

# Putting the “Learning” in Learning-Augmented Algorithms for Frequency Estimation

Elbert Du\*, **Franklyn Wang\***, and Michael Mitzenmacher

Harvard University

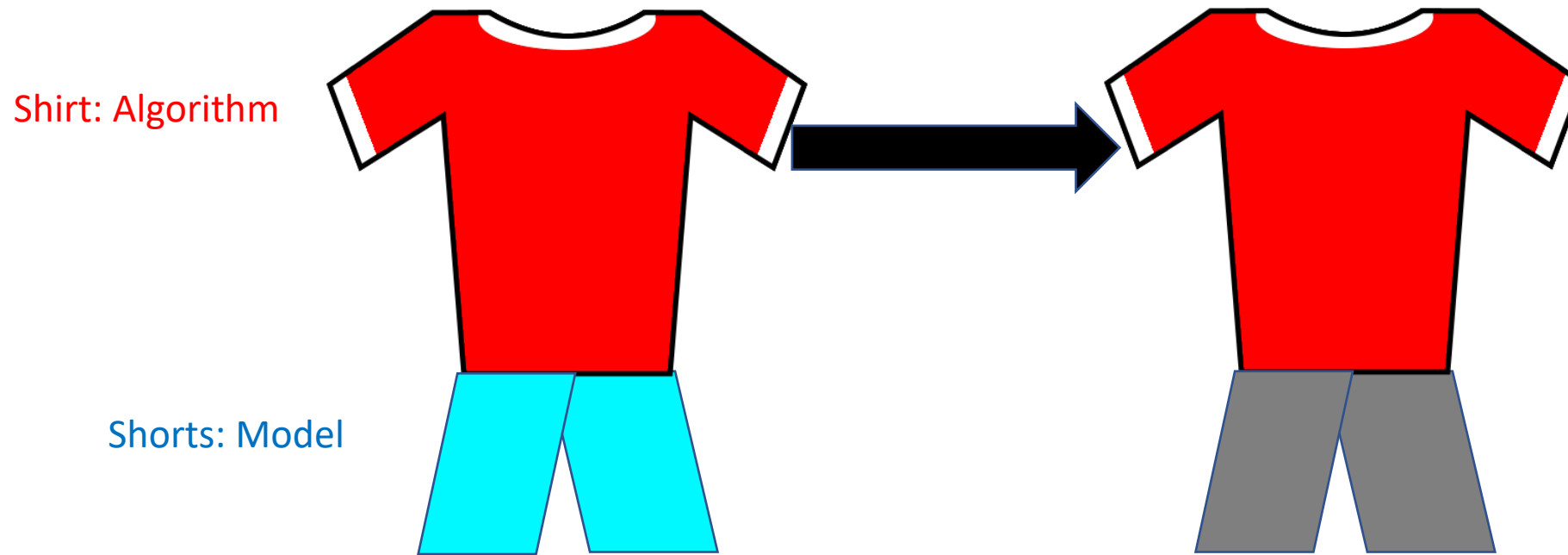
# Learning-Augmented Algorithms

- Learning-Augmented Algorithms are fusions of machine learning models and classical algorithms.
- In the present paper, we address the Learned Count-Min sketch [Hsu et. al 2019]

# Contribution, in a nutshell

One can think of learning-augmented algorithms as an outfit.

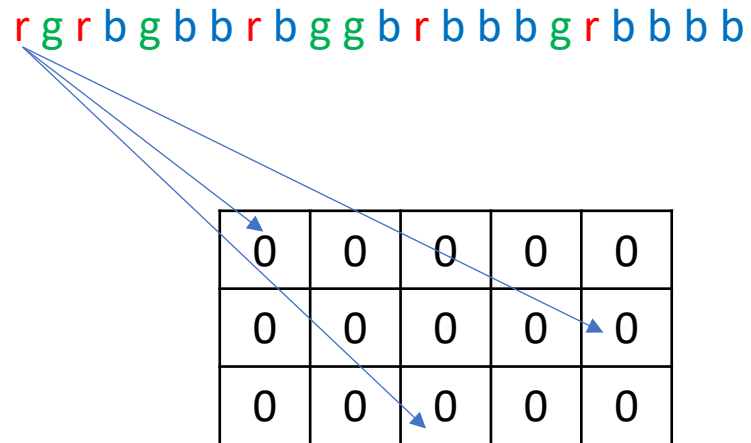
**Given the shirt, we choose the shorts.**



