = 5$ .**



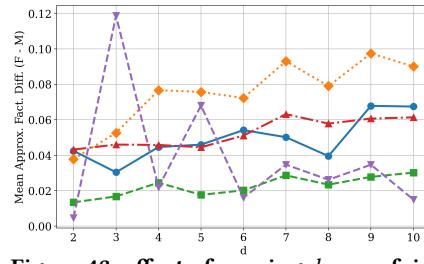
**Figure 40: effect of varying  $d$  on unfairness, NYC BUS,  $w = 512$ .**



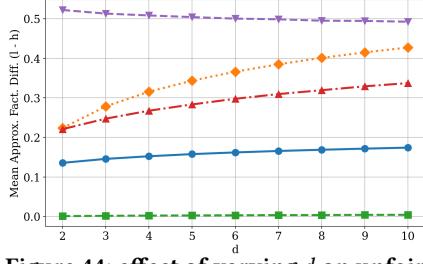
**Figure 41: effect of varying  $d$  on unfairness, DDoS,  $w = 64$ .**



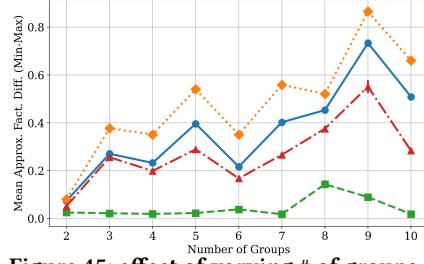
**Figure 42: effect of varying  $d$  on unfairness, REDDIT-TWITTER,  $w = 512$ .**



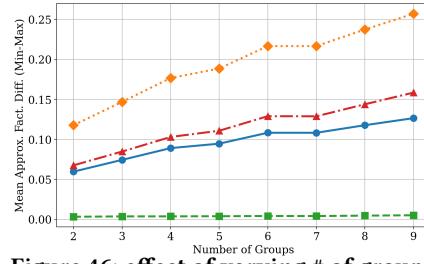
**Figure 43: effect of varying  $d$  on unfairness, CENSUS,  $w = 64$ .**



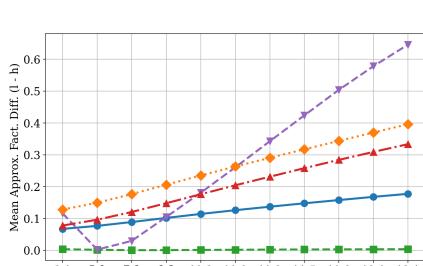
**Figure 44: effect of varying  $d$  on unfairness, GOOGLE NGRAM,  $w = 65536$ .**



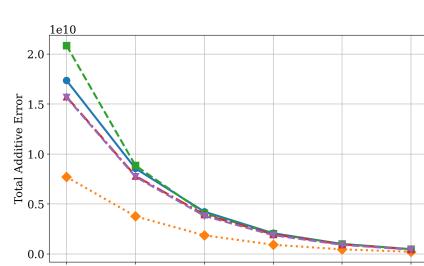
**Figure 45: effect of varying # of groups  $\ell$  on unfairness, CENSUS,  $d = 5$ .**



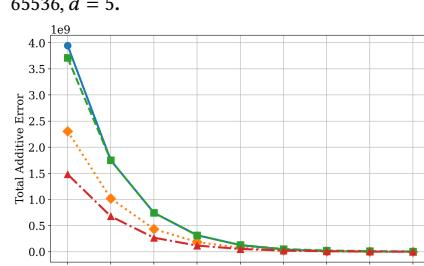
**Figure 46: effect of varying # of groups  $\ell$  on unfairness, GOOGLE NGRAM,  $w = 65536, d = 5$ .**



