

Mergeable Summaries and the DataSketches Library



Edo Liberty,
Principal Scientist
Head of Amazon AI Labs



Justin Thaler
Assistant Professor
Georgetown University.



Lee Rhodes
Distinguished Architect
Oath, Inc.



Alexander Saydakov
Senior Software Engineer
Oath, Inc



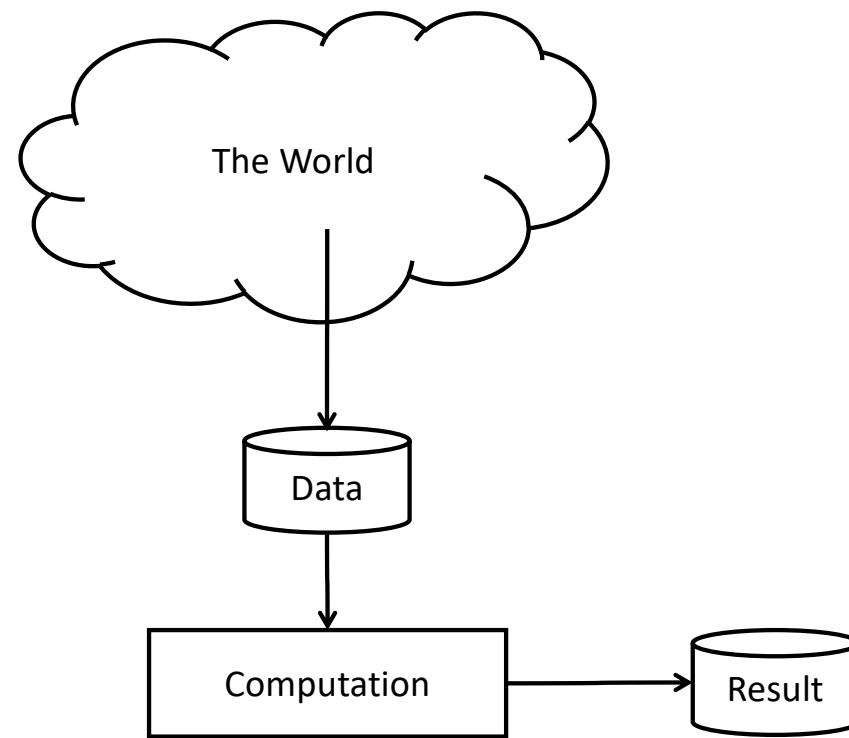
Kevin Lang
Principal Scientist
Oath, Inc



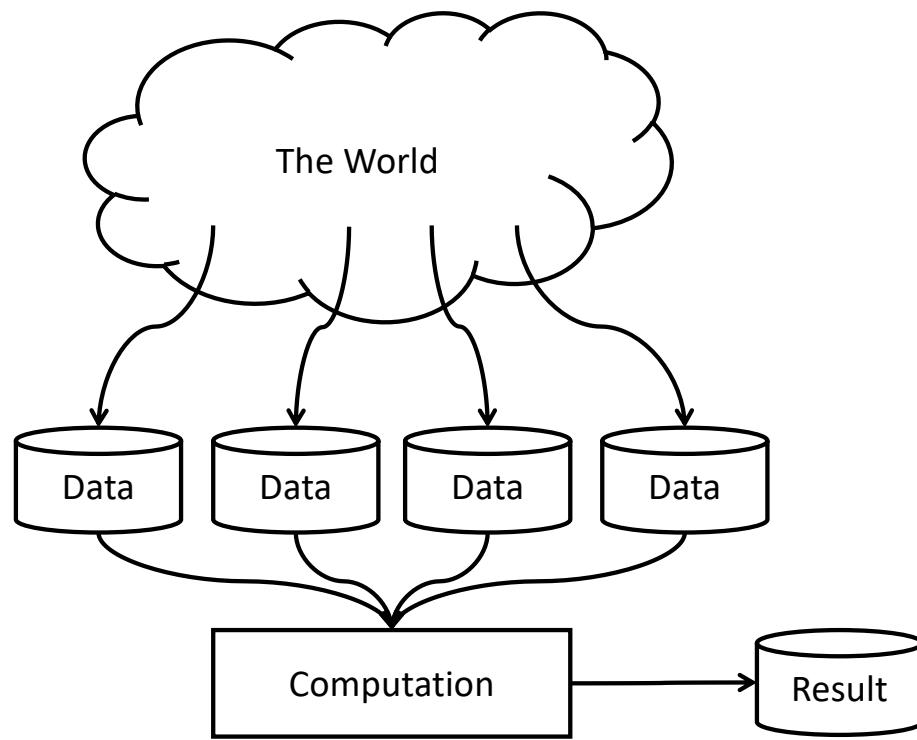
Jon Malkin
Senior Scientist
Oath, Inc



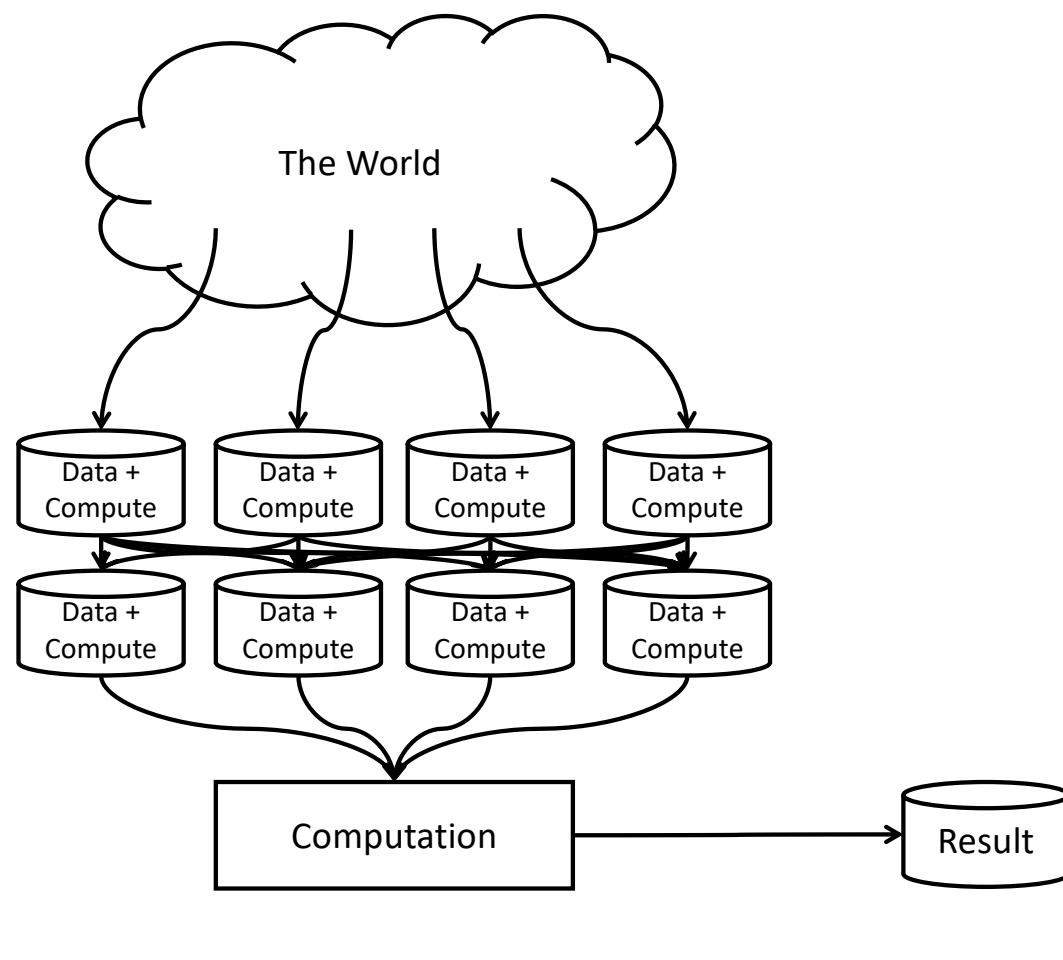
Single machine data processing



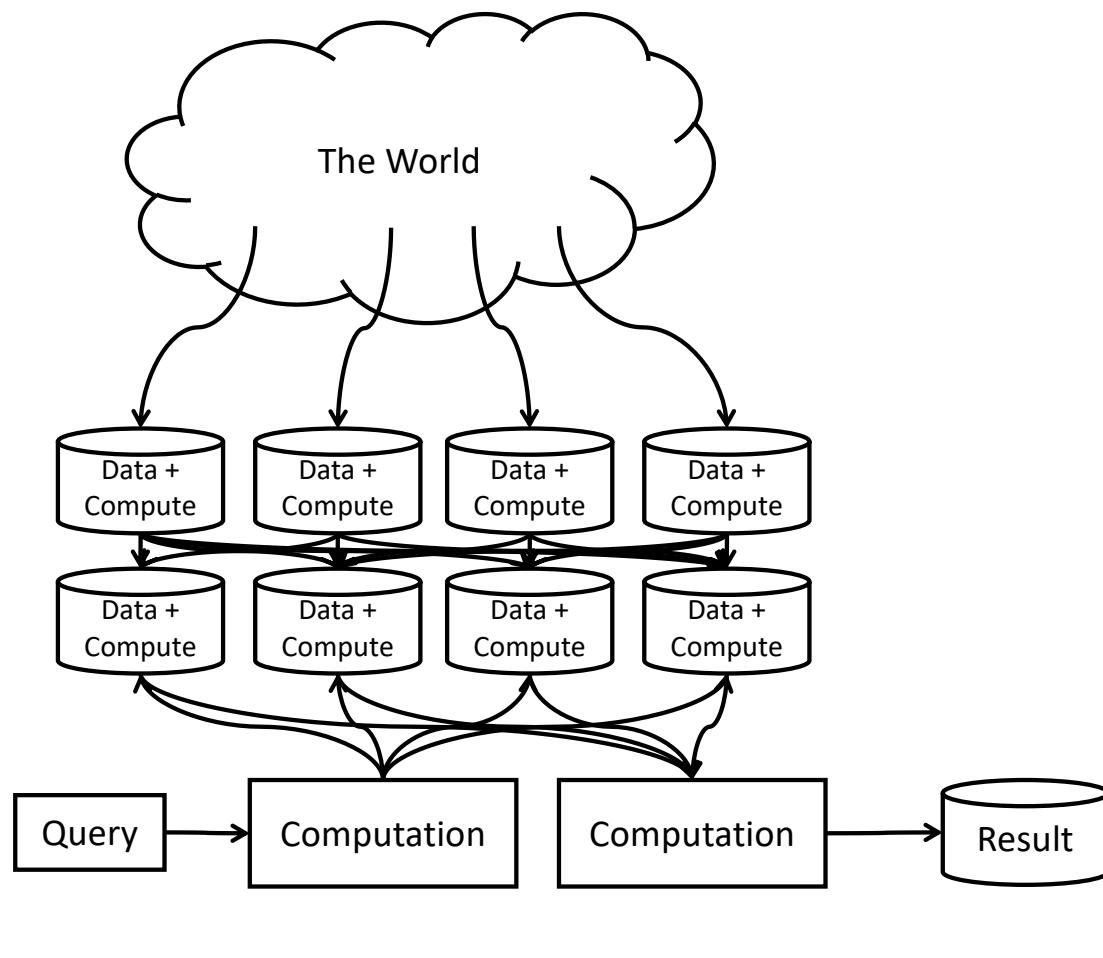
Distributed storage



Distributed compute (map/reduce, MPI, ...)