# Learned Count—Min Sketch

The Count-Min sketch is an algorithm which estimates frequencies in data streams with small space. It first hashes each element several times and increments a table of counters.

r g r b g b b r b g g b r b b b g r b b b b

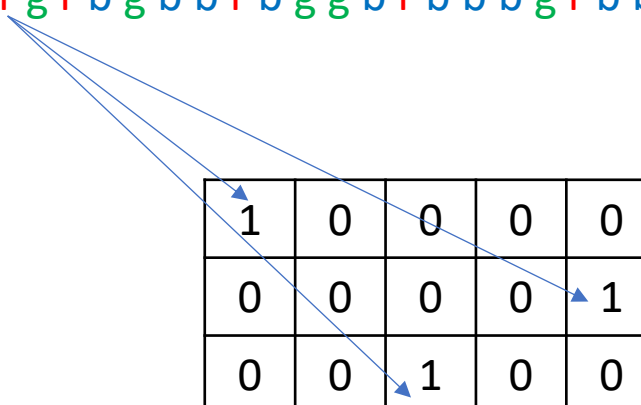


0	0	0	0	0
0	0	0	0	0
0	0	0	0	0

# Learned Count—Min Sketch

The Count-Min sketch is an algorithm which estimates frequencies in data streams with small space. It first hashes each element several times and increments a table of counters.

r g r b g b b r b g g b r b b b g r b b b b

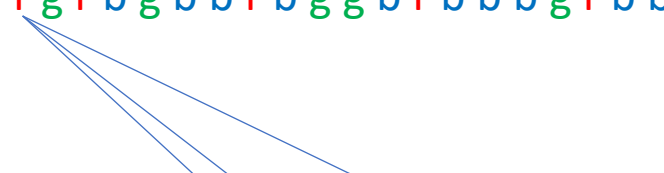


1	0	0	0	0
0	0	0	0	1
0	0	1	0	0

# Learned Count—Min Sketch

The Count-Min sketch is an algorithm which estimates frequencies in data streams with small space. It first hashes each element several times and increments a table of counters.

r g r b g b b r b g g b r b b b g r b b b b



5	0	0	0	0
0	0	0	0	5
0	0	5	0	0

# Learned Count—Min Sketch

The Count-Min sketch is an algorithm which estimates frequencies in data streams with small space. It first hashes each element several times and increments a table of counters.

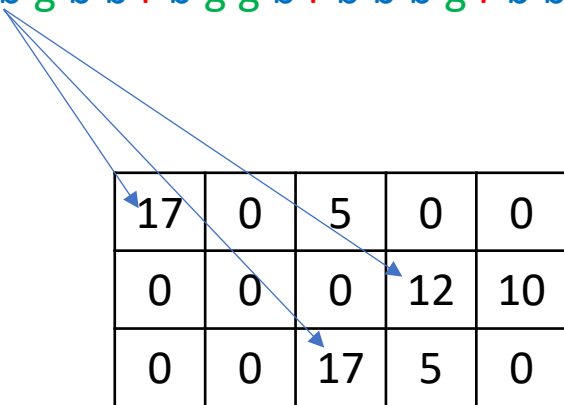
r g r b g b b r b g g b r b b b g r b b b b

5	0	5	0	0
0	0	0	0	10
0	0	5	5	0

# Learned Count—Min Sketch

The Count-Min sketch is an algorithm which estimates frequencies in data streams with small space. It first hashes each element several times and increments a table of counters.

r g r b g b b r b g g b r b b b g r b b b b



17	0	5	0	0
0	0	0	12	10
0	0	17	5	0



# Learned Count—Min Sketch

Finally, we estimate the frequency of an element by taking the *minimum* of all counters it corresponds to.

r g r b g b b r b g g b r b b b g r b b b b      r: 10

<b>17</b>	0	5	0	0
0	0	0	12	<b>10</b>
0	0	<b>17</b>	5	0

# Learned Count—Min Sketch

Finally, we estimate the frequency of an element by taking the *minimum* of all counters it corresponds to.

r g r b g b b r b g g b r b b b g r b b b b

r: 10

g: 5

17	0	5	0	0
0	0	0	12	10
0	0	17	5	0

# Learned Count—Min Sketch

Finally, we estimate the frequency of an element by taking the *minimum* of all counters it corresponds to.

