# A. Omitted Proofs

## A.1. Proof of Theorem 3.1

Recall the statement of Theorem 3.1.

**Theorem.** If $\|\boldsymbol{f}'\|_2 < \alpha \|\boldsymbol{f}'\|_1$, the expected error of the count-min sketch (with one row) lies between $(1 - \alpha^2) \|\boldsymbol{f}'\|_1^2 / m$ and $\|\boldsymbol{f}'\|_1^2 / m$.

The expected error is

$$\mathbb{E}\left[\sum_{i=1}^n f_i' R_{f',i}\right] = \sum_{i=1}^n f_i' \mathbb{E}[R_{f',i}] = \sum_{i=1}^n f_i' \frac{\|\boldsymbol{f}'\|_1 - f_i'}{m}$$

$$= \frac{\|\boldsymbol{f}'\|_1^2 - \|\boldsymbol{f}'\|_2^2}{m}.$$

The result follows.

## A.2. Proof of Theorem 3.2

Recall the statement of Theorem 3.2.

**Theorem.** Let $A$ and $B$ have domination number at most $d$ with respect to each other, and let $p = (1 - 1/w)^d$.

$$\mathbb{E}[R_A^\ell] - \mathbb{E}[R_B^\ell] \leq (1 - p)\left(\sum_{i=1}^{\ell-1} \binom{\ell}{i} p^i \right.$$

$$(1 - p)^{\ell-1-i}\left(\mathbb{E}[R_B^i] - \mathbb{E}[R_B^\ell]\right) + (1 - p)^{\ell-1}$$

$$\left.\left(\mathbb{E}[R_B^1] + \frac{1/w}{1 - (1 - 1/w)^{|Q|}} S(Q) - \mathbb{E}[R_B^\ell]\right)\right)$$

where $Q$ is the multiset of keys that should be removed from $A$ in order for $B$ to pointwise dominate $A$ and $S(Q)$ is the sum of the elements of $Q$.

**Lemma A.1.** *Let $S = L \cup H$, where $L \cap H = \emptyset$.*

*Letting $p' = \mathbb{P}(R_H = 0) = (1 - 1/w)^{|H|}$, we find*

$$\mathbb{E}[\min(R_S^{(1)}, R_S^{(2)}, \ldots R_S^{(\ell)})] \leq$$

$$\sum_{i=1}^\ell \binom{\ell}{i} p'^i (1 - p')^{\ell-i} \mathbb{E}[\min(R_L^{(1)}, \ldots R_L^{(i)})] +$$

$$(1 - p')^\ell \left(\mathbb{E}[R_L] + \frac{1/w}{1 - (1 - 1/w)^{|H|}} S(H)\right)$$

*where $S(H)$ is the sum of the values of $H$.*

*Proof.* The idea is to use the law of total probability and condition on how many $R_H$ are zero. Letting $N_0$ be the number of $R_H$ equal to zero, we can write

$$\mathbb{E}[\min(R_S^{(1)}, R_S^{(2)}, \ldots R_S^{(\ell)})]$$

$$= \sum_{i=0}^\ell \binom{\ell}{i} p'^i (1 - p')^{\ell-i} \mathbb{E}[\min(R_S^{(1)}, R_S^{(2)}, \ldots R_S^{(\ell)}) | N_0 = i]$$

If $N_0 = i > 0$ (WLOG $R_H^{(1)}, R_H^{(2)}, \ldots R_H^{(i)} = 0$), then

$$\mathbb{E}[\min(R_S^{(1)}, \ldots R_S^{(\ell)}) | R_H^{(1)} = 0, \ldots, R_H^{(i)} = 0]$$

$$\leq \mathbb{E}[\min(R_S^{(1)}, \ldots R_S^{(i)}) | R_H^{(1)} = 0, \ldots, R_H^{(i)} = 0]$$

$$= \mathbb{E}[\min(R_L^{(1)}, \ldots R_L^{(i)})]$$

which addresses the case of $N_0 > 0$.

Now suppose that $N_0 = 0$, or that all of $R_H$ are greater than zero.

Then we can write

$$\mathbb{E}[\min(R_S^{(1)}, \ldots R_S^{(\ell)}) | R_H^{(1)} \neq 0, \ldots, R_H^{(\ell)} \neq 0]$$

$$\leq \mathbb{E}[R_S^{(1)} | R_H^{(1)} \neq 0, R_H^{(2)} \neq 0, \ldots R_H^\ell \neq 0]$$

$$= \mathbb{E}[R_S^{(1)} | R_H^{(1)} \neq 0]$$

$$= \mathbb{E}[R_L^{(1)} | R_H^{(1)} \neq 0] + \mathbb{E}[R_H^{(1)} | R_H^{(1)} \neq 0]$$

$$= \mathbb{E}[R_L] + \mathbb{E}[R_H]/\mathbb{P}(R_H \neq 0)$$

which simplifies to the desired expression, as $\mathbb{E}[R_H] = S(H)/w$ and $\mathbb{P}(R_H \neq 0) = 1 - (1 - 1/w)^{|H|}$. $\qquad\square$

Now, write $A = Q \cup A'$, where $A'$ is point-wise dominated by $B$. Then, by using the equation above we have

$$\mathbb{E}[R_A^\ell] \leq \sum_{i=1}^\ell \binom{\ell}{i} p'^i (1 - p')^{\ell-i} \mathbb{E}[R_{A'}^i]$$

$$+ (1 - p')^\ell \left(\mathbb{E}[R_{A'}] + \frac{1/w}{1 - (1 - 1/w)^{|Q|}} S(Q)\right).$$

where $p' = (1 - 1/w)^{|Q|}$.

Since $A'$ is pointwise dominated by $B$, we have that

$$\mathbb{E}[R_A^\ell] \leq \sum_{i=1}^\ell \binom{\ell}{i} p'^i (1 - p')^{\ell-i} \mathbb{E}[R_B^i]$$

$$+ (1 - p')^\ell \left(\mathbb{E}[R_B] + \frac{1/w}{1 - (1 - 1/w)^{|Q|}} S(Q)\right).$$

where $p' = (1 - 1/w)^{|Q|}$. It remains to show that we can replace $p'$ with $p = (1 - 1/w)^d \leq p'$. To see this, let $c(i) = \mathbb{E}[R_B^i]$ for $i \neq 0$ and $c(i) = \mathbb{E}[R_B] + \frac{1/w}{1-(1-1/w)^{|Q|}} S(Q)$. It is straightforward to see that $c(i)$ is monotonically decreasing. Also note that the RHS can be expressed as $\mathbb{E}_{i \sim \text{Bin}(\ell, p')}[c(i)]$. Since $\text{Bin}(\ell, p)$ is stochastically dominated by $\text{Bin}(\ell, p')$ (as $p' > p$), we have that $\mathbb{E}_{i \sim \text{Bin}(\ell, p')}[c(i)] \leq \mathbb{E}_{i \sim \text{Bin}(\ell, p)}[c(i)]$, which yields

$$\mathbb{E}[R_A^\ell] \leq \sum_{i=1}^\ell \binom{\ell}{i} p^i (1 - p)^{\ell-i} \mathbb{E}[R_{A'}^i]$$

$$+ (1 - p)^\ell \left(\mathbb{E}[R_{A'}] + \frac{1/w}{1 - (1 - 1/w)^{|Q|}} S(Q)\right).$$

and the result now follows by subtracting $\mathbb{E}[R_B^\ell]$ from each side.

### A.3. Proof of Theorem 3.3

Recall the statement of Theorem 3.3.