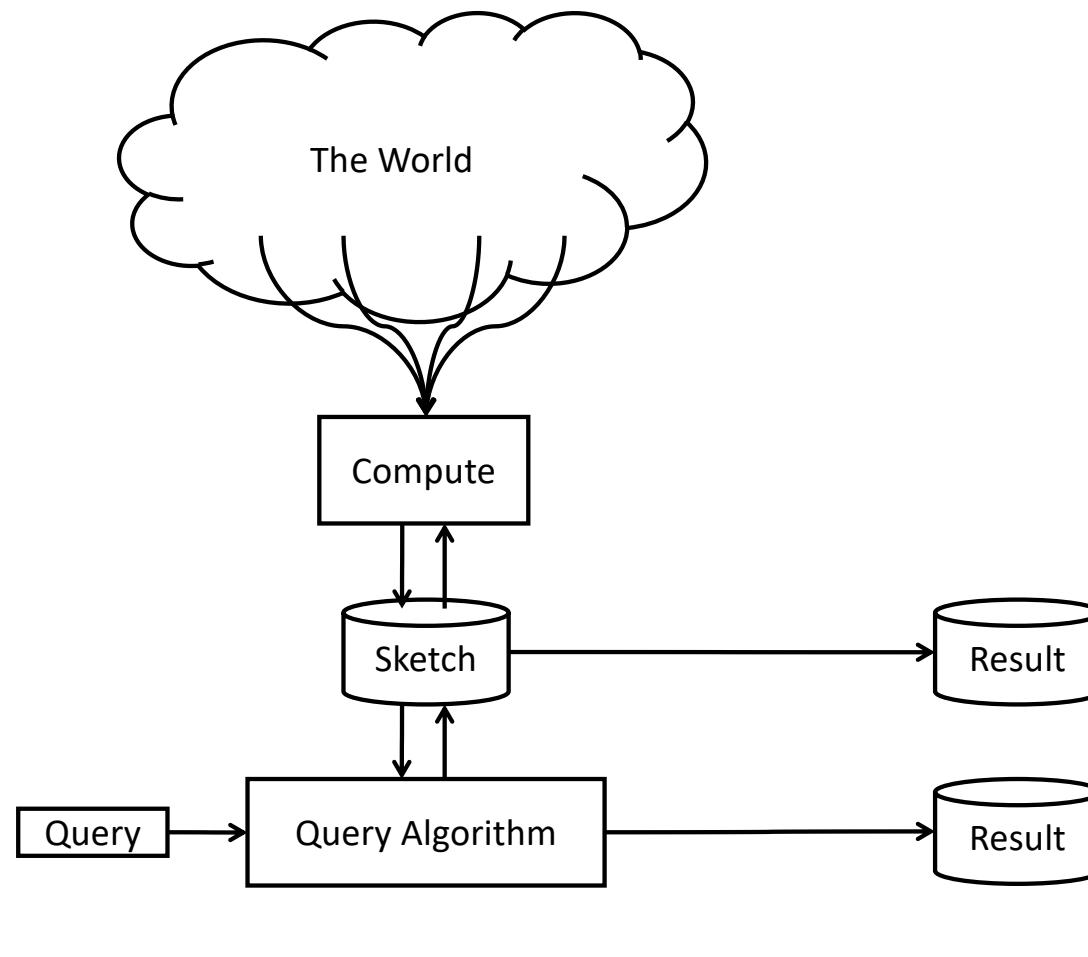
Distributed model (indexes, databases, Spark...)



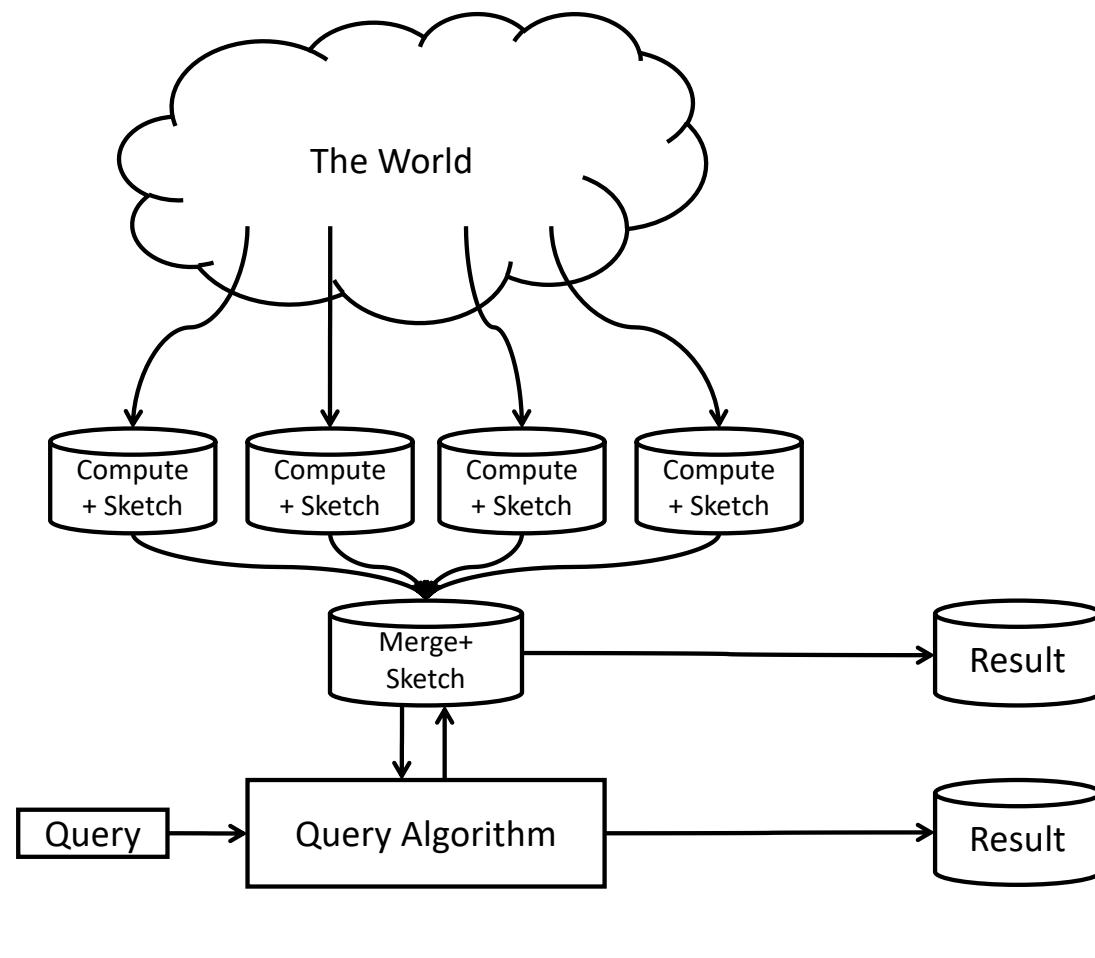
207 big-data infographics (a meta infographic)