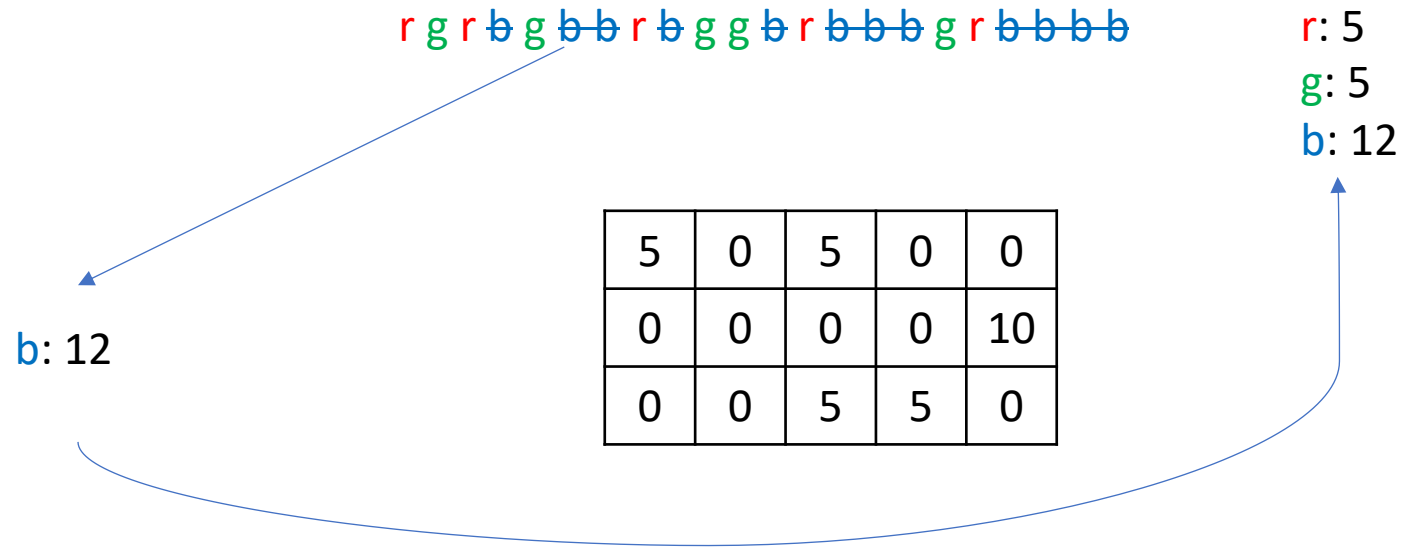
r g r b g b b r b g g b r b b b g r b b b b