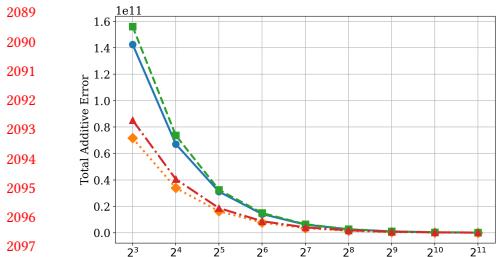
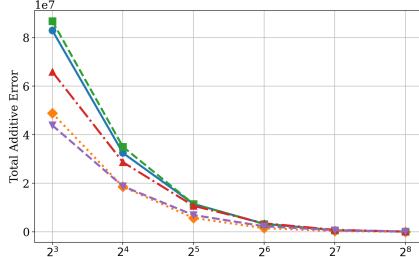
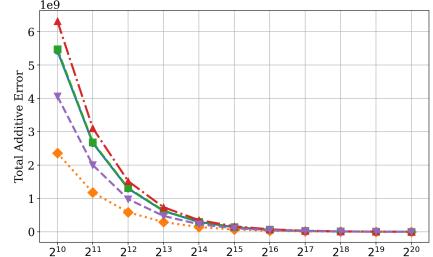
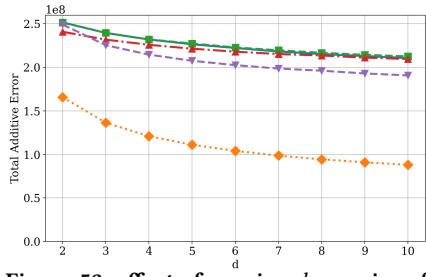
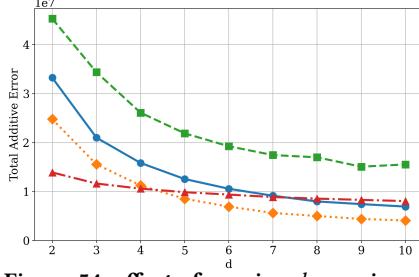
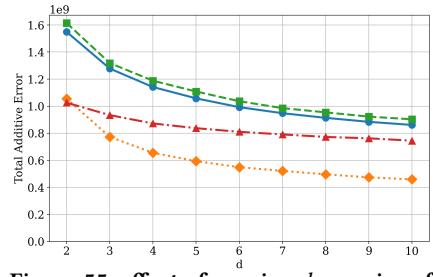
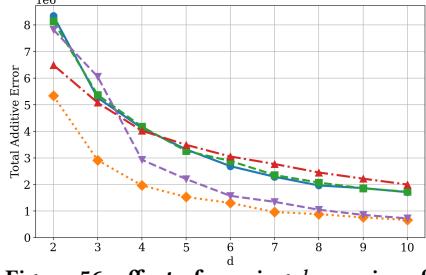
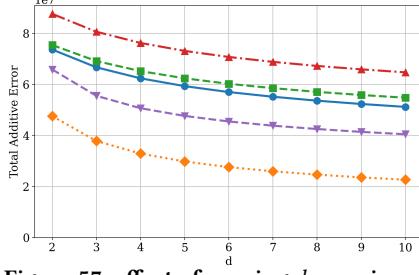
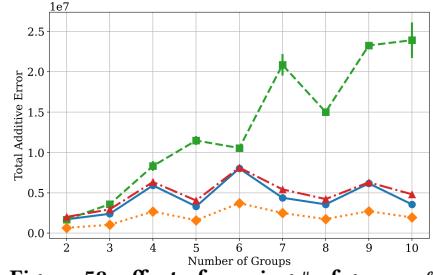
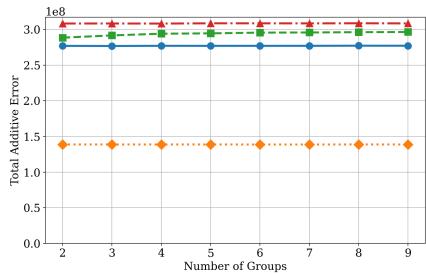
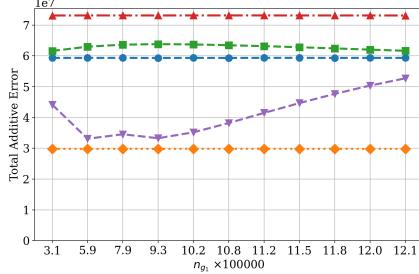
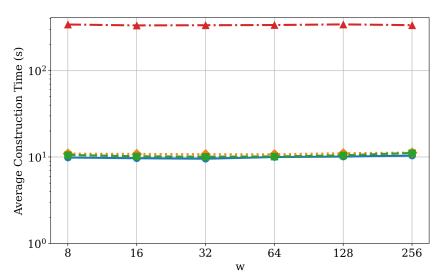
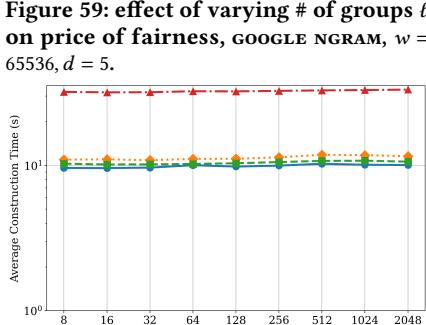
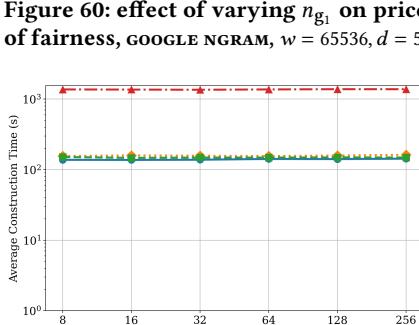
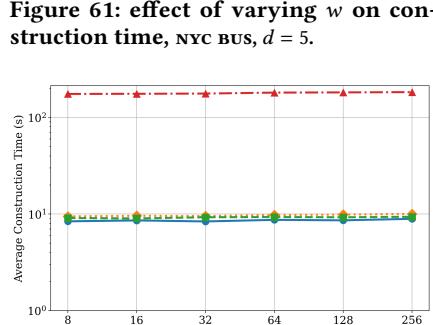
**Figure 47: effect of varying  $n_{g_1}$  on unfairness, GOOGLE NGRAM,  $w = 65536, d = 5$ .**