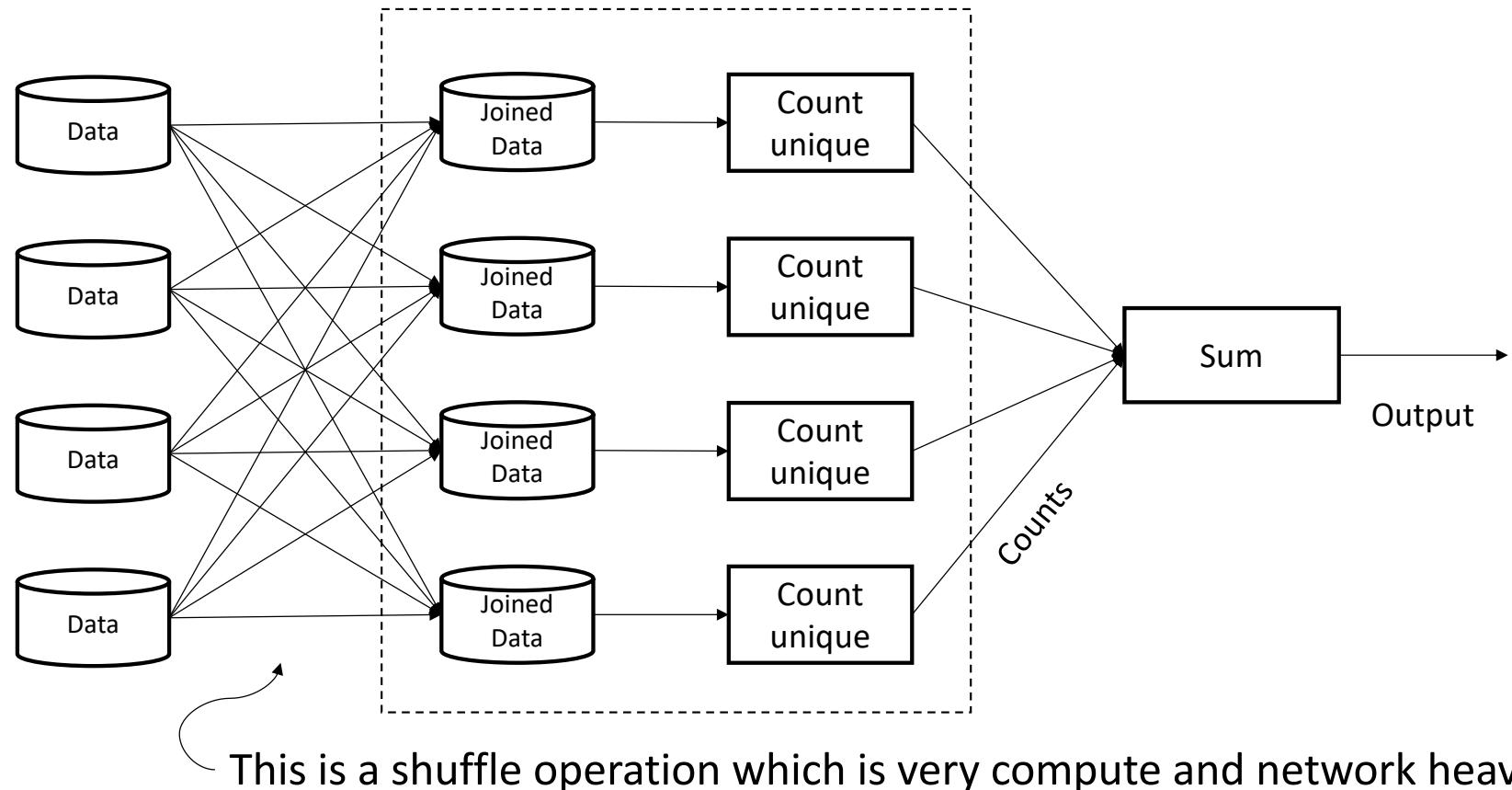
The streaming model



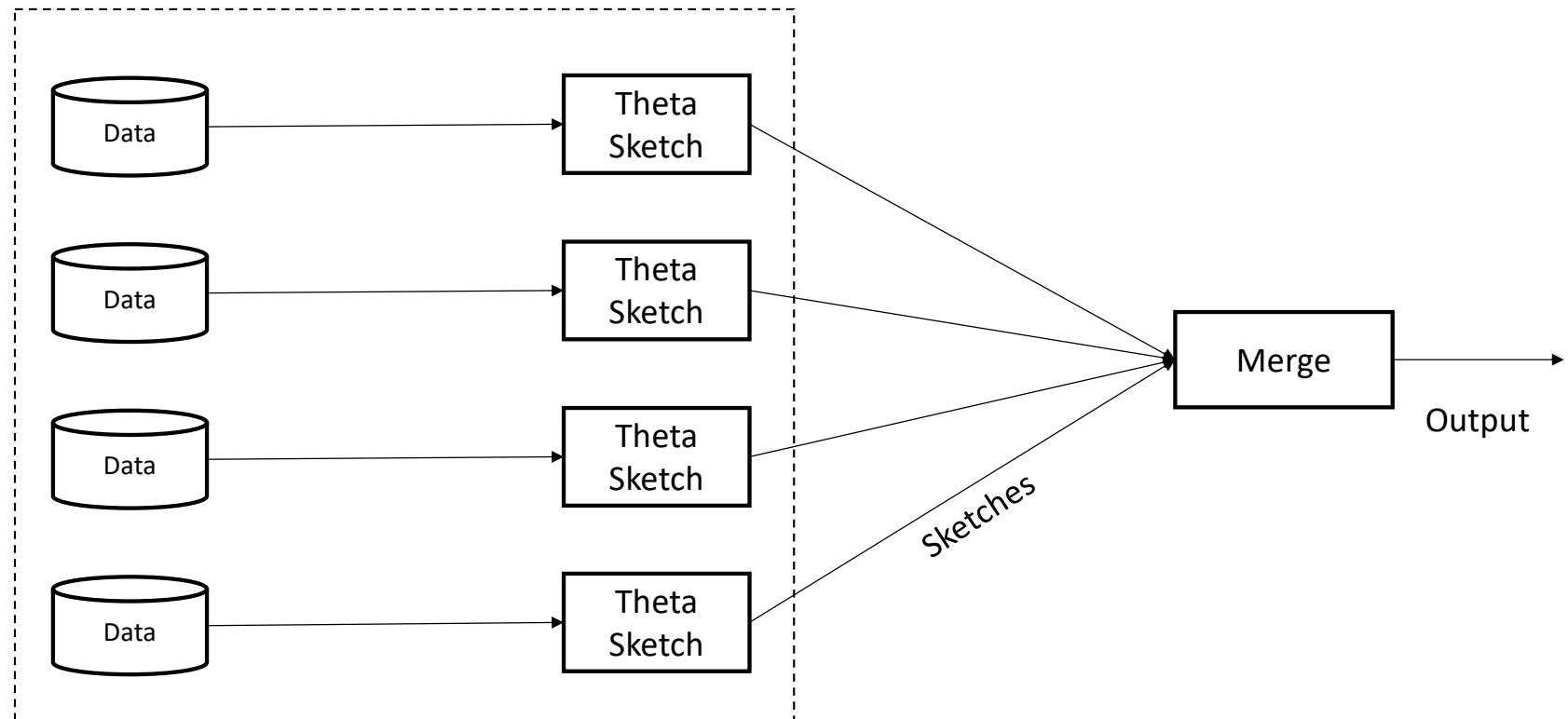
Mergeable Summaries



Unique Counting with Map Reduce

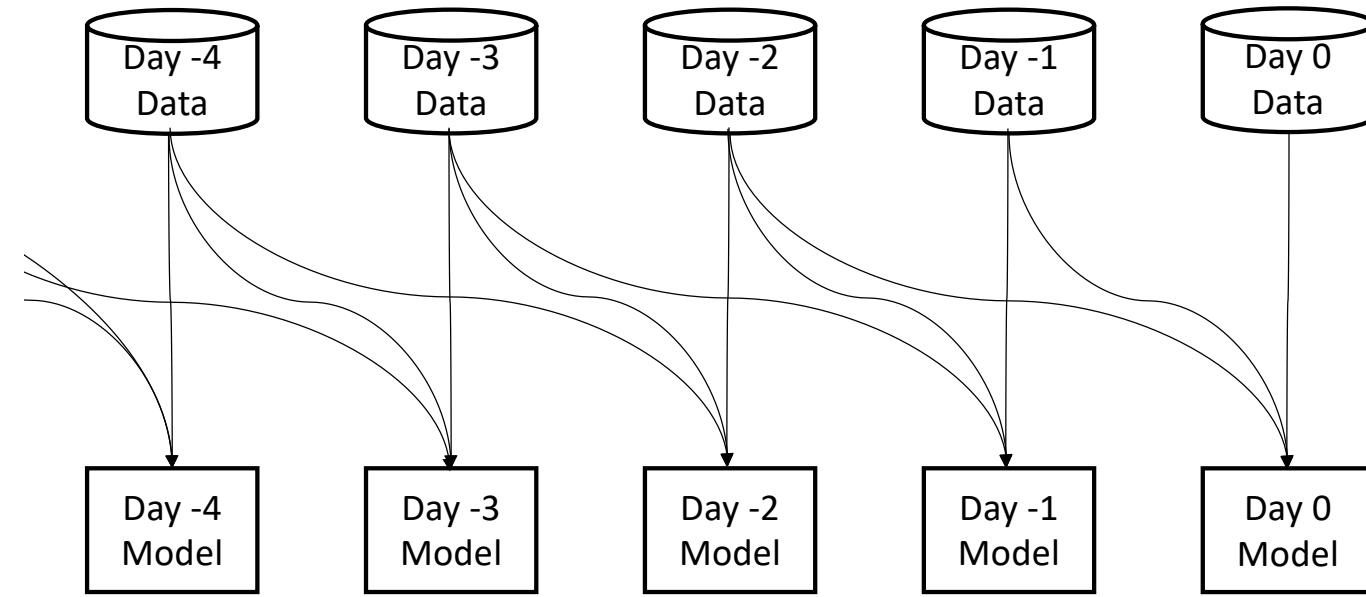


Unique Counting with Mergeable Summaries



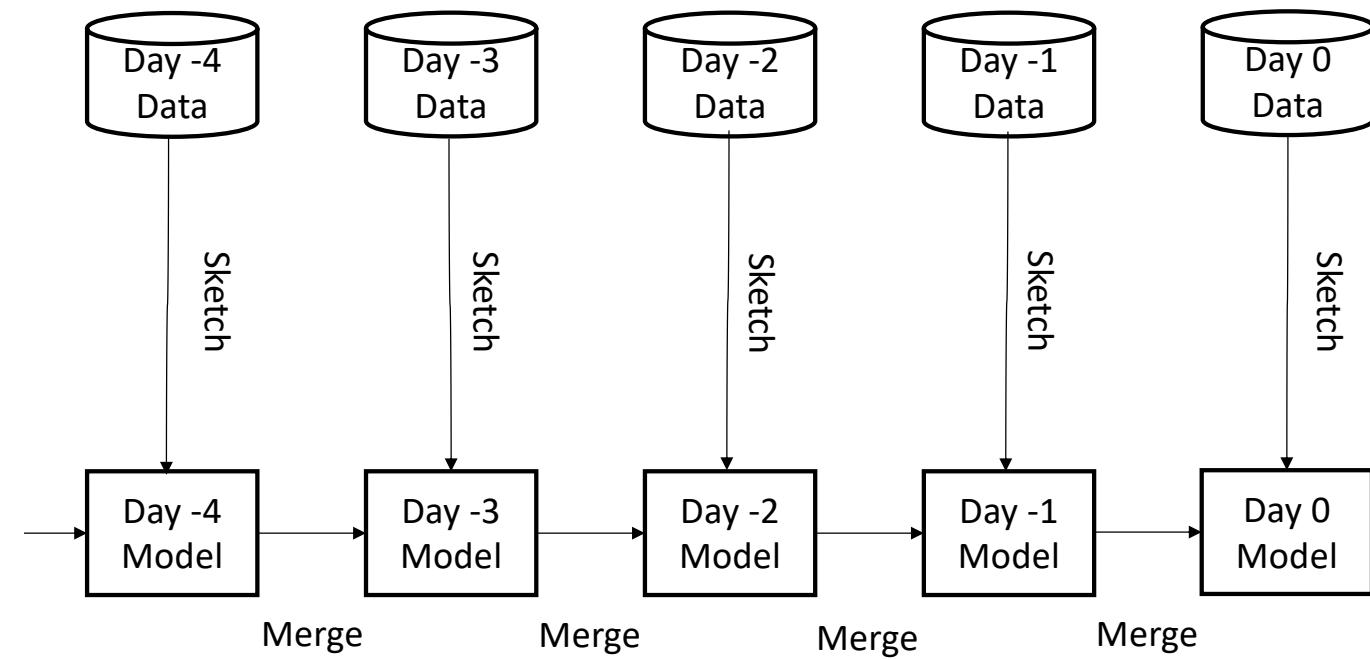
No shuffle operation is needed!

Data Mining with Traditional Windowing



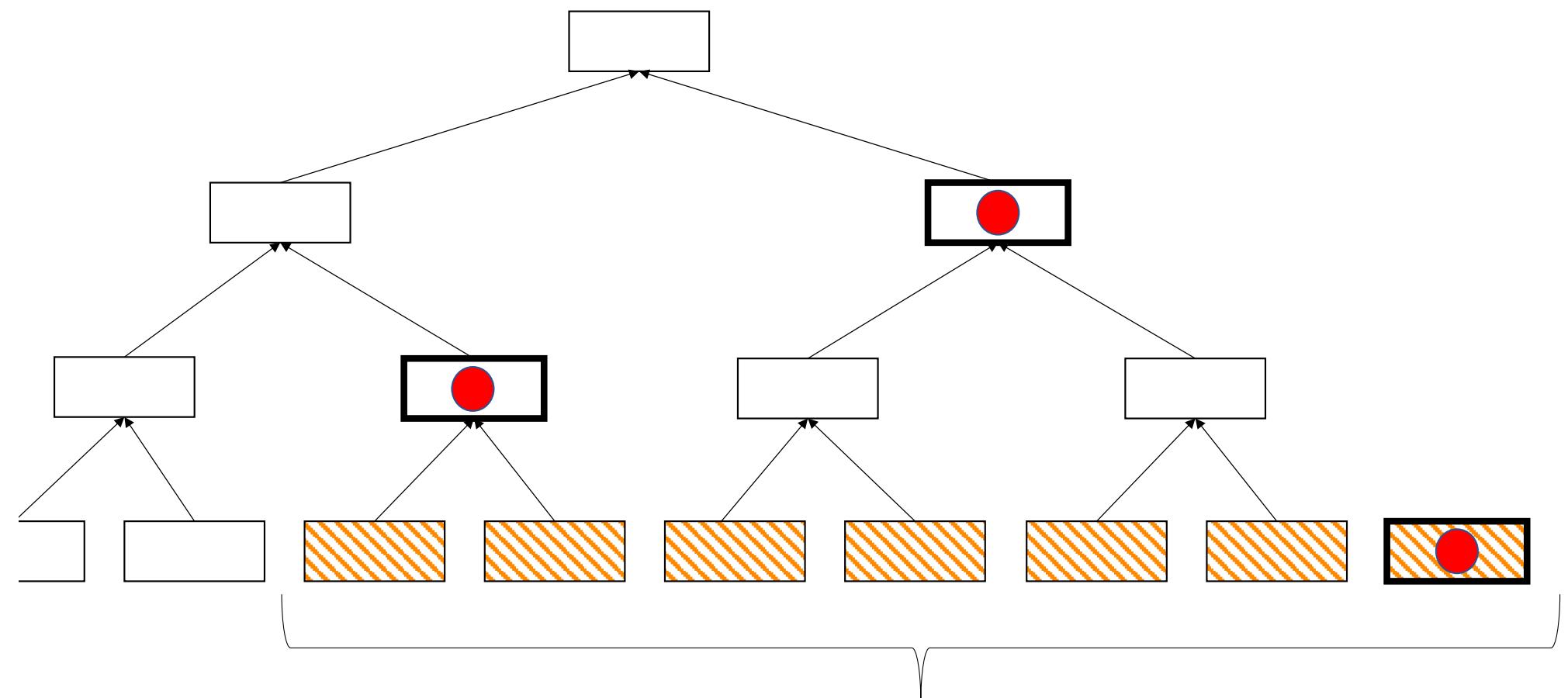
Every dataset is processed 3 times for a model consisting of 3 days

Data Mining with Mergeable Summaries



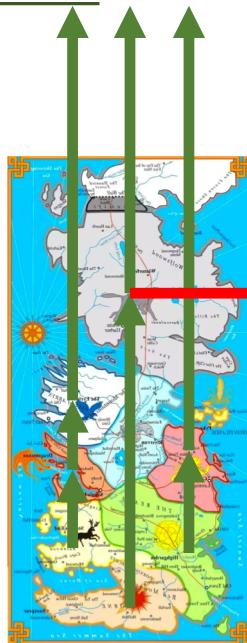
Every dataset is processed once for a model consisting of the entire history