**Theorem.** For each $p, c > 0$, given Zipfian input with parameter $p$, and the largest $n' = cn$ frequencies removed (say, by an oracle), there exists a constant $c' > 0$ such that a Count-Min Sketch with total space at most $c'n$ has larger expected error with $k$ rows than 1 row for all $k \geq 2$.

To prove this, we invoke a slightly stronger result.

**Theorem A.2.** *Assume that $f$ consists of $z$ keys whose frequencies are bounded between $x$ and $Cx$ for some constant $C > 1$ and some $x > 0$, and $m \leq \frac{z}{36C}$. Then, $\mathbb{E}[R_f^\ell]$ is minimized at $\ell = 1$ when $m = w\ell$ is held constant.*

*Proof.* The key idea is to bound the standard deviation-to-mean ratio of $R_f^1$. This ratio is bounded as follows.

$$
\begin{aligned}
\frac{\sqrt{\sum f_i^2 (1/m)(1 - 1/m)}}{\sum f_i (1/m)} &\leq \sqrt{m} \frac{\sqrt{\sum f_i^2}}{\sum f_i} \\
&\leq \sqrt{m} \frac{\sqrt{Cx \sum f_i}}{\sum f_i} \\
&= \sqrt{m} \frac{\sqrt{Cx}}{\sqrt{\sum f_i}} \\
&\leq \sqrt{\frac{Cm}{z}} \leq \frac{1}{6}.
\end{aligned}
$$

where the last step follows because $\sum f_i \geq zx$.

In other words, when the width is considerably smaller than the number of keys, the ratio of the standard deviation to the mean is low.

We now use the following lemma.

**Lemma A.3.** *If $X_1, X_2, \ldots X_k$ are nonnegative i.i.d. random variables with mean $\mu$ and variance $\sigma^2$, we have that*

$$
\mathbb{E}[\min(X_1, X_2, \ldots X_k)] \leq \frac{3}{4}(\mu - 2\sigma\sqrt{k}).
$$

*Proof.* By Chebyshev's inequality, the probability that $X_i > \mu - 2\sigma\sqrt{k}$ is at most $1/4k$. Thus, by a union bound the probability that $X_i < \mu - 2\sigma\sqrt{k}$ for at least one value of $i$ is at most $1/4$ and thus the expected value (remembering that the variable is nonnegative) is at least $3/4(\mu - 2\sigma\sqrt{k})$. $\square$

Now, in terms of the number of rows $\ell$, note that the variance of $R_f$ (keeping in mind that $R_f$ depends on $\ell$) is $\sum f_i^2/w = \sum f_i^2 \ell/m = \ell\sigma_1^2$, and the mean of $R_f$ is $\sum f_i/w = \sum f_i \ell/m = \ell\mu_1$, where $\mu_1$ and $\sigma_1$ are the

mean and standard deviation of $R_f^\ell$ when $\ell = 1$. Thus, the expected value of $R_f^\ell$ is at least

$$
\frac{3(\ell\mu_1 - 2\ell\sigma_1)}{4}.
$$

Note that $\mu_1 > 6\sigma_1$, so for $\ell > 2$, then this is minimized at $\ell = 2$.

Since $\mu_1 \geq 6\sigma_1$, we have that $\mu_1$ is less than $\frac{3}{4}(2\mu_1 - 4\sigma_1) = \frac{3\mu_1}{2} - 3\sigma_1$, which implies that $\ell = 1$ is optimal.

$\square$

In our setting, our input is a truncated Zipf distribution of parameter $p$ supported on $n' + 1$ to $n$ (as the $n'$ greatest frequencies are removed), in which the ratio between the largest and smallest frequencies is at most $(1/c)^p$, and $z = n - n' - 1 = n(1 - c) - 1$. Thus, as long as $m \leq \frac{n(1-c)-1}{36c^p}$, the desired conclusion follows.

### A.4. Proof of Corollary 6.1.1

Recall the statement of Corollary 6.1.1.

**Corollary.** The error of the Learned Count-Min Sketch under the oracle model from 6.1, with $B_r = c'B$ for some $0 < c' < 1$, and with Zipfian input with parameter 1, is $\Theta\left(\frac{\ln^2(n/B)}{B}\right)$.

*Proof.* The error of the Learned Count-Min Sketch with a perfect oracle is $\Omega\left(\frac{\ln^2(n/B)}{B}\right)$ by Theorem 10.4 in (Hsu et al., 2019). We note that the bound in 6.1 is minimized when $B_r$ and $B - B_r$ are both $\Theta(B)$, in which case the upper and lower bounds match asymptotically so the error is $\Theta\left(\frac{\ln^2(n/B)}{B}\right)$. $\square$

### A.5. Proof of Theorem 6.1

Recall the statement of Theorem 6.1.

**Theorem.** For any constant $c > 0$, if the input is Zipfian with parameter 1 and the heavy hitter oracle screens key $i$ with probability $p(f_i)$ where for $1 \leq i \leq (1 + 1/c)B_r$,

$$
p(f_i) = 1 - \left(\frac{B_r/c}{\sum_{j=1}^{B_r(1+1/c)} j^c}\right) i^c
$$

and $p(f_i) = 0$ otherwise, then the error for the Learned Count-Min Sketch is $O\left(\frac{1/c^2 + \ln^2(n/B_r)}{B - B_r}\right)$.

In this theorem, the term $\frac{B_r/c}{\sum_{j=1}^{B_r(1+1/c)} j^c}$ is a normalizing constant to make sure that we screen $B_r$ keys on expectation and also bound all probabilities between 0 and 1.

*Proof.* First, we check that $p(f_i)$ is a valid probability distribution.

Note that we can lower bound the denominator of the fraction in $p(f_i)$ with

$$\int_0^{B_r(1+1/c)} x^c dx = \frac{B_r}{c}\left(B_r(1+1/c)\right)^c,$$

so

$$p(f_i) \geq 1 - \frac{B_r/ci^c}{B_r/c(B_r(1+1/c))^c} \geq 0$$

since $i \leq B_r(1+1/c)$. Thus, this is a valid probability distribution, and the sum of the probabilities is $B_r(1+1/c) - B_r/c = B_r$ so this gives us the correct number of heavy hitters $B_r$.

We now note that we can upper bound the error of the Learned Count-Min Sketch by upper bounding the error of the Learned Count-Min Sketch with one row. Letting $X_1, X_2, \ldots X_n$ be independent Bernoulli random variables that are 1 with probability $\frac{1}{B - B_r}$,

$$\mathbb{E}\left[\sum_{i=1}^n f_i|\tilde{f}_i - f_i|\right]$$

$$\leq \sum_{i=1}^{B_r(1+1/c)} \left(\frac{B_r i^c}{\Omega(B_r^{c+1})}\right) f_i \mathbb{E}\left[\sum_{j=B_r(1+1/c)+1}^n f_j X_j\right.$$

$$\left. + \sum_{j=1}^{B_r(1+1/c)} f_j X_j \left(\frac{B_r j^c}{\Omega(B_r^{c+1})}\right)\right]$$

$$+ \sum_{i=B_r(1+1/c)+1}^n f_i \mathbb{E}\left[\sum_{j=1}^{B_r(1+1/c)} f_j X_j \left(\frac{B_r j^c}{\Omega(B_r^{c+1})}\right)\right.$$

$$\left. + \sum_{j=B_r(1+1/c)+1}^n f_j X_j\right]$$

$$= \sum_{i=1}^{B_r(1+1/c)} \left(\frac{i^{c-1}}{\Omega(B_r^c)}\right)\left[\sum_{j=1}^{B_r(1+1/c)} \frac{j^{c-1}}{\Omega(B_r^c)(B-B_r)}\right.$$