**Figure 48: effect of varying  $w$  on price of fairness, NYC BUS,  $d = 5$ .**

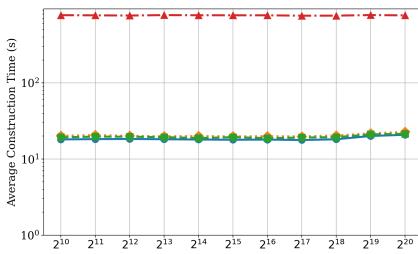


**Figure 49: effect of varying  $w$  on price of fairness, DDoS,  $d = 5$ .**

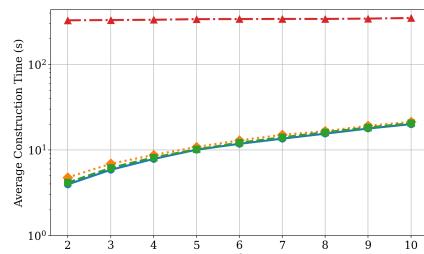
Figure 50: effect of varying  $w$  on price of fairness, REDDIT-TWITTER,  $d = 5$ .Figure 51: effect of varying  $w$  on price of fairness, CENSUS,  $d = 5$ .Figure 52: effect of varying  $w$  on price of fairness, GOOGLE NGRAM,  $d = 5$ .Figure 53: effect of varying  $d$  on price of fairness, NYC BUS,  $w = 512$ .Figure 54: effect of varying  $d$  on price of fairness, DDOS,  $w = 64$ .Figure 55: effect of varying  $d$  on price of fairness, REDDIT-TWITTER,  $w = 512$ .Figure 56: effect of varying  $d$  on price of fairness, CENSUS,  $w = 64$ .Figure 57: effect of varying  $d$  on price of fairness, GOOGLE NGRAM,  $w = 65536$ .Figure 58: effect of varying # of groups  $\ell$  on price of fairness, CENSUS,  $w = 64, d = 5$ .Figure 59: effect of varying # of groups  $\ell$  on price of fairness, GOOGLE NGRAM,  $w = 65536, d = 5$ .Figure 60: effect of varying  $n_{g_i}$  on price of fairness, GOOGLE NGRAM,  $w = 65536, d = 5$ .Figure 61: effect of varying  $w$  on construction time, NYC BUS,  $d = 5$ .Figure 62: effect of varying  $w$  on construction time, DDOS,  $d = 5$ .Figure 63: effect of varying  $w$  on construction time, REDDIT-TWITTER,  $d = 5$ .Figure 64: effect of varying  $w$  on construction time, CENSUS,  $d = 5$ .