Dynamic Windowing with Mergeable Summaries



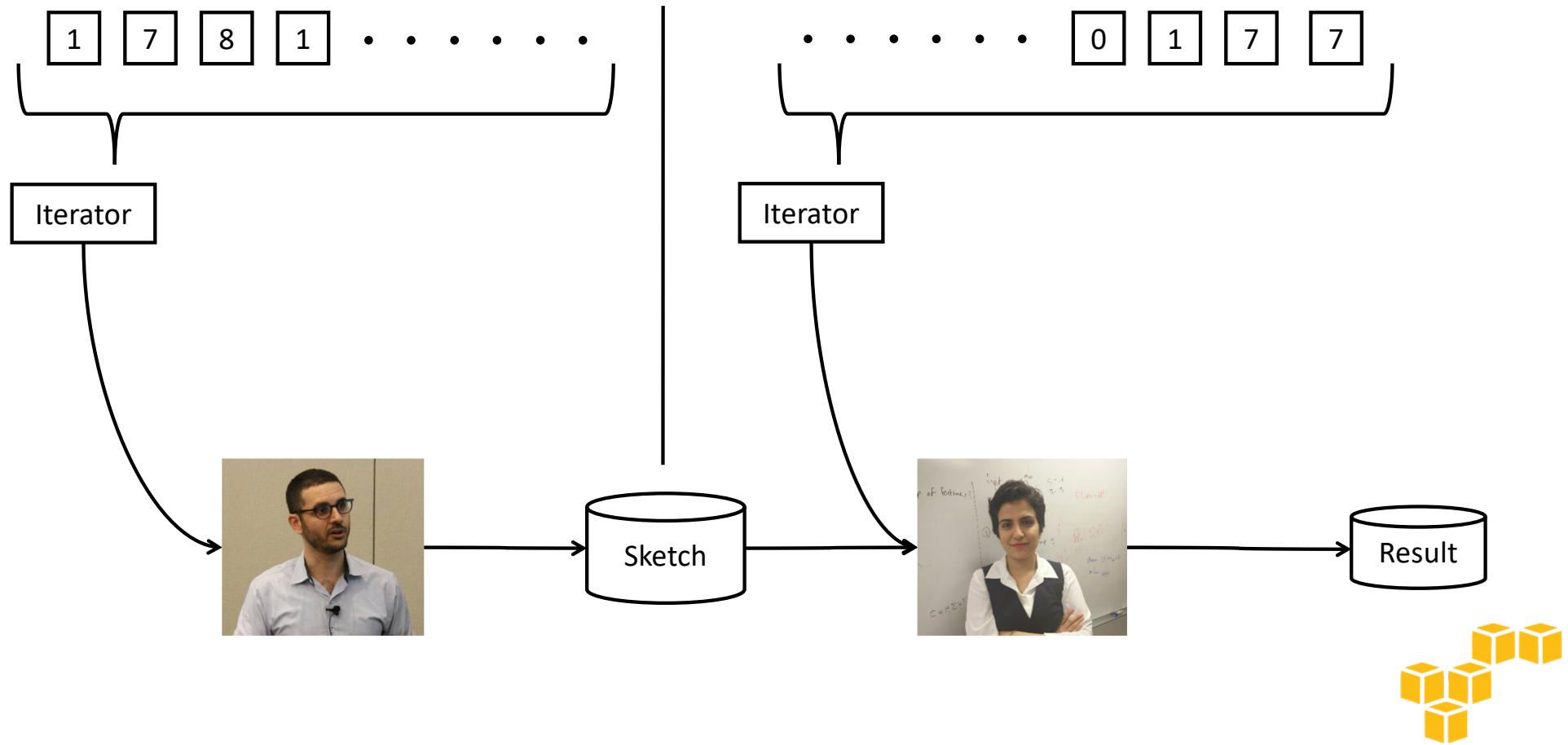
OLAP with Mergeable Summaries

“Median latency of
IoT device call in
Westeros January 2018”



“Median latency of
IoT device call in
The North Q1 2018”

Some Basic Problems are Impossible



Others Have Great Solutions

Items

(words, IP-addresses, events, clicks,...)

- Counting distinct elements
- Item frequencies
- Approximate Quantiles
- Moment and entropy estimation
- Approximate set operations
- Sampling

Matrices

(text corpora, recommendations, ...)

- Covariance estimation matrix
- Low rank approximation
- Sparsification

Vectors

(text documents, images, example features,...)

- Dimensionality reduction
- Clustering (k-means, k-median,...)
- Linear Regression
- Machine learning (some of it at least)
- Density Estimation / Anomaly detection

Graphs*

(social networks, communications, ...)

- Connectivity
- Cut Sparsification
- Weighted Matching



Data Sketches – Open Source Project



sketches-core

Core Sketch Library.

Java ★ 500 130 Apache-2.0 Updated an hour ago



```
>> brew tap DataSketches/sketches-cmd  
>> brew install data-sketches
```

L. Rhodes, K. Lang, A. Saydakov, J. Thaler, E. Liberty, and J. Malkin. DataSketches: A Java software library of stochastic streaming algorithms, 2017. <https://datasketches.github.io>.



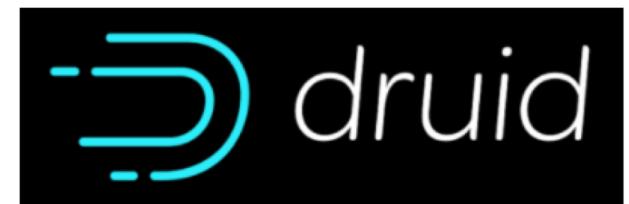
Production ready



YAHOO!



splunk>



New Research

- [ABL+17] Daniel Anderson, Pryce Bevan, Kevin J. Lang, Edo Liberty, Lee Rhodes, and Justin Thaler.
A high-performance algorithm for identifying frequent items in data streams.
In *ACM IMC 2017 (To Appear)*, 2017. Preliminary version available at <https://arxiv.org/abs/1705.07001>.
- [DLRT16] Anirban Dasgupta, Kevin J. Lang, Lee Rhodes, and Justin Thaler.
A framework for estimating stream expression cardinalities. In **ACM ICDT Proceedings '16 **, pages 6:1–6:17, 2016.
- [KLL16] Zohar S. Karnin, Kevin J. Lang, and Edo Liberty.
Optimal quantile approximation in streams. In *IEEE FOCS Proceedings '16*, pages 71–78, 2016.
- [Lan17] Kevin J Lang. Back to the future: an even more nearly optimal cardinality estimation algorithm.
arXiv preprint <https://arxiv.org/abs/1708.06839>, 2017.
- [Lib13] Edo Liberty. Simple and deterministic matrix sketching.
In *ACM KDD Proceedings '13*, pages 581–588, 2013.
- [LMTU16] Edo Liberty, Michael Mitzenmacher, Justin Thaler, and Jonathan Ullman.
Space lower bounds for itemset frequency sketches. In *ACM PODS Proceedings '16*, pages 441–454, 2016.
- [MST12] Michael Mitzenmacher, Thomas Steinke, and Justin Thaler.
Hierarchical heavy hitters with the space saving algorithm. In *SIAM ALENEX Proceedings '12*, pages 160–174, 2012.

In this talk

Items

(words, IP-addresses, events, clicks,...)

- Counting distinct elements
- Item frequencies
- Approximate Quantiles
- Moment and entropy estimation
- Approximate set operations
- Sampling

Matrices

(text corpora, recommendations, ...)

- Covariance estimation matrix
- Low rank approximation
- Sparsification

Vectors

(text documents, images, example features,...)

- Dimensionality reduction
- Clustering (k-means, k-median,...)
- Linear Regression
- Machine learning (some of it at least)
- Density Estimation / Anomaly detection

Graphs*

(social networks, communications, ...)

- Connectivity
- Cut Sparsification
- Weighted Matching



Counting Distinct Elements

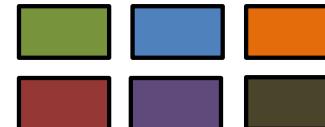
- N. Alon, Y. Matias, and M. Szegedy. The space complexity of approximating the frequency moments
Z. Bar-Yossef, T. Jayram, R. Kumar, D. Sivakumar, and L. Trevisan. Counting distinct elements in a data stream
E. Cohen. All-distances sketches, revisited: HIP estimators for massive graphs analysis
E. Cohen and H. Kaplan. Summarizing data using bottom-k sketches
G. Cormode. Sketch techniques for massive data
P. Flajolet and G. Nigel Martin. Probabilistic counting algorithms for data base applications
D. M. Kane, J. Nelson, and D. P. Woodruff. An optimal algorithm for the distinct elements problem
M. Thorup. Bottom-k and priority sampling, set similarity and subset sums with minimal independence
A. Dasgupta, K Lang, L Rhodes, J. Thaler, A Framework for Estimating Stream Expression Cardinalities



Problem Definition

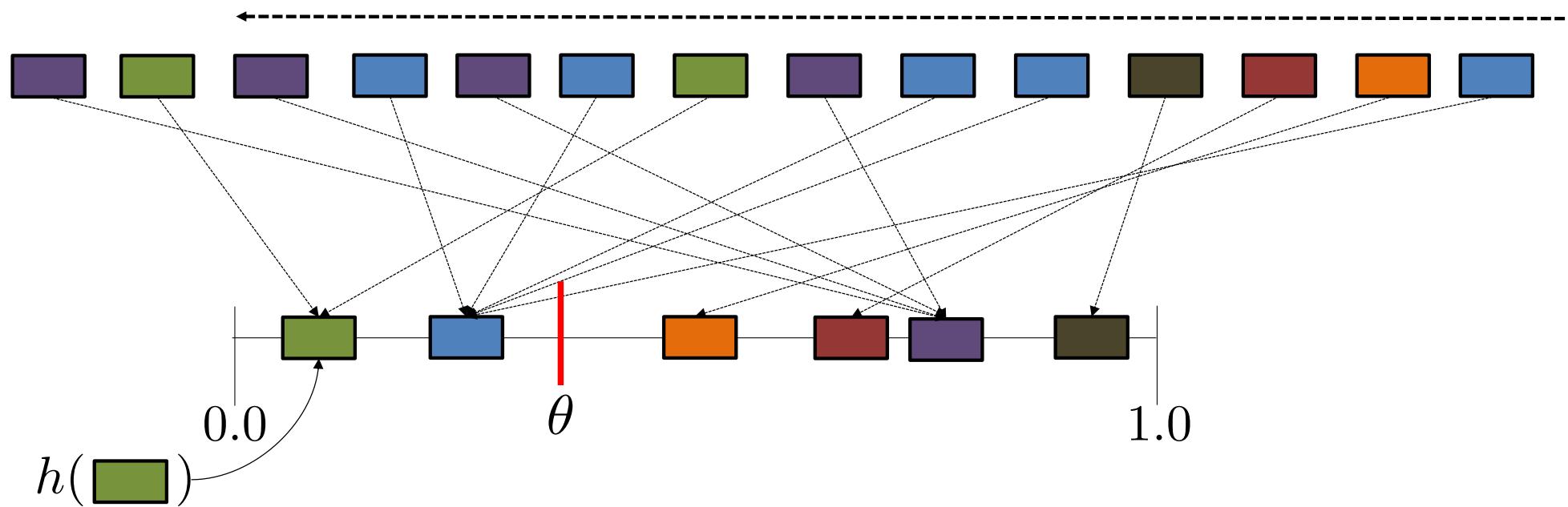


Approximate the number of
distinct items in the stream



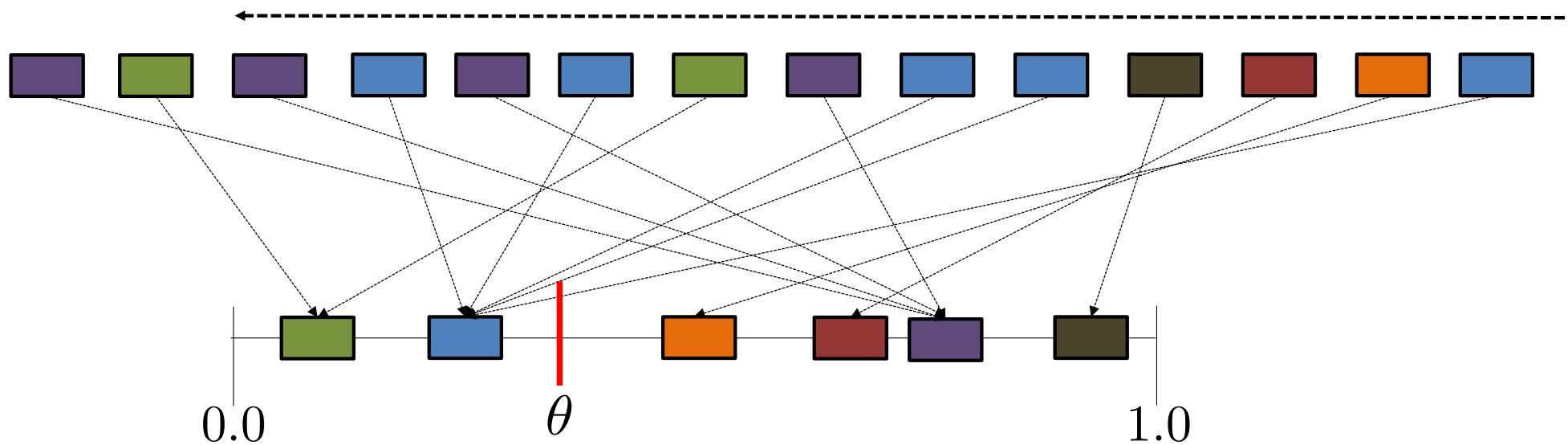
- # of unique IPs are important statistics of networks
- # of different customers using web services
- # of unique keys in a database join table
- ...