$$\left. + \sum_{j=B_r(1+1/c)+1}^n \frac{1}{j(B-B_r)}\right]$$

$$+ \sum_{i=B_r(1+1/c)+1}^n \frac{1}{i}\left[\sum_{j=1}^{B_r(1+1/c)} \frac{j^{c-1}}{\Omega(B_r^c)(B-B_r)}\right.$$

$$\left. + \sum_{j=B_r(1+1/c)+1}^n \frac{1}{j(B-B_r)}\right] \tag{5}$$

Now, note that

$$\sum_{i=1}^{B_r(1+1/c)} i^{c-1} \geq \int_0^{B_r(1+1/c)} x^{c-1} dx = \frac{(B_r(1+1/c))^c}{c}$$

so

$$\sum_{i=1}^{B_r(1+1/c)} i^{c-1} = O\left(\frac{B_r^c}{c}\right). \tag{6}$$

A similar bound with integrals tells us that

$$\sum_{j=B_r(1+1/c)+1}^n \frac{1}{j} \geq \int_{B_r(1+1/c)}^n \frac{1}{x} dx$$

$$= \int_0^n \frac{1}{x} dx - \int_0^{B_r(1+1/c)} \frac{1}{x} dx$$

$$= \ln(n) - \ln(B_r(1+1/c) + 1) = O\left(\ln(n/B_r)\right).$$

Thus,

$$\sum_{j=B_r(1+1/c)+1}^n \frac{1}{j} = O\left(\ln(n/B_r)\right). \tag{7}$$

Plugging (6) and (7) into Equation (5), we get

$$= O\left(\frac{1}{c}\right) O\left(\frac{1}{c(B-B_r)} + \frac{\ln(n/B_r)}{B-B_r}\right)$$

$$+ O(\ln(n/B_r)) O\left(\frac{1}{c(B-B_r)} + \frac{\ln(n/B_r)}{B-B_r}\right)$$

$$= O\left(\frac{1}{c^2(B-B_r)} + \frac{\ln(n/B_r)}{c(B-B_r)}\right)$$

$$+ O\left(\frac{\ln(n/B_r)}{c(B-B_r)} + \frac{\ln^2(n/B_r)}{B-B_r}\right)$$

$$= O\left(\frac{1/c^2 + \ln^2(n/B_r)}{B-B_r}\right).$$

Thus, we know that our error is

$$O\left(\frac{1/c^2 + \ln^2(n/B_r)}{B-B_r}\right).$$

$\square$

### A.6. Generalization of Theorem 6.1

The proof of this fact follows directly from plugging in $i^{-p}$ instead of $\frac{1}{i}$ in for the frequency $f_i$ in our proof in Appendix A.5.

Now, we will prove the similar result when the Zipf parameter $p \neq 1$:

*Proof.* When the Zipf parameter $p \neq 1$, the only difference is that the sums over the frequencies become

$$\sum_{i=1}^{B_r(1+1/c)} i^{c-p}$$

and

$$\sum_{j=B_r(1+1/c)+1}^{n} \frac{1}{j^p}$$

so when $c = p+1$ the first bound from the integrals becomes $O(\ln(B_r))$, and otherwise the bounds are

$$O\left(\frac{B_r^{c-p+1}}{c-p+1}\right)$$

and

$$O\left(\frac{n^{1-p} - B_r^{1-p}}{1-p}\right) = O\left(\frac{n^{1-p}}{1-p}\right)$$

respectively. Plugging in our new probabilities and frequencies, equation (6) becomes

$$= \sum_{i=1}^{B_r(1+1/c)} \left(\frac{i^{c-p}}{\Omega(B_r^c)}\right) \left[\sum_{j=1}^{B_r(1+1/c)} \frac{j^{c-p}}{\Omega(B_r^c)(B-B_r)}\right.$$
$$+ \sum_{j=B_r(1+1/c)+1}^{n} \frac{1}{j^p(B-B_r)}\right]$$
$$+ \sum_{i=B_r(1+1/c)+1}^{n} \frac{1}{i^p}\left[\sum_{j=1}^{B_r(1+1/c)} \frac{j^{c-p}}{\Omega(B_r^c)(B-B_r)}\right.$$
$$+ \sum_{j=B_r(1+1/c)+1}^{n} \frac{1}{j^p(B-B_r)}\right],$$

so now when we use our integral approximations, we get:

When $p = c + 1$:

$$= O\left(\frac{\ln(B_r)}{B_r^c}\right) O\left(\frac{\ln(B_r)}{B_r^c(B-B_r)} + \frac{n^{1-p}}{(1-p)(B-B_r)}\right)$$
$$+ O\left(\frac{n^{1-p}}{1-p}\right) O\left(\frac{\ln(B_r)}{B_r^c(B-B_r)} + \frac{n^{1-p}}{(1-p)(B-B_r)}\right)$$
$$= O\left(\frac{\ln^2(B_r)}{B_r^{2c}(B-B_r)} + \frac{\ln(B_r)n^{1-p}}{(1-p)(B-B_r)B_r^c}\right)$$
$$+ O\left(\frac{\ln(B_r)n^{1-p}}{(1-p)B_r^c(B-B_r)} + \frac{n^{2-2p}}{B-B_r(1-p)^2}\right)$$
$$= O\left(\frac{\ln^2(B_r) + 1/c^2}{B_r^{2c}(B-B_r)}\right)$$

since we note that $B_r \le n$ so $\frac{1}{B_r} > \frac{1}{n}$. For all other $p$, it is exactly the same except the $\ln(B_r)$ terms get replaced by $\frac{B_r^{c-p+1}}{c-p+1}$ so making that replacement, we get

$$O\left(\frac{B_r^{2-2p}/(c-p+1)^2}{B-B_r} + \frac{n^{2-2p}/(1-p)^2}{B-B_r}\right).$$

$\square$

## A.7. Discussion of the Bounds in Section 6

When we assumed the frequency of the $i^{th}$ most frequent item was $i^{-p}$, the total frequency of the input $\|\mathbf{f}\|_1$ became $\Theta\left(1 + \frac{n^{1-p}}{1-p}\right)$. Thus, if we normalize the inputs by dividing the frequency of each key by $1 + \frac{n^{1-p}}{1-p}$ (or equivalently dividing the error by the square of this) and let $B_r = \Theta(n)$ as is often the case, simply comparing terms allows us to see that the error bound in (4) is

$$O\left(\frac{n^{2-2p}/(1-p)^2}{B-B_r}\right),$$

then when $p < 1$, $\|\mathbf{f}\|_1$ is dominated by $\frac{n^{1-p}}{1-p}$ so the normalized error we get is

$$O\left(\frac{1}{B-B_r}\right).$$

On the other hand, when $p > 1$, $\|\mathbf{f}\|_1$ is instead dominated by 1 so the normalized error we get is

$$O\left(\frac{n^{2-2p}/(1-p)^2}{B-B_r}\right).$$

This result makes sense as when the input gets more skewed, there is more weight in the removed heavy hitters so we get less error.

# B. AOL Dataset Results

We show here the results on the AOL dataset. The AOL dataset is very small, so both methods being ended up predicting similar sets of heavy hitters, and the performance of our method and theirs are very close.