r: 10  
g: 5  
b: 12

<b>17</b>	0	5	0	0
0	0	0	<b>12</b>	10
0	0	<b>17</b>	5	0

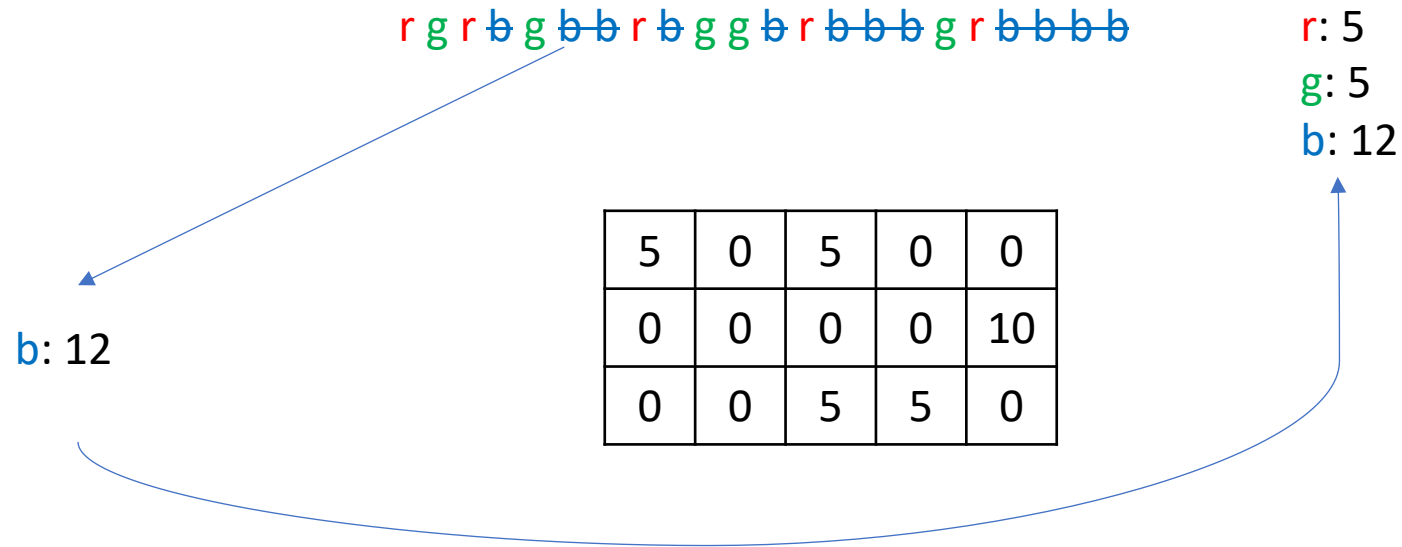
# Learned Count—Min Sketch

The learned count-min sketch is the same, except we *screen* some of the keys which are predicted to be heavy and track their frequencies individually.



# Optimizing the Sketch

Note that intuitively, screening more frequent keys is more valuable. We call the sum of all screened keys the *coverage*, which we seek to maximize.



# How do we optimize the model?

- We seek to optimize the model for downstream performance, **without changing architecture!**
- We will maximize the coverage.
- **Theorem (Informal):** The estimation error is essentially increasing in the complementary coverage.

# Optimizing for Coverage

- Hsu et. al uses the squared loss function:

$$\mathbb{E}_{i \sim F^0} \left[ \underbrace{(g_\theta(i))}_{\text{predictor}} - \underbrace{\ln f_i}_{\text{log frequency}} \right]^2$$

- Since we are optimizing for coverage, frequent elements are more important than infrequent elements.

$$\mathbb{E}_{i \sim F^1} \left[ (g_\theta(i) - \ln f_i)^2 \right]$$

# Optimizing for Coverage (cont.)

- Note that the coverage is proportional to

$$\mathbb{E}_{i \sim F^0} [1_{i \text{ screened}} f_i]$$

- Another method: **BatchRank**. We split each batch into sub-batches of size- $K$  and normalize, so we have

$$g'_\theta(i) = \frac{g_\theta(i) - \mu}{\sigma}.$$

- Then, our final loss function is:

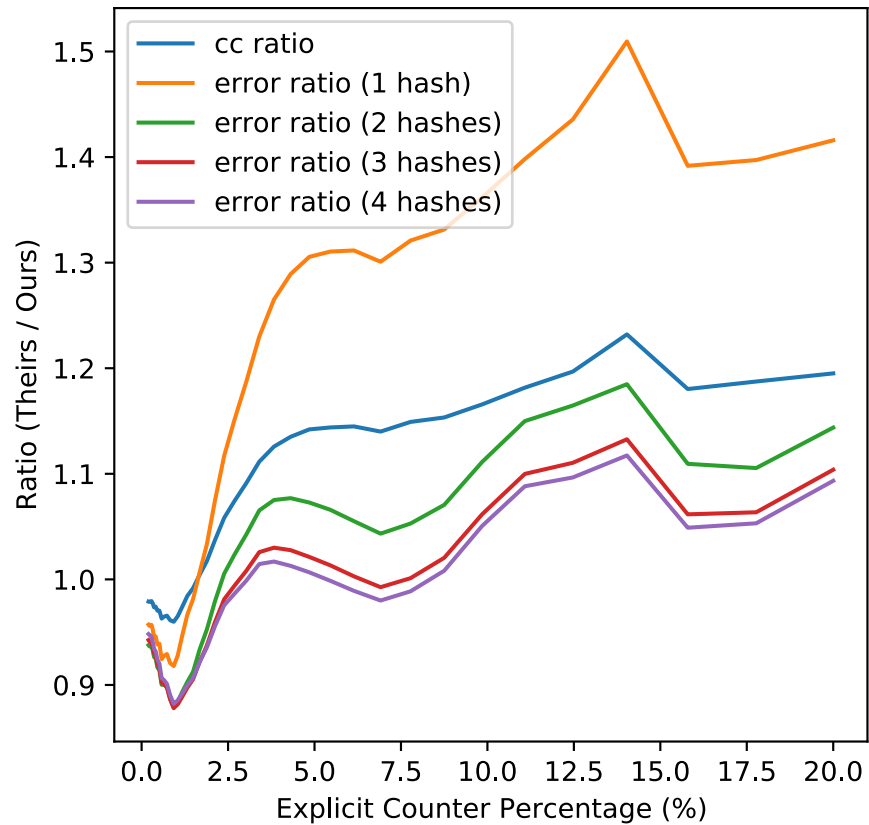
$$\mathbb{E}_{i \sim F^0} [g'_\theta(i) f_i]$$