General Hashing Idea



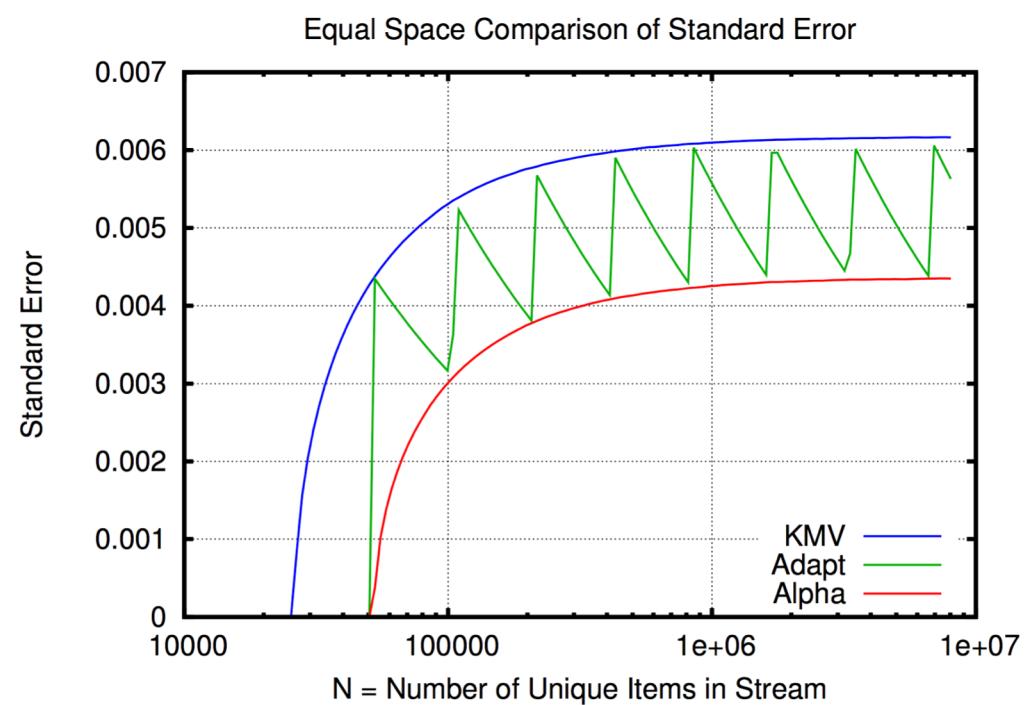
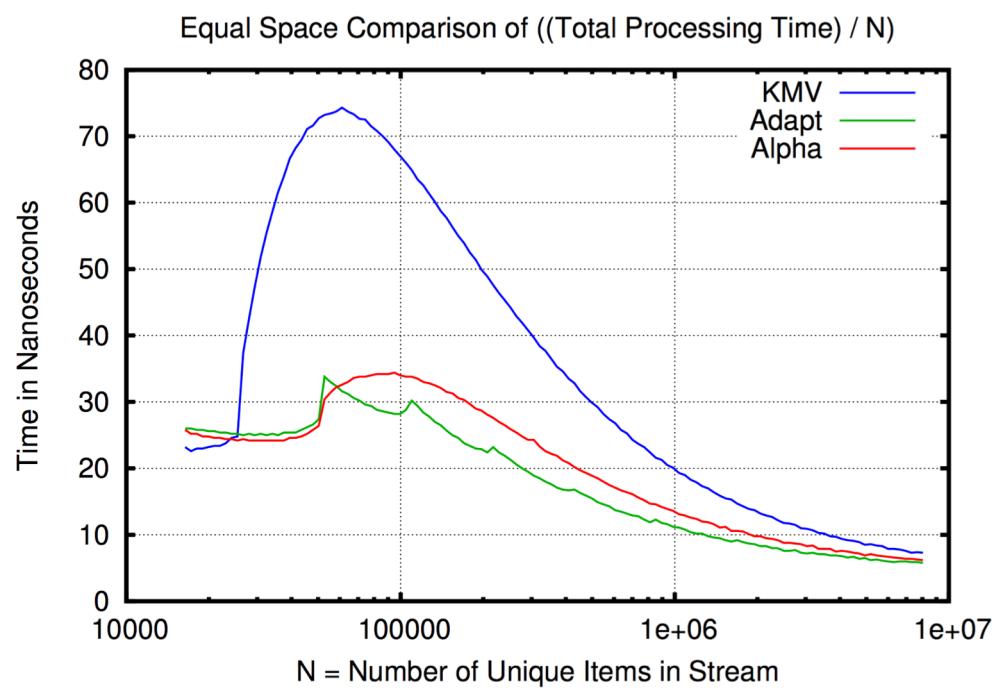
- Map all entries to the interval $(0,1)$ using a hash function
- Keep only k values smaller than some threshold θ .
- From (k, θ) we can approximate the number of unique items.

Our results



- Generalize a family of algorithms (Adaptive sampling, KMV)
- New Variance bounds for all such algorithms
- New tradeoffs between accuracy, space, and update time (alpha alg')
- Very careful implementation

Experimental Results



Quick Demo

Counting Distinct Elements



Assume you need to estimate the number of **unique** numbers in a file

```
>>head data.csv  
0  
1  
0  
3  
0  
2  
3  
7  
3  
2
```

In this file, row i tasks integer value from $[0,i]$ uniformly at random.



Some stats: there are 10,000,000 such numbers in this ~76Mb file.

```
>>time wc -lc data.csv  
10000000 76046666 data.csv  
  
real 0m0.101s  
user 0m0.072s  
sys 0m0.021s
```

Reading the file take ~1/10 seconds. We don't foresee IO being an issue.



To count the number of distinct items you might try this:

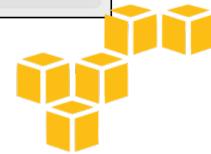
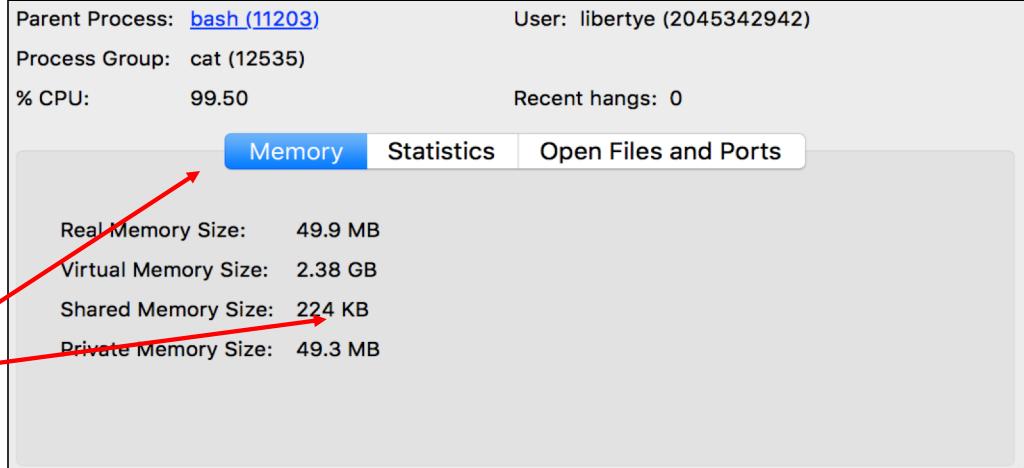
```
>>sort data.csv | uniq | wc -l
```

However, it is faster to “unify” while sorting.

```
>>sort data.csv -u | wc -l
```

```
>>time sort data.csv -u | wc  
-l  
5001233
```

```
real 2m37.071s  
user 2m36.587s ←  
sys 0m0.376s
```



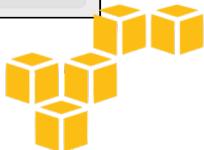
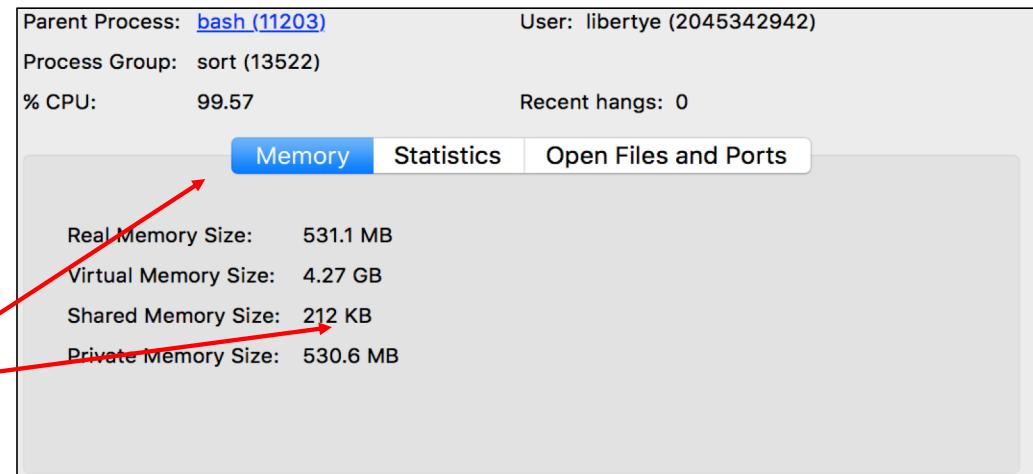
Still, most of the time is spent on comparing strings....

```
>>sort data.csv -u -n -S 100% | wc -l
```

This is much better!

```
>>time sort data.csv -u -n |  
wc -l  
5001233
```

```
real 0m11.809s  
user 0m11.587s ←  
sys 0m0.228s
```



This is the way to do this with the sketching library

```
>>cat data.csv | ds theta
```

```
>>time cat data.csv | ds  
theta  
4835874.828244  
4974249.044005  
5116568.411181  
real 0m1.539s  
user 0m2.184s ←  
sys 0m0.310s
```

Too fast for system monitor UI...

It uses ~ 32k of memory!



Weighted Item frequencies

Space-optimal heavy hitters with strong error bounds.] R. Berinde, P. Indyk, G. Cormode, and M. J. Strauss

An optimal algorithm for `1-heavy hitters in insertion streams and related problems A. Bhattacharyya, P. Dey, and D. P. Woodruff

Finding frequent items in data streams M. Charikar, K. Chen, and M. Farach-Colton

Methods for finding frequent items in data streams G. Cormode and M. Hadjieleftheriou

An improved data stream summary: The count-min sketch and its applications. G. Cormode and S. Muthukrishnan.

Approximate frequency counts over data streams G. S. Manku and R. Motwani.

Efficient computation of frequent and top-k elements in data streams A. Metwally, D. Agrawal, and A. El Abbadi.

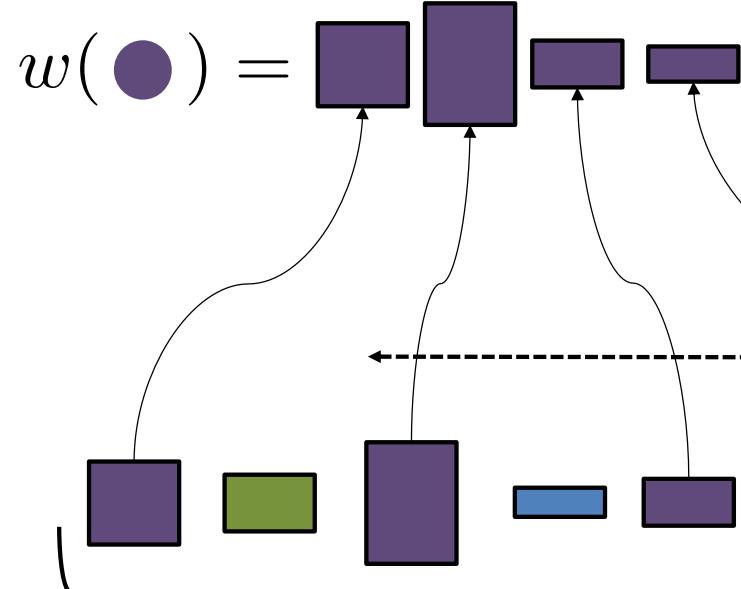
Finding repeated elements J. Misra and D. Gries.