Day 50, BatchRank (K = 8), width 100, Count-Min



Day 50, BatchRank (K = 8), width 300, Count-Min



Day 50, BatchRank (K = 8), width 1000, Count-Min



Day 50, BatchRank (K = 8), width 3000, Count-Min



Day 50, BatchRank (K = 8), width 10000, Count-Min



Day 50, BatchRank (K = 8), width 100, Count-Sketch



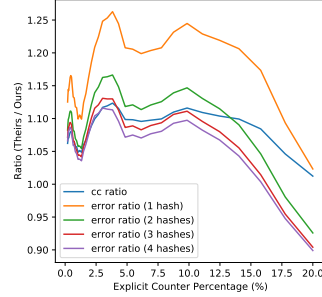Day 50, BatchRank (K = 8), width 300, Count-Sketch



Day 50, BatchRank (K = 8), width 1000, Count-Sketch



Day 50, BatchRank (K = 8), width 3000, Count-Sketch



Day 50, BatchRank (K = 8), width 10000, Count-Sketch

# C. All Tables and Graphs for Section 5

In this section, we show all the graphs for each minute's results for each of the loss functions in table 4. As expected from the results in the table, BatchRank with K=64 performs very well when the number of explicit counters is small, and then its performance falls off. Unweighted L1 Loss performs poorly, and Weighted Log Loss performs comparably to BatchRank with K=8 when the explicit counter percentage is sufficiently large. The results for BatchRank with K=8 are similar for each of the minutes shown here to the results we saw earlier in figure 1.

## C.1. Minute 9 results



Minute 9, BatchRank (K = 8), width 1000, Count-Min



Minute 9, BatchRank (K = 8), width 3000, Count-Min



Minute 9, BatchRank (K = 8), width 10000, Count-Min



Minute 9, BatchRank (K = 8), width 30000, Count-Min



Minute 9, BatchRank (K = 8), width 100000, Count-Min

Minute 9, BatchRank (K = 8), width 1000, Count-Sketch
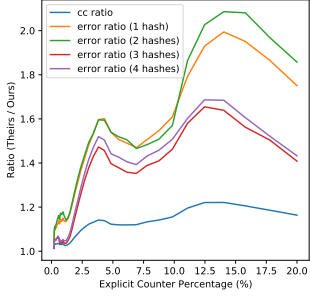
Minute 9, BatchRank (K = 64), width 1000, Count-Min
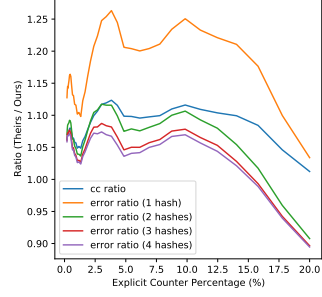
Minute 9, BatchRank (K = 8), width 3000, Count-Sketch
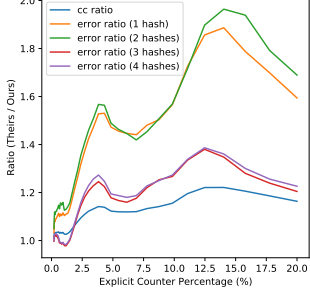
Minute 9, BatchRank (K = 64), width 3000, Count-Min
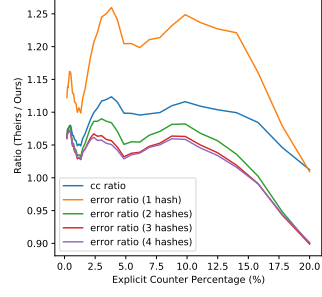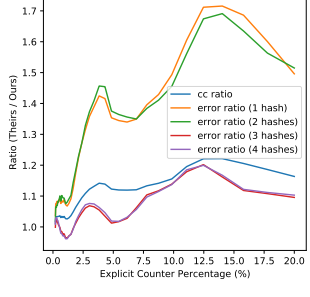
Minute 9, BatchRank (K = 8), width 10000, Count-Sketch

Minute 9, BatchRank (K = 64), width 10000, Count-Min
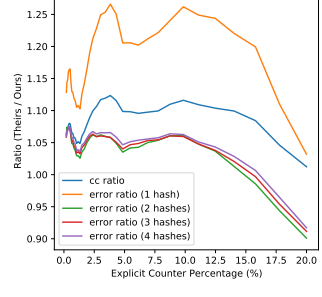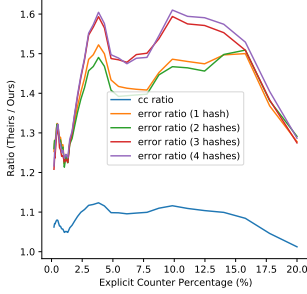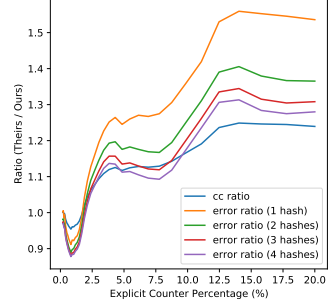
Minute 9, BatchRank (K = 8), width 30000, Count-Sketch

Minute 9, BatchRank (K = 64), width 30000, Count-Min

Minute 9, BatchRank (K = 8), width 100000, Count-Sketch
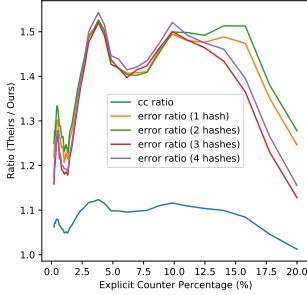
Minute 9, BatchRank (K = 64), width 100000, Count-Min
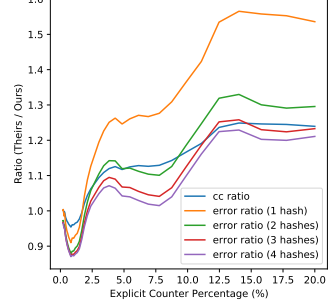
Minute 9, BatchRank (K = 64), width 1000, Count-Sketch
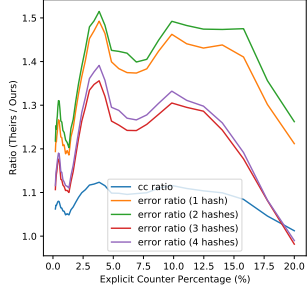
Minute 9, Weighted Log Loss, width 1000, Count-Min
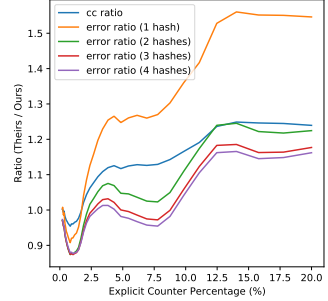
Minute 9, BatchRank (K = 64), width 3000, Count-Sketch
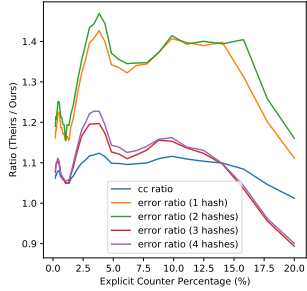
Minute 9, Weighted Log Loss, width 3000, Count-Min
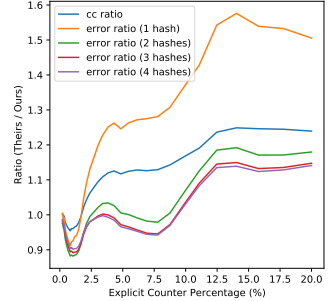
Minute 9, BatchRank (K = 64), width 10000, Count-Sketch
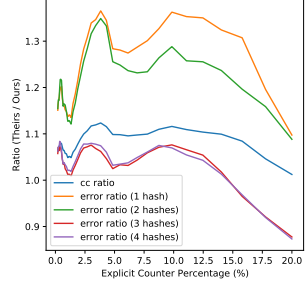
Minute 9, Weighted Log Loss, width 10000, Count-Min