2089  
2090  
2091  
2092  
2093  
2094  
2095  
2096  
2097  
2098  
2099  
2100  
2101  
2102  
2103  
2104  
2105  
2106  
2107  
2108  
2109  
2110  
2111  
2112  
2113  
2114  
2115  
2116  
2117  
2118  
2119  
2120  
2121  
2122  
2123  
2124  
2125  
2126  
2127  
2128  
2129  
2130  
2131  
2132  
2133  
2134  
2135  
2136  
2137  
2138  
2139  
2140  
2141  
2142  
2143  
2144  
2145  
2146  
2147  
2148  
2149  
2150  
2151  
2152  
2153  
2154  
2155  
2156  
2157  
2158  
2159  
2160  
2161  
2162  
2163  
2164  
2165  
2166  
2167  
2168  
2169  
2170  
2171  
2172  
2173  
2174  
2175  
2176  
2177  
2178  
2179  
2180  
2181  
2182  
2183  
2184  
2185  
2186  
2187  
2188  
2189  
2190  
2191  
2192  
2193  
2194  
2195  
2196  
2197  
2198  
2199  
2200  
2201  
2202  
2203  
2204

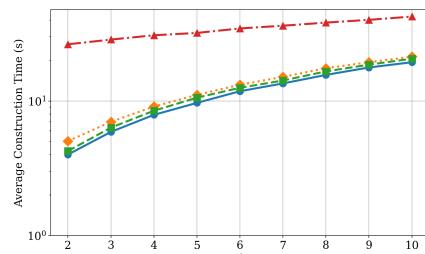
2205  
2206  
2207  
2208  
2209  
2210  
2211  
2212  
2213  
2214  
2215  
2216  
2217

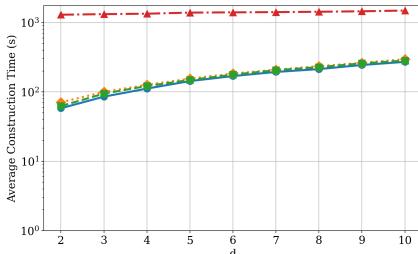
**Figure 65: effect of varying  $w$  on construction time, GOOGLE NGRAM,  $d = 5$ .**



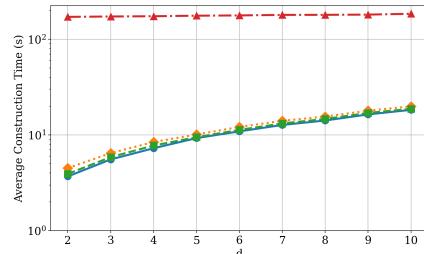
**Figure 66: effect of varying  $d$  on construction time, NYC BUS,  $w = 512$ .**



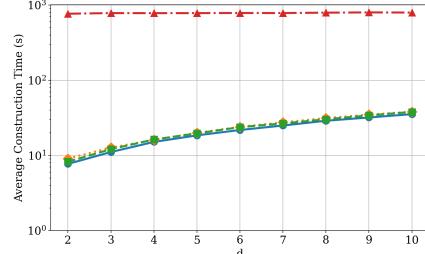
**Figure 67: effect of varying  $d$  on construction time, DDOS,  $w = 64$ .**



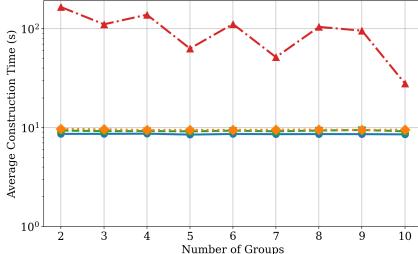
**Figure 68: effect of varying  $d$  on construction time, REDDIT-TWITTER,  $w = 512$ .**



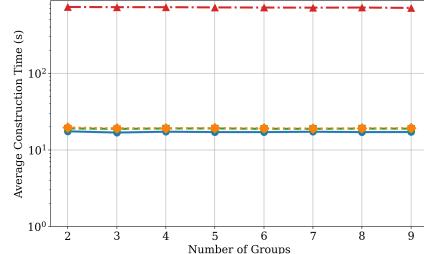
**Figure 69: effect of varying  $d$  on construction time, CENSUS,  $w = 64$ .**



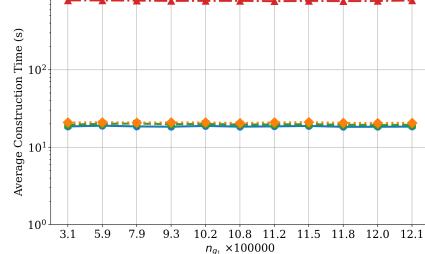
**Figure 70: effect of varying  $d$  on construction time, GOOGLE NGRAM,  $w = 65536$ .**