A High-Performance Algorithm for Identifying Frequent Items in Data Streams

Daniel Anderson, Pryce Bevin, Kevin Lang, Edo Liberty, Lee Rhodes, Justin Thaler



Problem Definition



$$|w'(x) - w(x)| \leq \varepsilon W$$

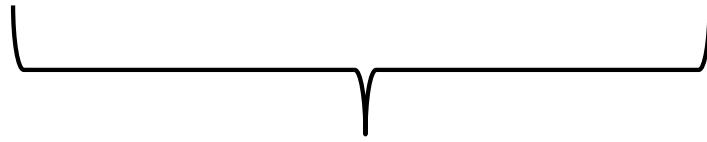
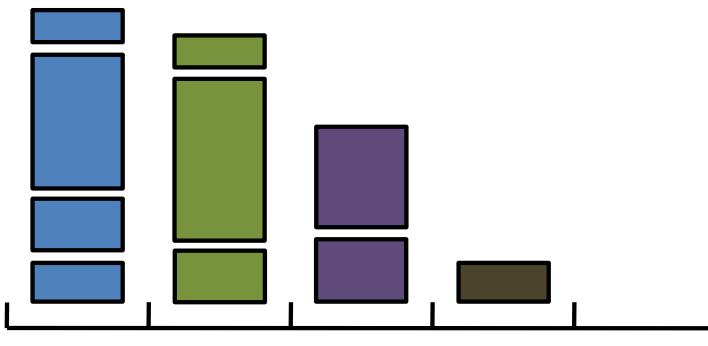
$$W = \sum_i w_i$$

n



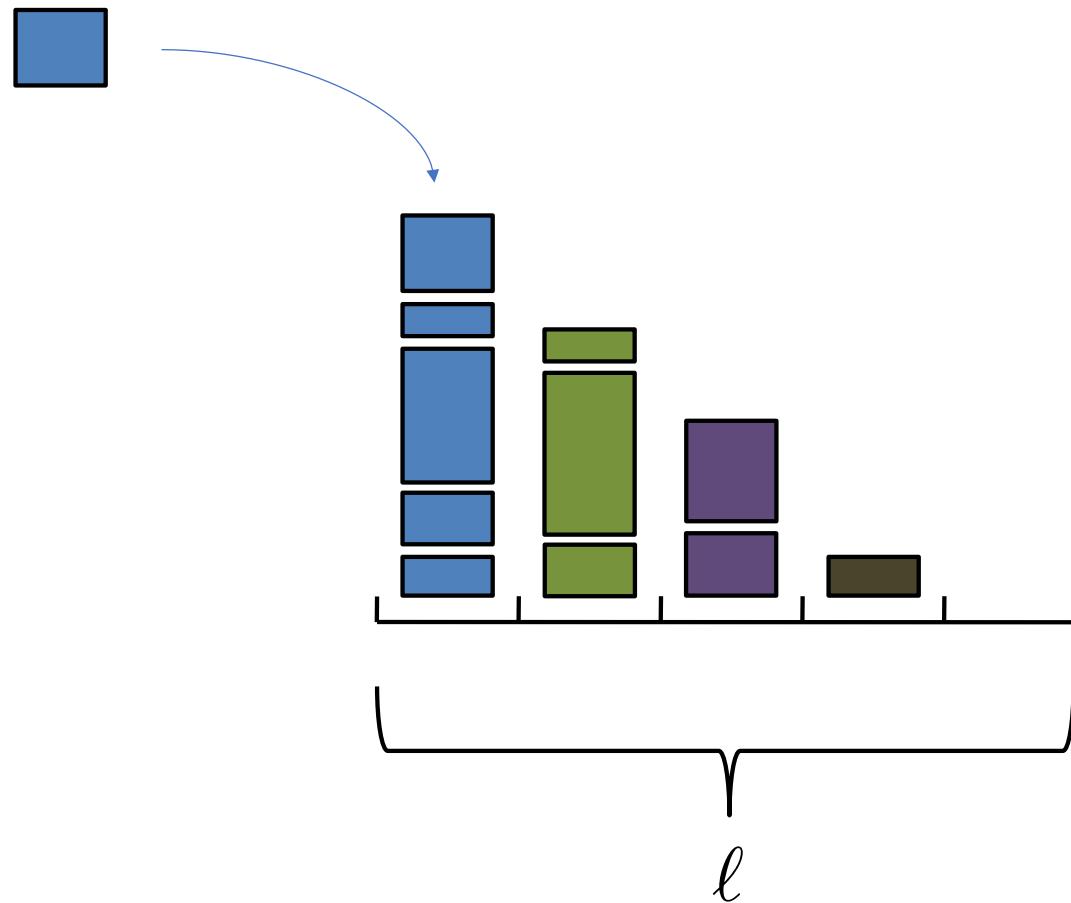
Our Contributions

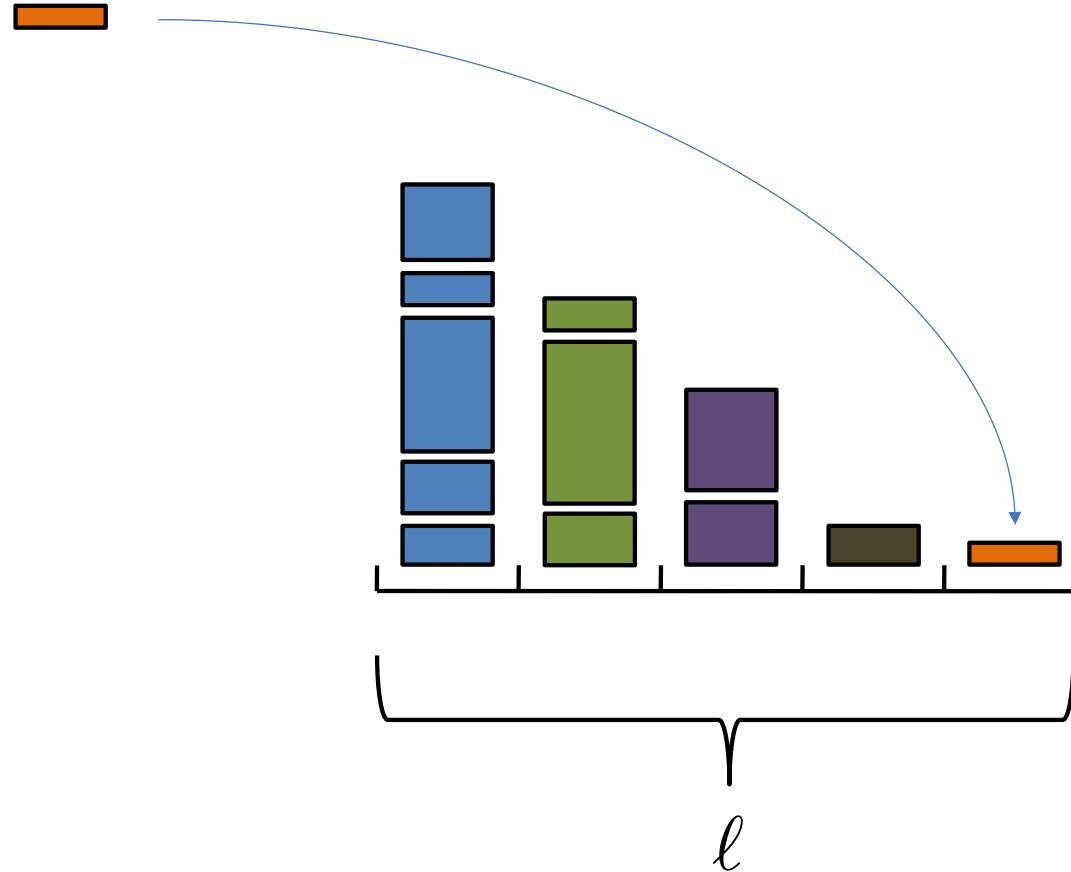
- Improved streaming algorithm for weighted updates
- Improved merging procedure
- Improved Estimator
- Careful implementation