Minute 9, BatchRank (K = 64), width 30000, Count-Sketch
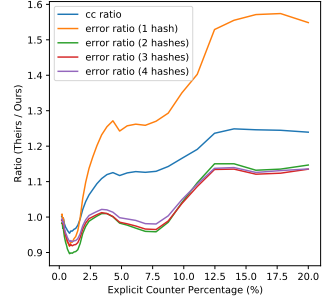
Minute 9, Weighted Log Loss, width 30000, Count-Min
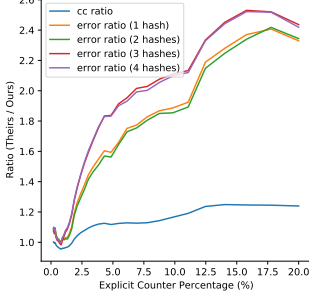
Minute 9, BatchRank (K = 64), width 100000, Count-Sketch
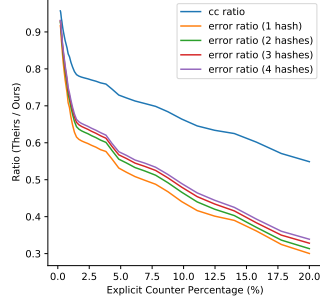
Minute 9, Weighted Log Loss, width 100000, Count-Min

Minute 9, Weighted Log Loss, width 1000, Count-Sketch

Minute 9, Unweighted L1 Loss, width 1000, Count-Min

Minute 9, Weighted Log Loss, width 3000, Count-Sketch

Minute 9, Unweighted L1 Loss, width 3000, Count-Min

Minute 9, Weighted Log Loss, width 10000, Count-Sketch

Minute 9, Unweighted L1 Loss, width 10000, Count-Min

Minute 9, Weighted Log Loss, width 30000, Count-Sketch
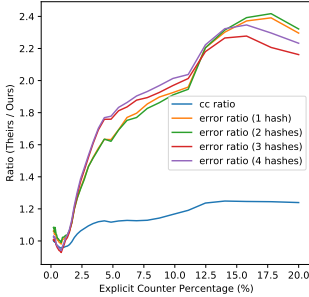
Minute 9, Unweighted L1 Loss, width 30000, Count-Min
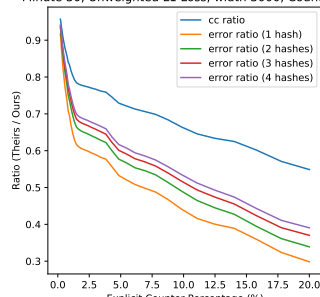
Minute 9, Weighted Log Loss, width 100000, Count-Sketch

Minute 9, Unweighted L1 Loss, width 100000, Count-Min

Minute 9, Unweighted L1 Loss, width 1000, Count-Sketch
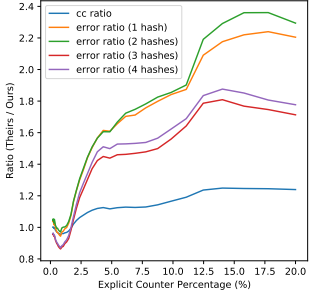
Minute 9, Weighted L1 Loss, width 1000, Count-Min
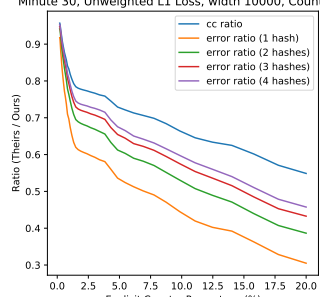
Minute 9, Unweighted L1 Loss, width 3000, Count-Sketch

Minute 9, Weighted L1 Loss, width 3000, Count-Min

Minute 9, Unweighted L1 Loss, width 10000, Count-Sketch

Minute 9, Weighted L1 Loss, width 10000, Count-Min

Minute 9, Unweighted L1 Loss, width 30000, Count-Sketch
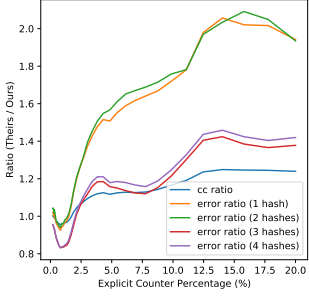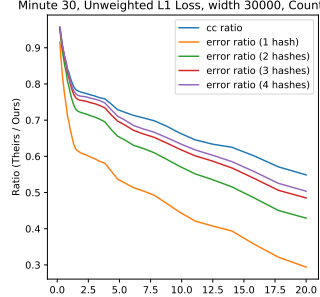
Minute 9, Weighted L1 Loss, width 30000, Count-Min

Minute 9, Unweighted L1 Loss, width 100000, Count-Sketch

Minute 9, Weighted L1 Loss, width 100000, Count-Min
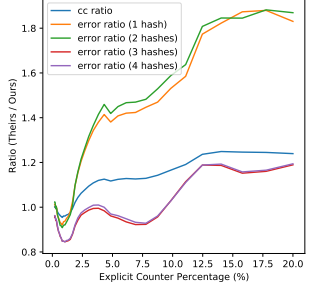
Minute 9, Weighted L1 Loss, width 1000, Count-Sketch

Minute 9, Weighted L1 Loss, width 3000, Count-Sketch

Minute 9, Weighted L1 Loss, width 10000, Count-Sketch

Minute 9, Weighted L1 Loss, width 30000, Count-Sketch

Minute 9, Weighted L1 Loss, width 100000, Count-Sketch
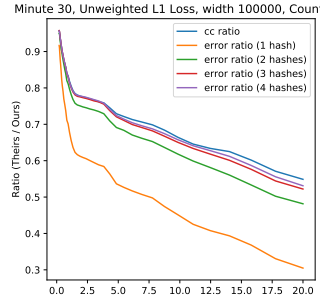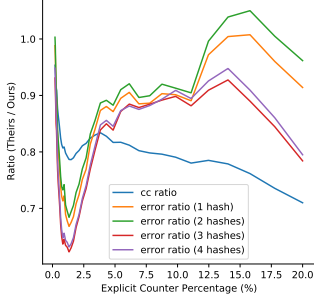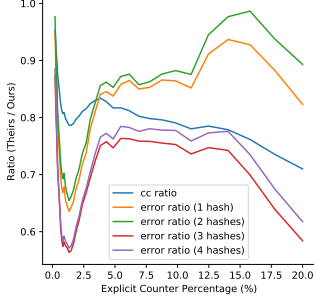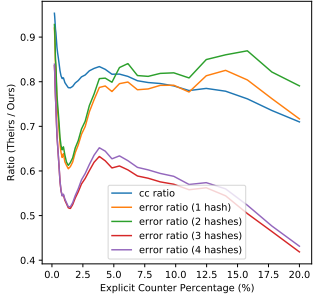
## C.2. Minute 30 results

Minute 30, BatchRank (K = 8), width 1000, Count-Min

Minute 30, BatchRank (K = 8), width 3000, Count-Min

Minute 30, BatchRank (K = 8), width 10000, Count-Min

Minute 30, BatchRank (K = 8), width 30000, Count-Min

Minute 30, BatchRank (K = 8), width 100000, Count-Min

Minute 30, BatchRank (K = 8), width 1000, Count-Sketch

Minute 30, BatchRank (K = 64), width 1000, Count-Min

Minute 30, BatchRank (K = 8), width 3000, Count-Sketch

Minute 30, BatchRank (K = 64), width 3000, Count-Min

Minute 30, BatchRank (K = 8), width 10000, Count-Sketch

Minute 30, BatchRank (K = 64), width 10000, Count-Min

Minute 30, BatchRank (K = 8), width 30000, Count-Sketch

Minute 30, BatchRank (K = 64), width 30000, Count-Min

Minute 30, BatchRank (K = 8), width 100000, Count-Sketch

Minute 30, BatchRank (K = 64), width 100000, Count-Min

Minute 30, BatchRank (K = 64), width 1000, Count-Sketch

Minute 30, Weighted Log Loss, width 1000, Count-Min

Minute 30, BatchRank (K = 64), width 3000, Count-Sketch

Minute 30, Weighted Log Loss, width 3000, Count-Min

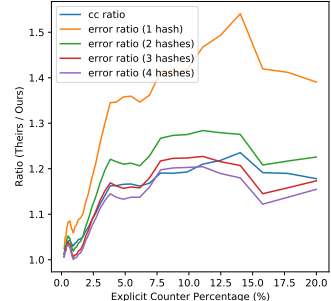Minute 30, BatchRank (K = 64), width 10000, Count-Sketch