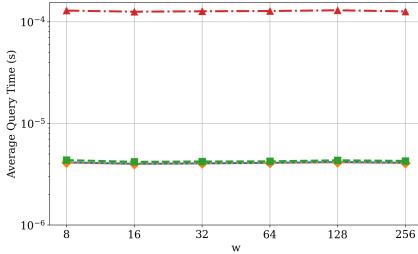
**Figure 71: effect of varying # of groups  $\ell$  on construction time, CENSUS,  $w = 64, d = 5$ .**



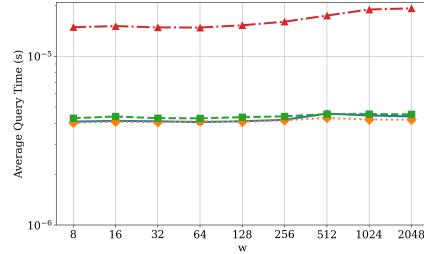
**Figure 72: effect of varying # of groups  $\ell$  on construction time, GOOGLE NGRAM,  $w = 65536, d = 5$ .**



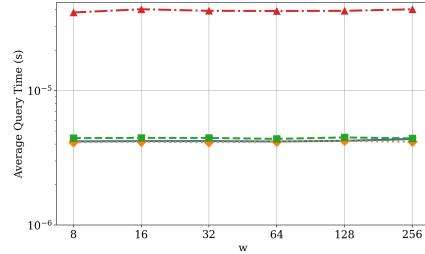
**Figure 73: effect of varying  $n_{g_1}$  on construction time, GOOGLE NGRAM,  $w = 65536, d = 5$ .**



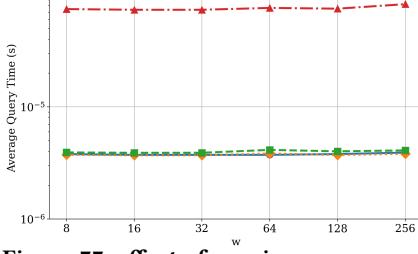
**Figure 74: effect of varying  $w$  on query time, NYC BUS,  $d = 5$ .**



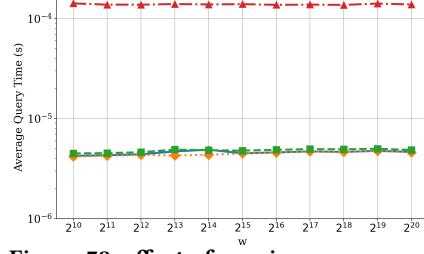
**Figure 75: effect of varying  $w$  on query time, DDOS,  $d = 5$ .**



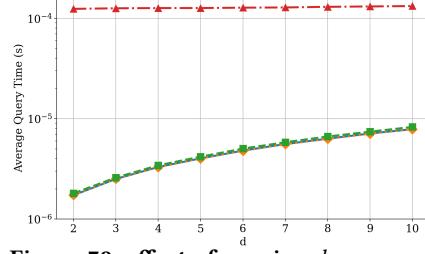
**Figure 76: effect of varying  $w$  on query time, REDDIT-TWITTER,  $d = 5$ .**



**Figure 77: effect of varying  $w$  on query time, CENSUS,  $d = 5$ .**

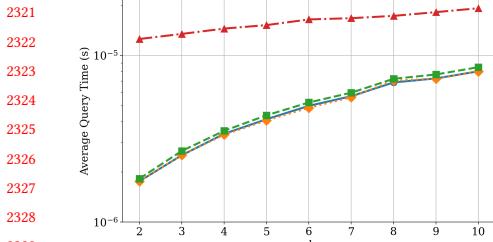
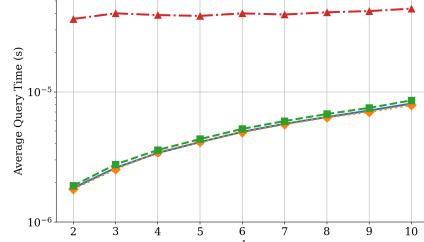
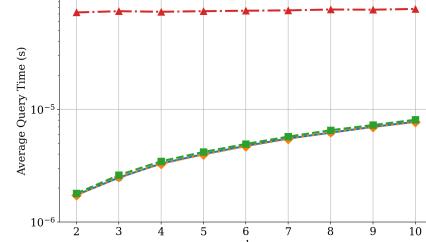
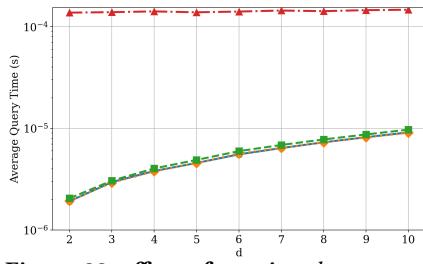
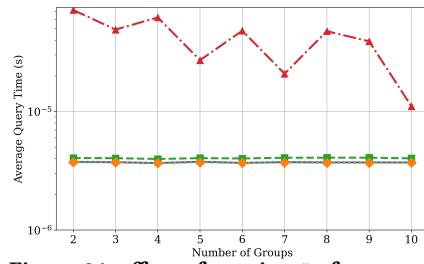
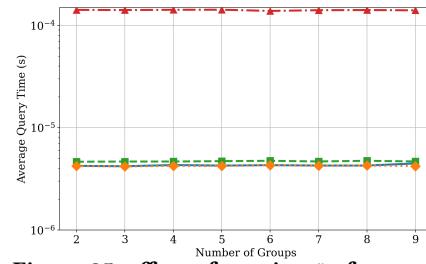
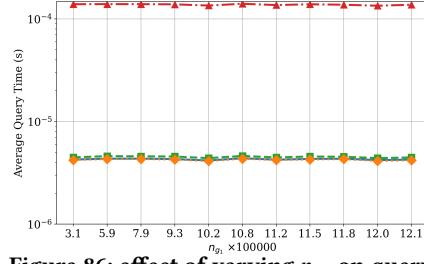