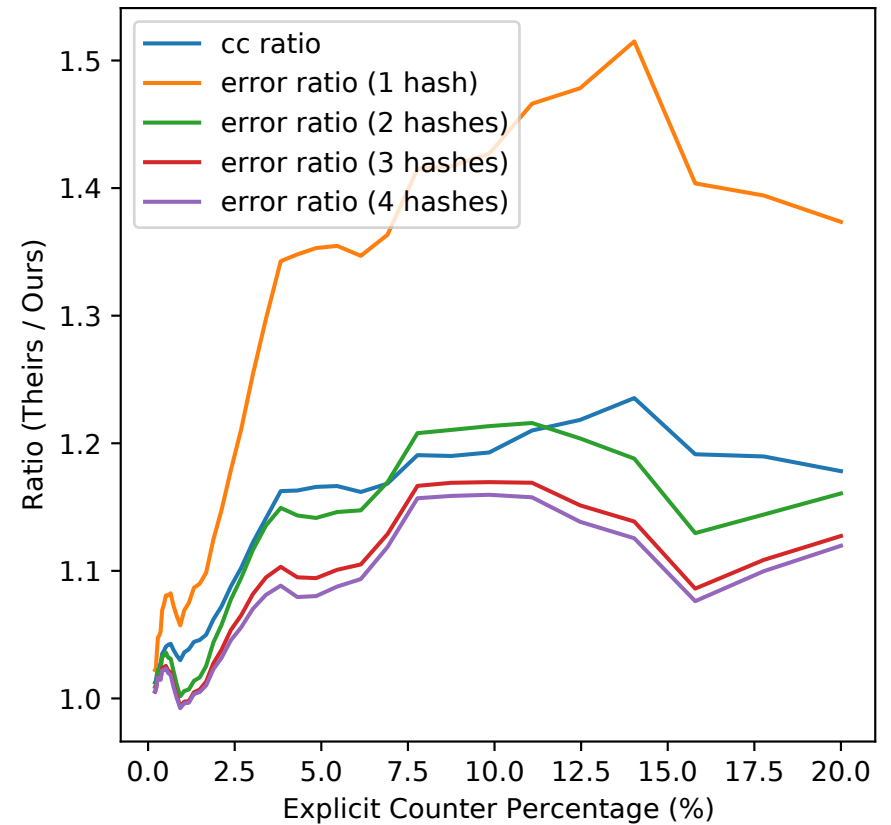
# Results

Ratios are theirs / ours

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



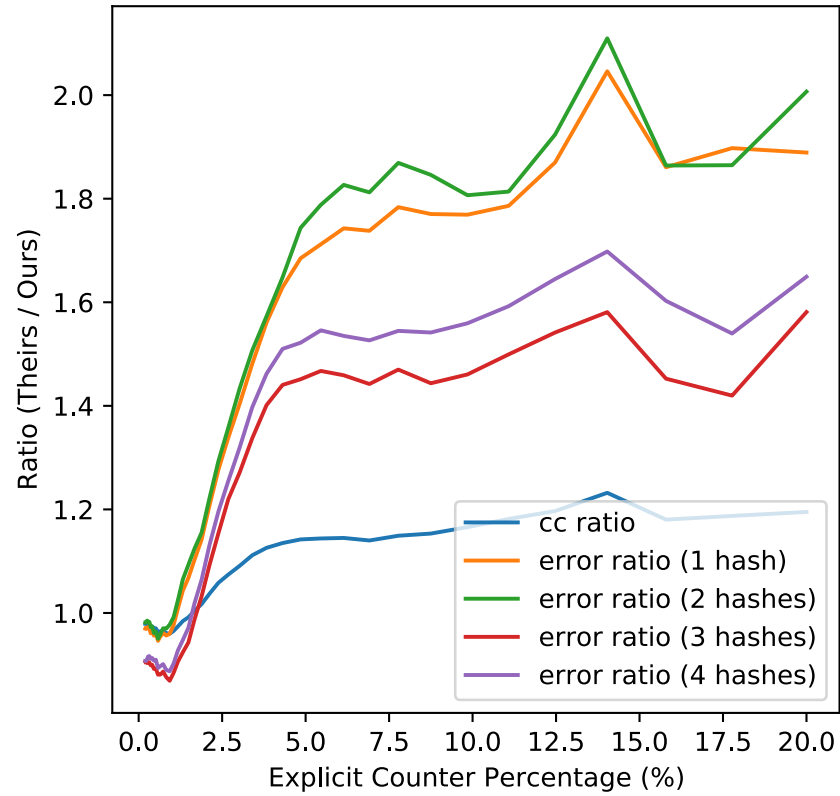
Minute 60, BatchRank (K = 8), width 10000, Count-Min