Minute 30, Weighted Log Loss, width 10000, Count-Min

Minute 30, BatchRank (K = 64), width 30000, Count-Sketch

Minute 30, Weighted Log Loss, width 30000, Count-Min

Minute 30, BatchRank (K = 64), width 100000, Count-Sketch

Minute 30, Weighted Log Loss, width 100000, Count-Min

Minute 30, Weighted Log Loss, width 1000, Count-Sketch

Minute 30, Unweighted L1 Loss, width 1000, Count-Min

Minute 30, Weighted Log Loss, width 3000, Count-Sketch

Minute 30, Unweighted L1 Loss, width 3000, Count-Min

Minute 30, Weighted Log Loss, width 10000, Count-Sketch

Minute 30, Unweighted L1 Loss, width 10000, Count-Min

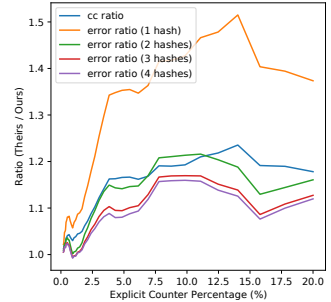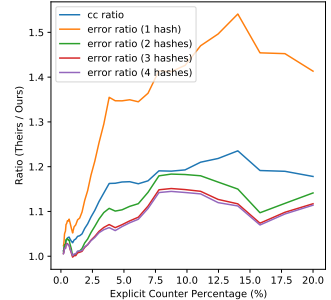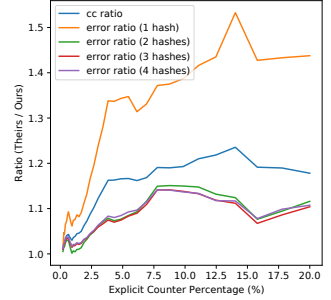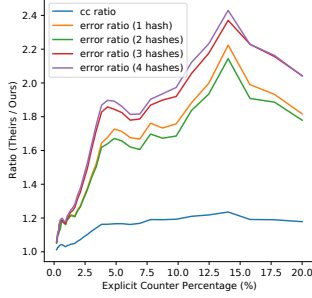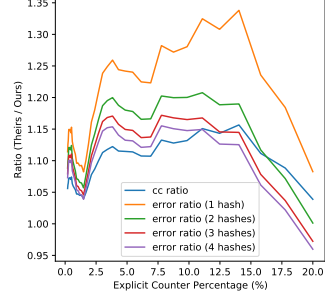Minute 30, Weighted Log Loss, width 30000, Count-Sketch

Minute 30, Unweighted L1 Loss, width 30000, Count-Min

Minute 30, Weighted Log Loss, width 100000, Count-Sketch

Minute 30, Unweighted L1 Loss, width 100000, Count-Min

Minute 30, Unweighted L1 Loss, width 1000, Count-Sketch

Minute 30, Weighted L1 Loss, width 1000, Count-Min

Minute 30, Unweighted L1 Loss, width 3000, Count-Sketch

Minute 30, Weighted L1 Loss, width 3000, Count-Min

Minute 30, Unweighted L1 Loss, width 10000, Count-Sketch

Minute 30, Weighted L1 Loss, width 10000, Count-Min

Minute 30, Unweighted L1 Loss, width 30000, Count-Sketch

Minute 30, Weighted L1 Loss, width 30000, Count-Min

Minute 30, Unweighted L1 Loss, width 100000, Count-Sketch

Minute 30, Weighted L1 Loss, width 100000, Count-Min

Minute 30, Weighted L1 Loss, width 1000, Count-Sketch



Minute 30, Weighted L1 Loss, width 3000, Count-Sketch



Minute 30, Weighted L1 Loss, width 10000, Count-Sketch



Minute 30, Weighted L1 Loss, width 30000, Count-Sketch



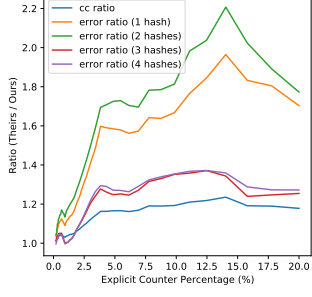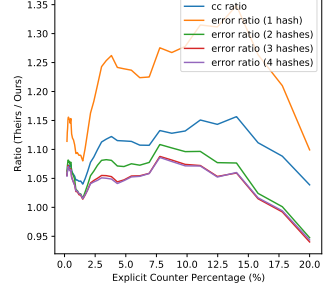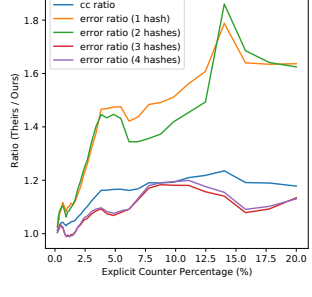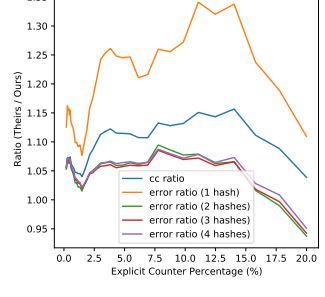Minute 30, Weighted L1 Loss, width 100000, Count-Sketch



## C.3. Minute 60 results

Minute 60, BatchRank (K = 8), width 1000, Count-Min



Minute 60, BatchRank (K = 8), width 3000, Count-Min



Minute 60, BatchRank (K = 8), width 10000, Count-Min



Minute 60, BatchRank (K = 8), width 30000, Count-Min
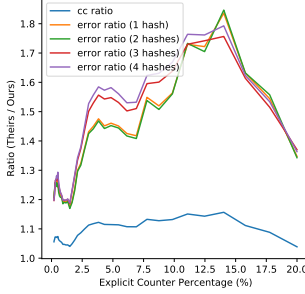


Minute 60, BatchRank (K = 8), width 100000, Count-Min
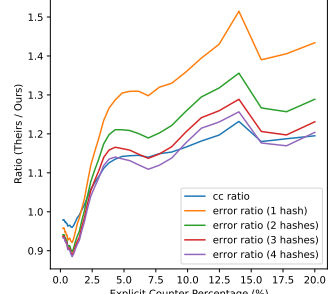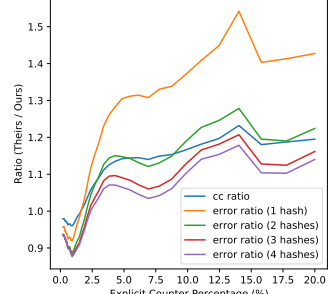
Minute 60, Weighted Log Loss, width 1000, Count-Sketch

Minute 60, Unweighted L1 Loss, width 1000, Count-Min

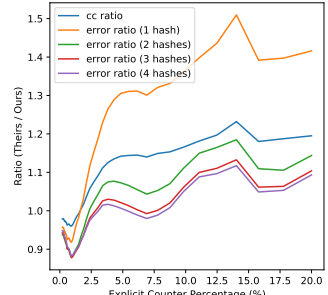Minute 60, Weighted Log Loss, width 3000, Count-Sketch