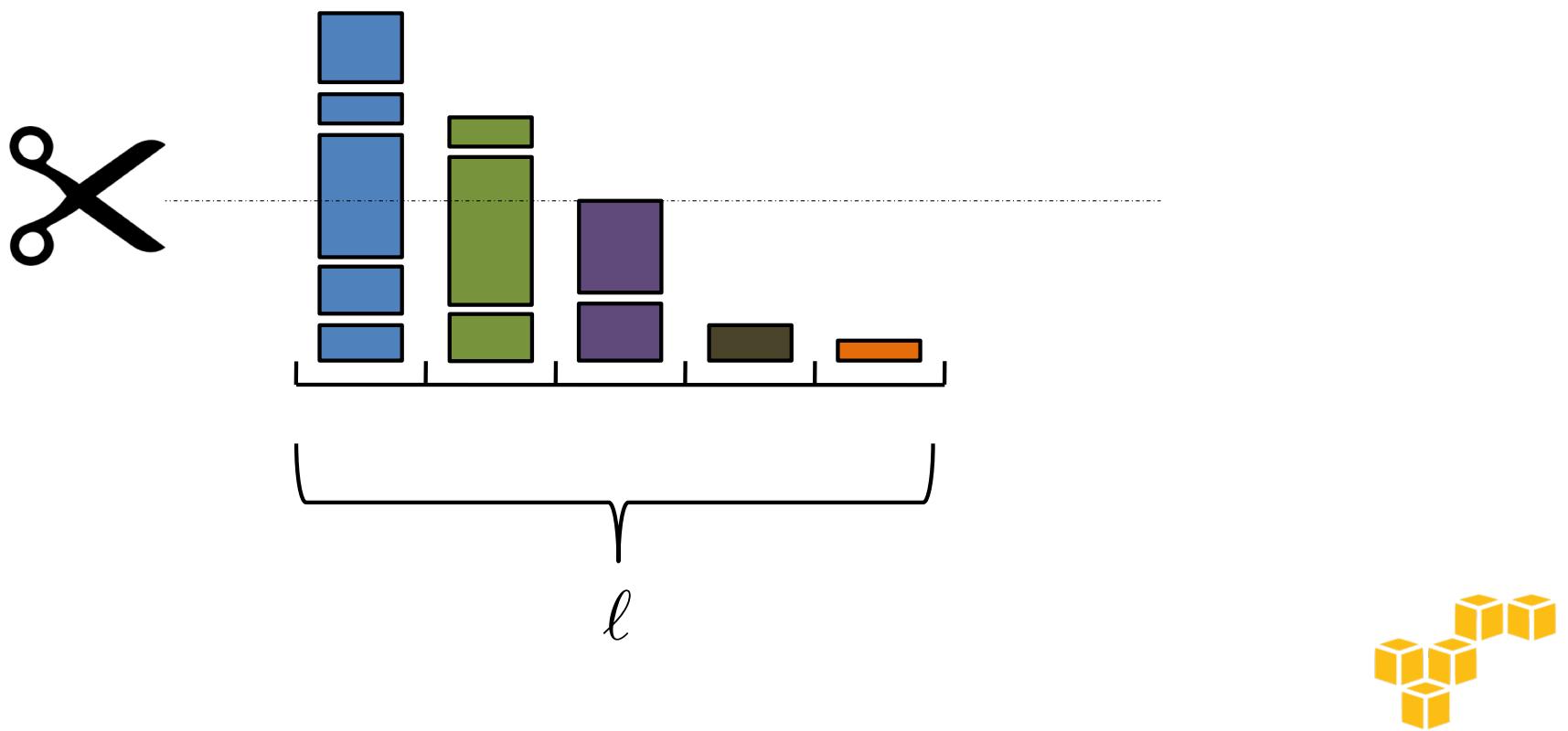


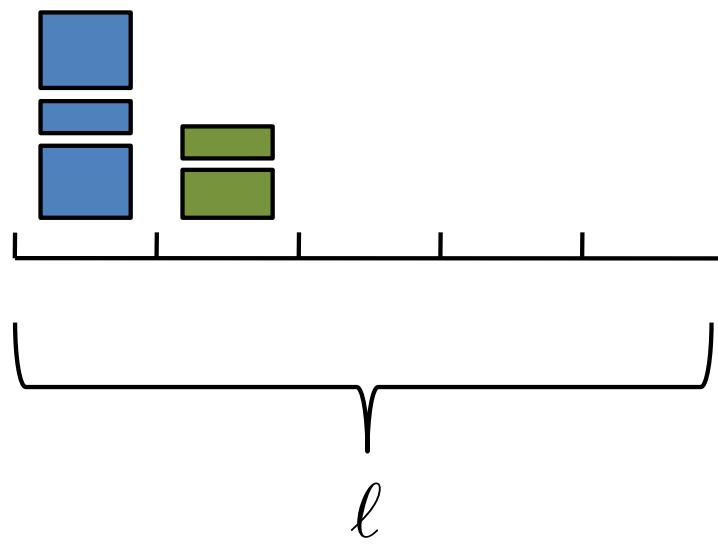
ℓ





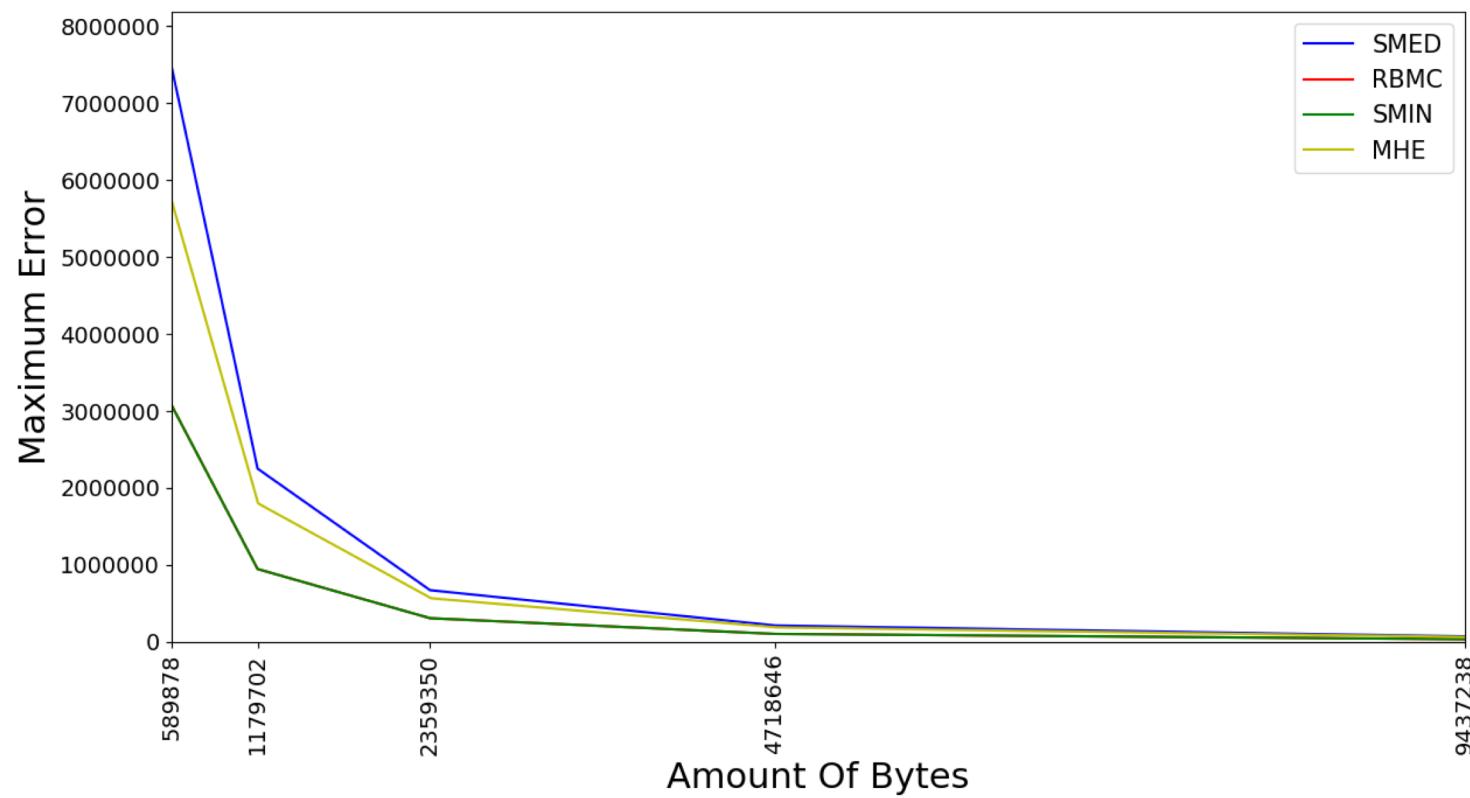






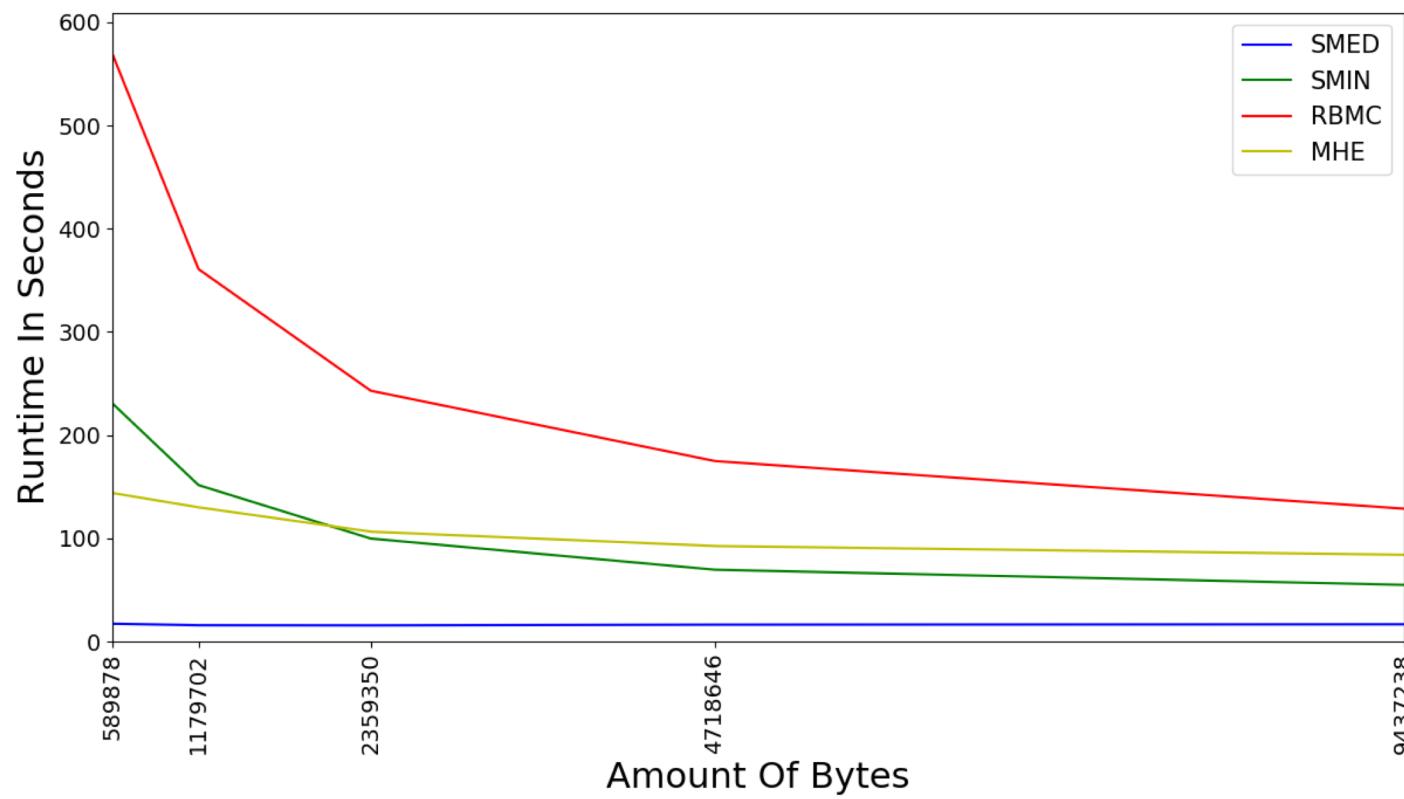
Comparable Error

Error With Equal Space



Significantly Faster Updates

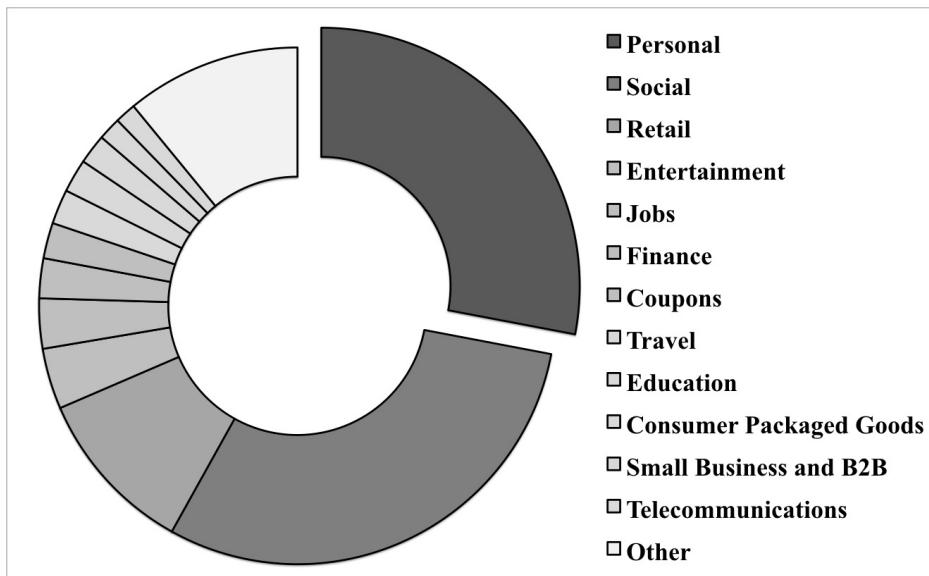
Runtime With Equal Space



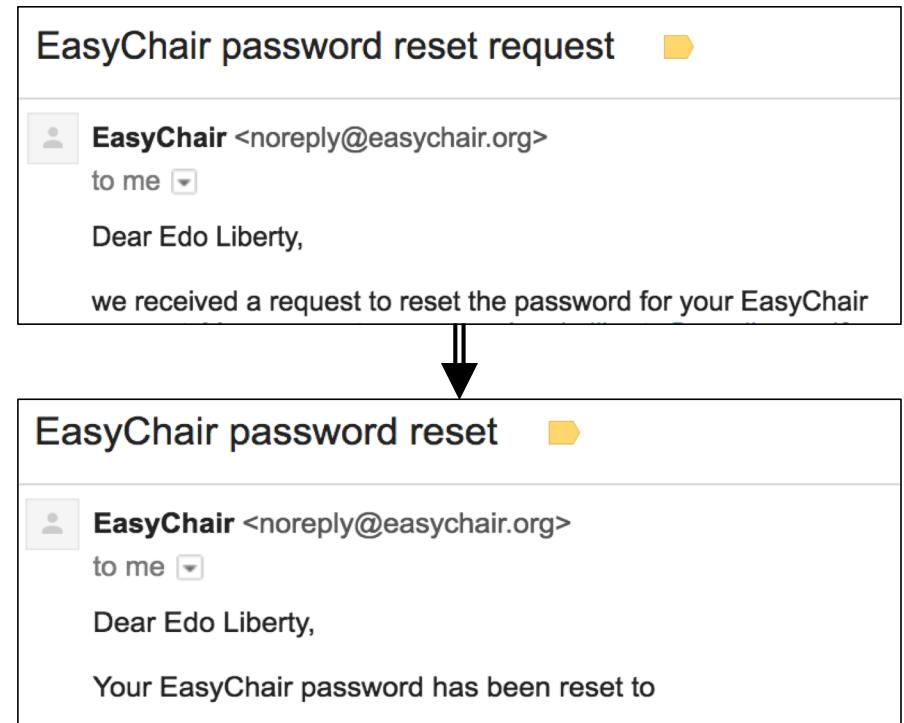
Weighted Item Frequencies Application



Threading Machine Generated Email



Ailon, Karnin, Maarek, Liberty,
Threading Machine Generated Email, WSDM 2013



Threading Machine Generated Email

If b should be threaded with a then the lift should be large

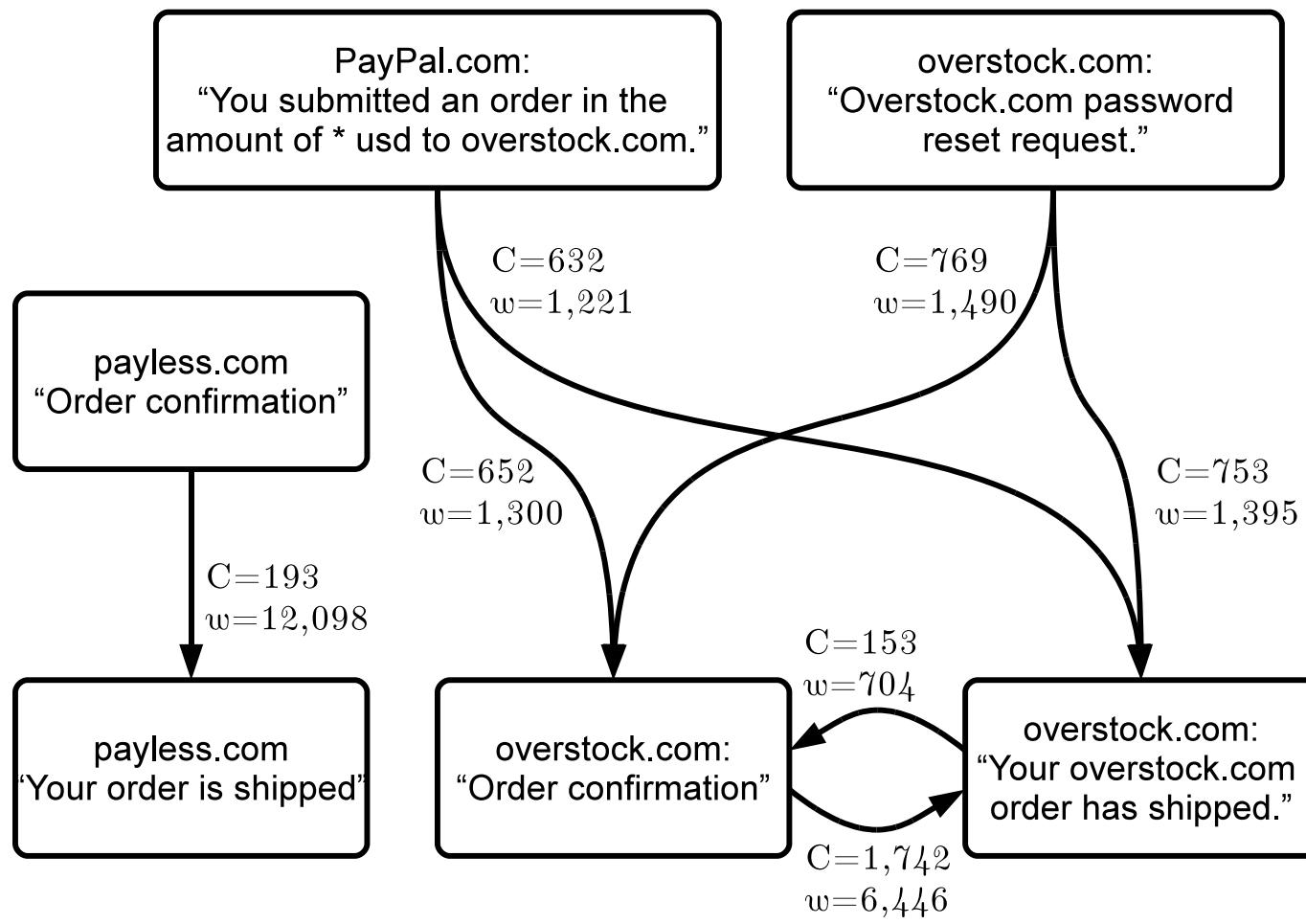
$$\text{lift}(a, b) = p(b|a)/p(b) \gg 1$$

Alas, computing all pair conditional probability is impossible!

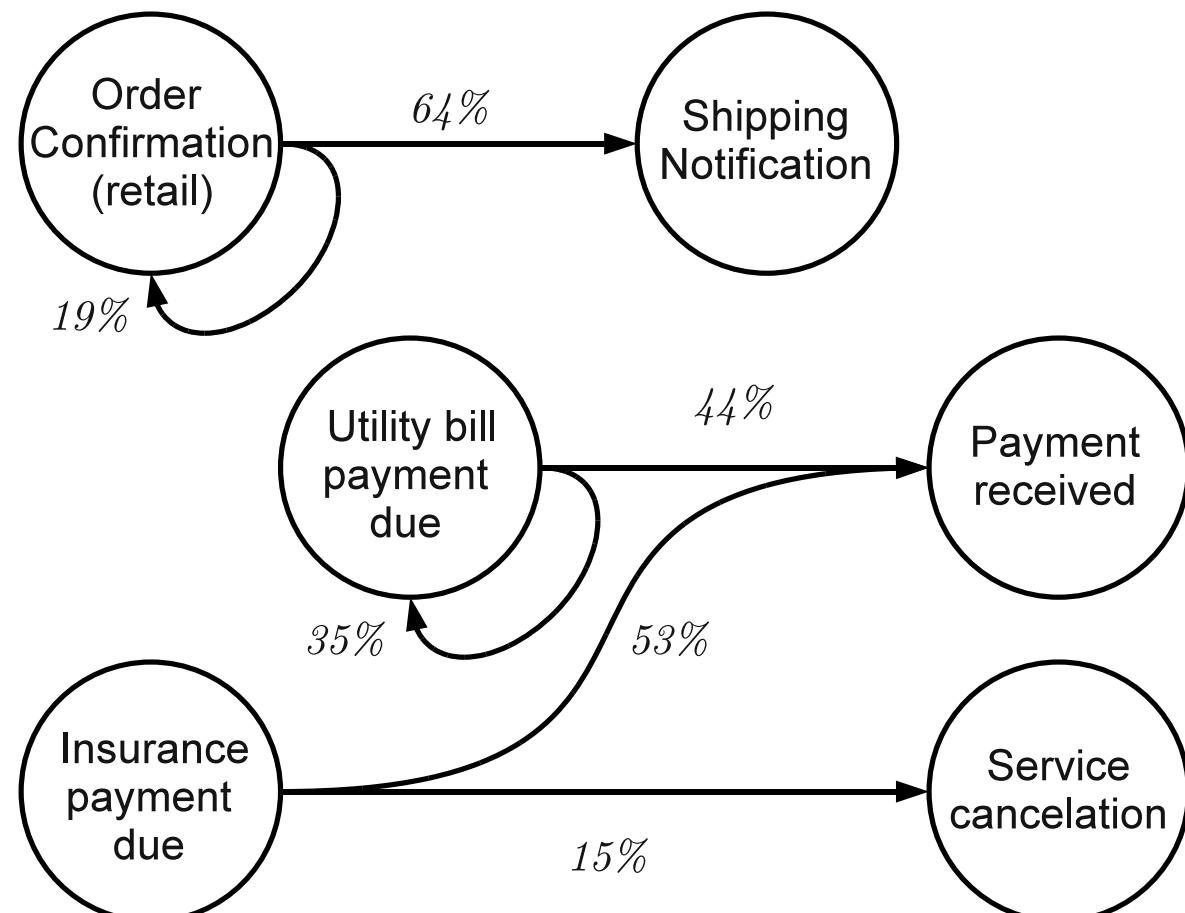
$$\text{lift}(a, b) = n(b, a) \frac{n}{n(b)n(a)} = n(b, a)w(b)$$

This is possible with weighted frequency sketching!

Threading Machine Generated Email



Threading Machine Generated Email



Streaming quantiles

Manku, Rajagopalan, Lindsay. Random sampling techniques for space efficient online computation of order statistics of large datasets.

Munro, Paterson. Selection and sorting with limited storage.