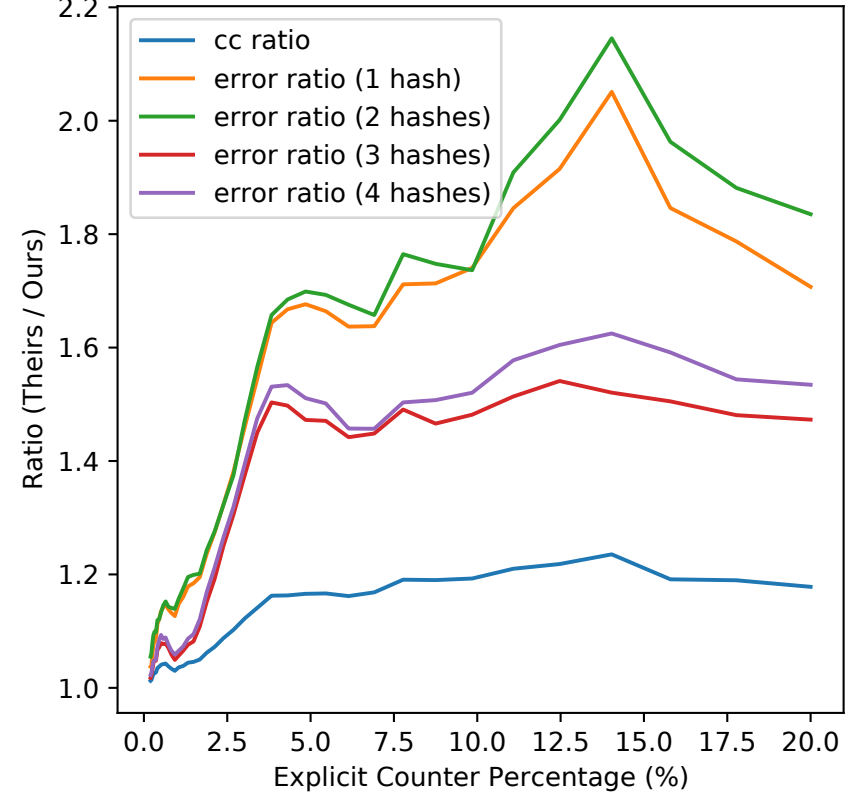
# Results

Ratios are theirs / ours

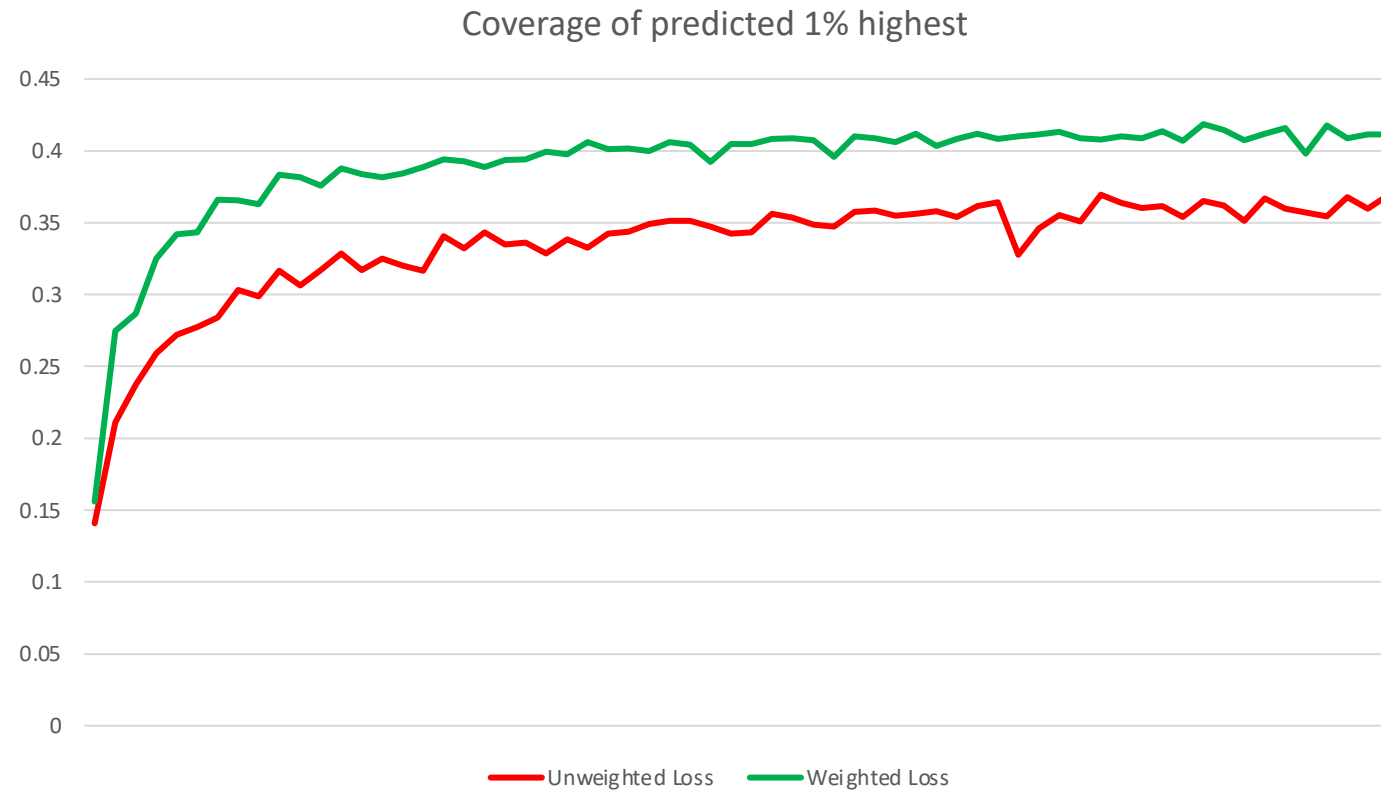
Minute 60, Weighted Log Loss, width 10000, Count-Sketch



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

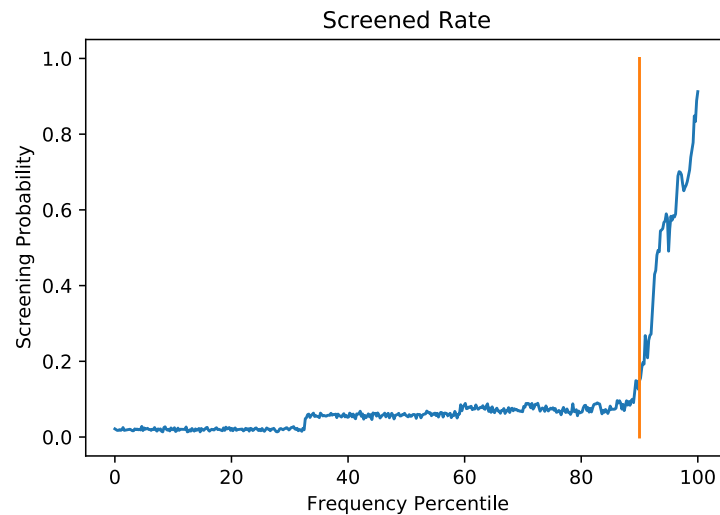


# Results



# A new model for classification errors

- We observe the rates by which keys are screened and note that it increases monotonically with frequency.



- This suggests a heterogeneous error model, allowing for stronger guarantees than a uniform error model.

# Conclusion

- We would like to invite the machine learning community to try out these tasks—our code can be found [here](#).
- When one can find an accurate proxy for performance (e.g. coverage), training to optimize the proxy can lead to significant improvements.

# References

- Hsu, Chen-Yu, et al. "Learning-Based Frequency Estimation Algorithms." International Conference on Learning Representations. 2019.