Minute 60, Unweighted L1 Loss, width 3000, Count-Min

Minute 60, Weighted Log Loss, width 10000, Count-Sketch

Minute 60, Unweighted L1 Loss, width 10000, Count-Min

Minute 60, Weighted Log Loss, width 30000, Count-Sketch

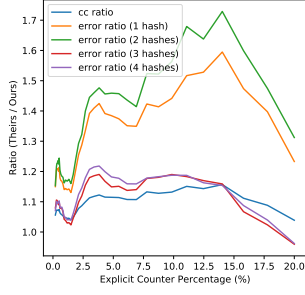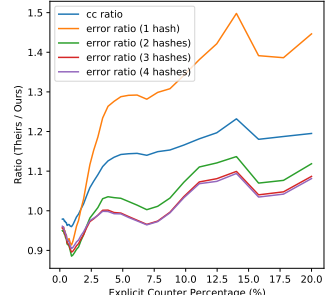Minute 60, Unweighted L1 Loss, width 30000, Count-Min

Minute 60, Weighted Log Loss, width 100000, Count-Sketch
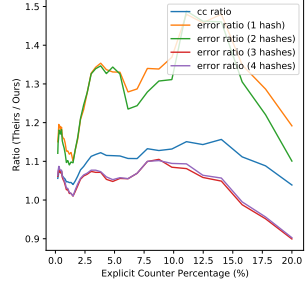
Minute 60, Unweighted L1 Loss, width 100000, Count-Min

Minute 60, Unweighted L1 Loss, width 1000, Count-Sketch

Minute 60, Weighted L1 Loss, width 1000, Count-Min
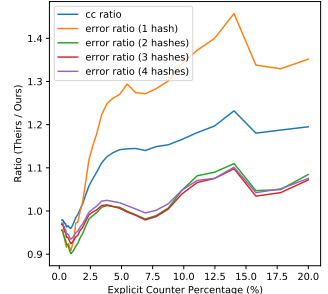
Minute 60, Unweighted L1 Loss, width 3000, Count-Sketch

Minute 60, Weighted L1 Loss, width 3000, Count-Min

Minute 60, Unweighted L1 Loss, width 10000, Count-Sketch

Minute 60, Weighted L1 Loss, width 10000, Count-Min

Minute 60, Unweighted L1 Loss, width 30000, Count-Sketch

Minute 60, Weighted L1 Loss, width 30000, Count-Min

Minute 60, Unweighted L1 Loss, width 100000, Count-Sketch
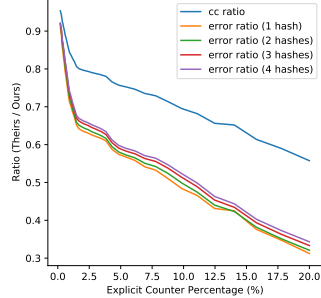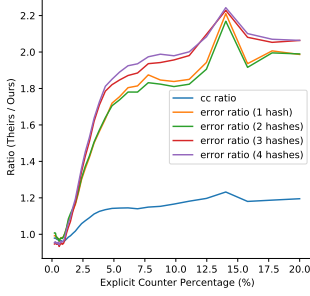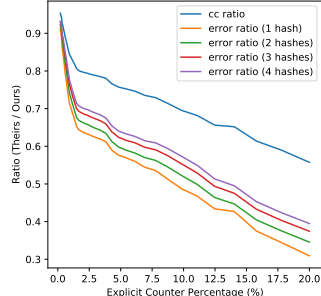
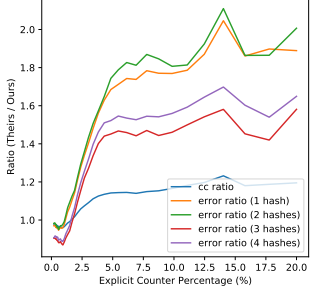Minute 60, Weighted L1 Loss, width 100000, Count-Min

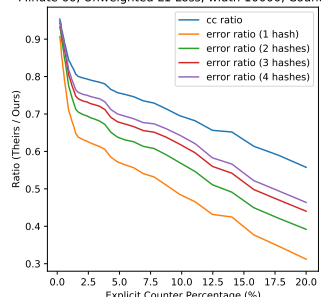Minute 60, Weighted L1 Loss, width 1000, Count-Sketch



Minute 60, Weighted L1 Loss, width 3000, Count-Sketch
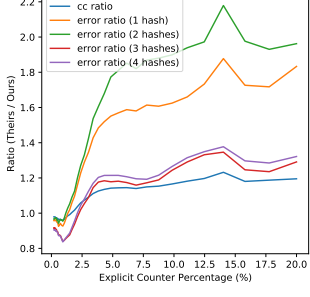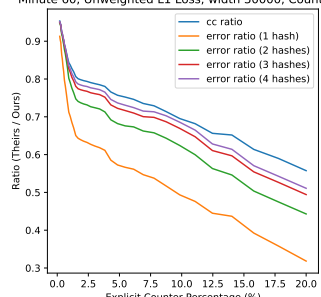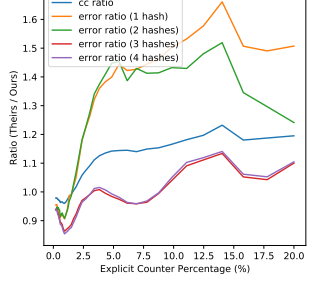


Minute 60, Weighted L1 Loss, width 10000, Count-Sketch



Minute 60, Weighted L1 Loss, width 30000, Count-Sketch
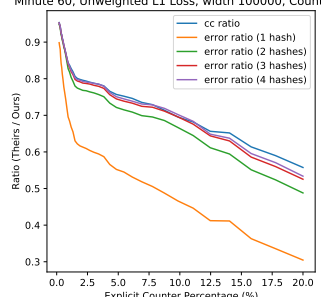


Minute 60, Weighted L1 Loss, width 100000, Count-Sketch

## D. All Screened Rates

Figures 3 and 4 show the screened rates of all the different algorithms on the 30th and 60th minute, arranged in the same order as in Table 1. They all exhibit the same behavior, where they are near 0% and all of a sudden begin increasing as the keys get heavier.



(a) Unweighted Log Loss



(b) Weighted Log Loss



(c) Unweighted $L^1$ Loss



(d) Weighted $L^1$ Loss



(e) BatchRank, $K = 64$



(f) BatchRank, $K = 8$

*Figure 3.* Screening Rates on Minute 30

(a) Unweighted Log Loss



(b) Weighted Log Loss



(c) Unweighted $L^1$ Loss



(d) Weighted $L^1$ Loss



(e) BatchRank, $K = 64$



(f) BatchRank, $K = 8$

*Figure 4.* Screening Rates on Minute 60

# E. Additional Runs

In this section, we show the coverage results of running our models five times in Tables 5, 6, 7. This verifies that the behavior seen in the run in Table 4 is consistent over multiple trials.

As a reminder, unweighted log loss is the original from (Hsu et al., 2019) and the other methods are from this paper. As we are now working with multiple trials and have standard errors, we bold the largest entry as well as all entries for which the largest entry is within its standard error.