**Figure 78: effect of varying  $w$  on query time, GOOGLE NGRAM,  $d = 5$ .**



**Figure 79: effect of varying  $d$  on query time, NYC BUS,  $w = 512$ .**

2205  
2206  
2207  
2208  
2209  
2210  
2211  
2212  
2213  
2214  
2215  
2216  
2217  
2218  
2219  
2220  
2221  
2222  
2223  
2224  
2225  
2226  
2227  
2228  
2229  
2230  
2231  
2232  
2233  
2234  
2235  
2236  
2237  
2238  
2239  
2240  
2241  
2242  
2243  
2244  
2245  
2246  
2247  
2248  
2249  
2250  
2251  
2252  
2253  
2254  
2255  
2256  
2257  
2258  
2259  
2260  
2261  
2262  
2263  
2264  
2265  
2266  
2267  
2268  
2269  
2270  
2271  
2272  
2273  
2274  
2275  
2276  
2277  
2278  
2279  
2280  
2281  
2282  
2283  
2284  
2285  
2286  
2287  
2288  
2289  
2290  
2291  
2292  
2293  
2294  
2295  
2296  
2297  
2298  
2299  
2300  
2301  
2302  
2303  
2304  
2305  
2306  
2307  
2308  
2309  
2310  
2311  
2312  
2313  
2314  
2315  
2316  
2317  
2318  
2319  
2320

Figure 80: effect of varying  $d$  on query time, DDOS,  $w = 64$ .Figure 81: effect of varying  $d$  on query time, REDDIT-TWITTER,  $w = 512$ .Figure 82: effect of varying  $d$  on query time, CENSUS,  $w = 64$ .Figure 83: effect of varying  $d$  on query time, GOOGLE NGRAM,  $w = 65536$ .Figure 84: effect of varying # of groups  $\ell$  on query time, CENSUS,  $w = 64, d = 5$ .Figure 85: effect of varying # of groups  $\ell$  on query time, GOOGLE NGRAM,  $w = 65536, d = 5$ .Figure 86: effect of varying  $n_{g_1}$  on query time, GOOGLE NGRAM,  $w = 65536, d = 5$ .

2340  
2341  
2342  
2343  
2344  
2345  
2346  
2347  
2348  
2349  
2350  
2351  
2352  
2353  
2354  
2355  
2356  
2357  
2358  
2359  
2360  
2361  
2362  
2363  
2364  
2365  
2366  
2367  
2368  
2369  
2370  
2371  
2372  
2373  
2374  
2375  
2376  
2377  
2378

2379  
2380  
2381  
2382  
2383  
2384  
2385  
2386  
2387  
2388  
2389  
2390  
2391  
2392  
2393  
2394  
2395  
2396  
2397  
2398  
2399  
2400  
2401  
2402  
2403  
2404  
2405  
2406  
2407  
2408  
2409  
2410  
2411  
2412  
2413  
2414  
2415  
2416  
2417  
2418  
2419  
2420  
2421  
2422  
2423  
2424  
2425  
2426  
2427  
2428  
2429  
2430  
2431  
2432  
2433  
2434  
2435  
2436