*Table 5.* Coverage (%) on CAIDA dataset, 9th minute. The largest entries for each column are in **bold**.

| METHOD | COVERAGE SIZE | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | 1% | 2% | 5% | 10% | 20% | 30% | 50% | 75% |
| UNWEIGHTED LOG LOSS | 37.4% ($\pm$0.7%) | 45.6% ($\pm$0.7%) | 53.2% ($\pm$0.7%) | 64.5% ($\pm$0.4%) | 76.9% ($\pm$0.5%) | 83.9% ($\pm$0.4%) | 91.9% ($\pm$0.4%) | **98.6%** ($\pm$0.1%) |
| WEIGHTED LOG LOSS | 39.7% ($\pm$1.1%) | 49.7% ($\pm$0.9%) | **61.8%** ($\pm$0.5%) | **70.1%** ($\pm$0.5%) | **81.8%** ($\pm$0.2%) | **87.1%** ($\pm$0.3%) | **93.8%** ($\pm$0.2%) | **98.7%** ($\pm$0.0%) |
| UNWEIGHTED $L^1$ LOSS | 22.4% ($\pm$0.7%) | 25.5% ($\pm$0.7%) | 32.0% ($\pm$1.0%) | 41.1% ($\pm$0.9%) | 51.4% ($\pm$0.7%) | 65.3% ($\pm$2.2%) | 79.7% ($\pm$1.5%) | 90.8% ($\pm$1.6%) |
| WEIGHTED $L^1$ LOSS | 26.7% ($\pm$0.9%) | 34.0% ($\pm$0.9%) | 45.3% ($\pm$0.9%) | 55.5% ($\pm$0.7%) | 67.0% ($\pm$0.9%) | 74.5% ($\pm$1.3%) | 85.5% ($\pm$1.1%) | 93.6% ($\pm$0.7%) |
| BATCHRANK (K = 64) | **43.4%** ($\pm$0.5%) | 50.2% ($\pm$0.7%) | 58.7% ($\pm$0.7%) | 68.1% ($\pm$0.6%) | 77.1% ($\pm$0.4%) | 82.5% ($\pm$0.3%) | 90.3% ($\pm$0.3%) | 97.2% ($\pm$0.4%) |
| BATCHRANK (K = 8) | 39.6% ($\pm$3.2%) | **48.3%** ($\pm$3.5%) | **59.3%** ($\pm$3.5%) | **69.4%** ($\pm$3.3%) | **80.0%** ($\pm$1.8%) | 85.6% ($\pm$1.0%) | 92.7% ($\pm$0.4%) | 98.2% ($\pm$0.2%) |
| IDEAL | 62.3% | 69.6% | 78.5% | 84.9% | 90.5% | 93.8% | 97.3% | 99.1% |

*Table 6.* Coverage (%) on CAIDA dataset, 30th minute. The largest entries for each column are in **bold**.

| METHOD | COVERAGE SIZE | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | 1% | 2% | 5% | 10% | 20% | 30% | 50% | 75% |
| UNWEIGHTED LOG LOSS | 33.5% ($\pm$0.4%) | 40.1% ($\pm$0.5%) | 48.4% ($\pm$0.4%) | 59.5% ($\pm$0.3%) | 71.6% ($\pm$0.3%) | 80.4% ($\pm$0.2%) | 89.6% ($\pm$0.1%) | 97.7% ($\pm$0.1%) |
| WEIGHTED LOG LOSS | 31.1% ($\pm$0.6%) | 41.4% ($\pm$0.6%) | **54.3%** ($\pm$0.5%) | **64.5%** ($\pm$0.6%) | **77.7%** ($\pm$0.5%) | **84.2%** ($\pm$0.4%) | **92.1%** ($\pm$0.2%) | **98.1%** ($\pm$0.1%) |
| UNWEIGHTED $L^1$ LOSS | 21.7% ($\pm$0.9%) | 24.7% ($\pm$0.8%) | 31.7% ($\pm$1.0%) | 40.8% ($\pm$0.7%) | 50.7% ($\pm$0.5%) | 63.8% ($\pm$2.9%) | 78.7% ($\pm$1.6%) | 89.6% ($\pm$1.4%) |
| WEIGHTED $L^1$ LOSS | 18.6% ($\pm$0.8%) | 25.9% ($\pm$0.7%) | 38.5% ($\pm$1.1%) | 49.9% ($\pm$1.1%) | 63.5% ($\pm$1.2%) | 72.2% ($\pm$1.4%) | 84.3% ($\pm$1.1%) | 93.3% ($\pm$0.6%) |
| BATCHRANK (K = 64) | **36.8%** ($\pm$0.4%) | 43.7% ($\pm$0.6%) | 53.0% ($\pm$0.6%) | 63.8% ($\pm$0.7%) | 73.7% ($\pm$0.6%) | 79.8% ($\pm$0.4%) | 88.7% ($\pm$0.4%) | 96.7% ($\pm$0.4%) |
| BATCHRANK (K = 8) | 34.7% ($\pm$1.6%) | **42.6%** ($\pm$2.4%) | **53.9%** ($\pm$2.9%) | **64.5%** ($\pm$2.9%) | **76.1%** ($\pm$1.8%) | 82.4% ($\pm$1.2%) | 90.8% ($\pm$0.4%) | 97.4% ($\pm$0.4%) |
| IDEAL | 62.3% | 69.6% | 78.5% | 84.9% | 90.5% | 93.8% | 97.3% | 99.1% |

*Table 7.* Coverage (%) on CAIDA dataset, 60th minute. The largest entries for each column are in **bold**

| METHOD | COVERAGE SIZE | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | 1% | 2% | 5% | 10% | 20% | 30% | 50% | 75% |
| UNWEIGHTED LOG LOSS | 30.7% ($\pm$0.3%) | 37.0% ($\pm$0.4%) | 45.2% ($\pm$0.2%) | 55.9% ($\pm$0.2%) | 67.0% ($\pm$0.2%) | 78.2% ($\pm$0.3%) | 88.2% ($\pm$0.4%) | 96.6% ($\pm$0.2%) |
| WEIGHTED LOG LOSS | 28.3% ($\pm$0.8%) | 37.9% ($\pm$0.2%) | **51.5%** ($\pm$0.8%) | **62.2%** ($\pm$0.5%) | **75.1%** ($\pm$0.2%) | **82.2%** ($\pm$0.1%) | **91.0%** ($\pm$0.2%) | **97.6%** ($\pm$0.1%) |
| UNWEIGHTED $L^1$ LOSS | 19.7% ($\pm$0.9%) | 22.4% ($\pm$0.7%) | 29.1% ($\pm$1.2%) | 37.0% ($\pm$1.1%) | 46.3% ($\pm$0.6%) | 59.5% ($\pm$2.0%) | 77.1% ($\pm$1.5%) | 88.4% ($\pm$0.7%) |
| WEIGHTED $L^1$ LOSS | 16.4% ($\pm$1.0%) | 23.4% ($\pm$0.7%) | 35.9% ($\pm$1.0%) | 47.7% ($\pm$0.6%) | 61.7% ($\pm$0.8%) | 70.7% ($\pm$0.7%) | 83.1% ($\pm$0.9%) | 93.0% ($\pm$0.6%) |
| BATCHRANK (K = 64) | **33.9%** ($\pm$0.5%) | 40.5% ($\pm$0.4%) | 49.9% ($\pm$0.7%) | 61.4% ($\pm$0.5%) | 71.4% ($\pm$0.5%) | 78.0% ($\pm$0.2%) | 87.5% ($\pm$0.3%) | 95.9% ($\pm$0.3%) |
| BATCHRANK (K = 8) | 31.8% ($\pm$1.8%) | **39.8%** ($\pm$2.4%) | **51.3%** ($\pm$2.8%) | **62.1%** ($\pm$2.6%) | **73.7%** ($\pm$1.8%) | 80.6% ($\pm$1.3%) | 89.6% ($\pm$0.5%) | 96.8% ($\pm$0.2%) |
| IDEAL | 62.3% | 69.6% | 78.5% | 84.9% | 90.5% | 93.8% | 97.3% | 99.1% |