Greenwald, Khanna. Space-efficient online computation of quantile summaries.

Wang, Luo, Yi, Cormode. Quantiles over data streams: An experimental study.

Greenwald, Khanna. Quantiles and equidepth histograms over streams.

Agarwal, Cormode, Huang, Phillips, Wei, Yi. Mergeable summaries.

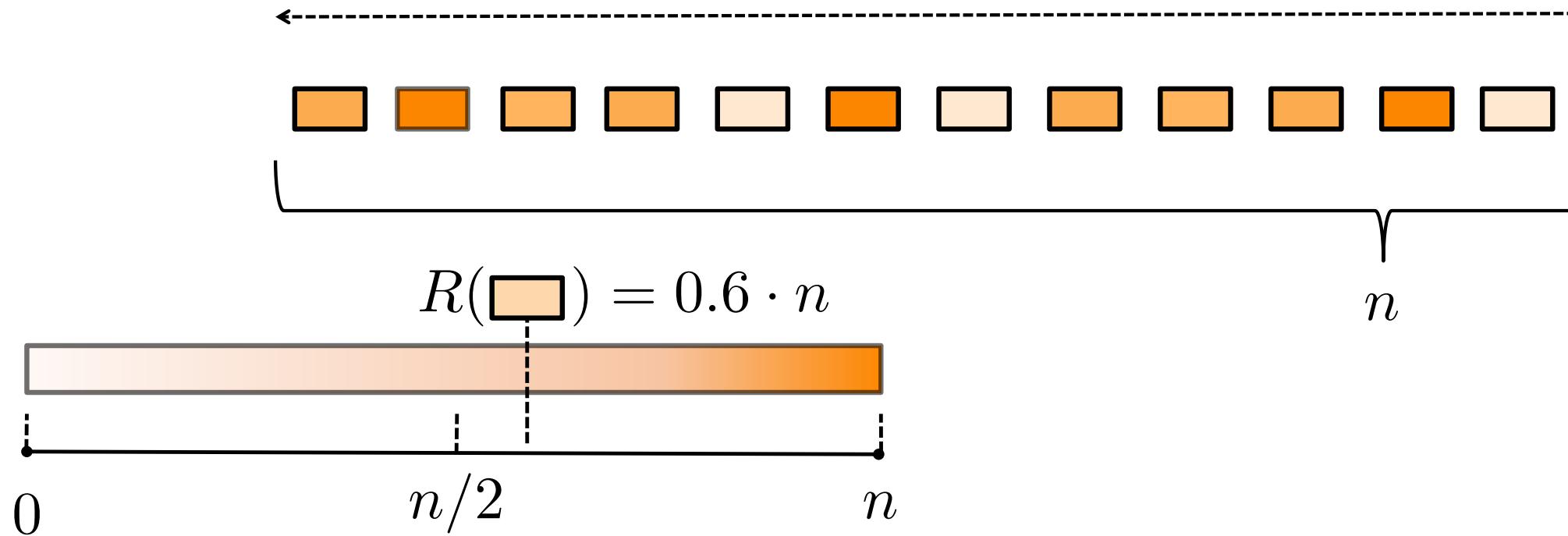
Felber, Ostrovsky. A randomized online quantile summary in $O((1/\varepsilon) \log(1/\varepsilon))$ words.

Lang, Karnin, Liberty, Optimal Quantile Approximation in Streams.

Ivkic, Lang, Karnin, Liberty, Braverman, Streaming quantiles algorithms with small space and update time



Problem Definition



Sketch the stream to estimate $|R' - R| < \varepsilon n$

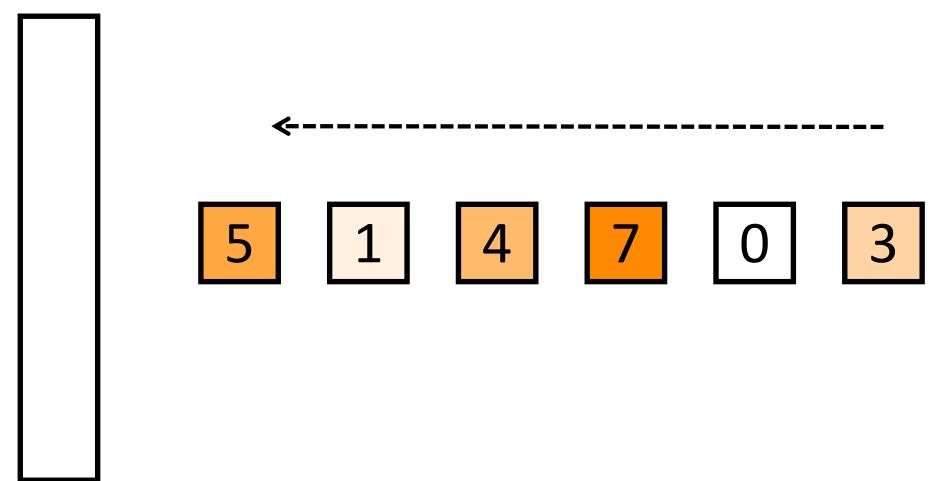


Our result

Algorithm	Simple?	Mergeable	Space Complexity	
Uniform sampling	✓	✓	$1/\varepsilon^2$	
Greenwald Khanna (GK)	✗	✗	$\log(n)/\varepsilon$	
Felber-Ostrovsky	✗	✗	$\log(1/\varepsilon)/\varepsilon$	
Manku-Rajagopalan-Lindsay (MRL)	✓	✓	$\log^2(n)/\varepsilon$	
Agarwal, Cormode, Huang, Phillips, Wei, Yi	✗	✓	$\log^{3/2}(1/\varepsilon)/\varepsilon$	Implemented
Karnin, Lang, Liberty	✓	✓	$\sqrt{\log(1/\varepsilon)}/\varepsilon$	Implemented
Karnin, Lang, Liberty	✗	✗	$\log^2 \log(1/\varepsilon)/\varepsilon$	
Karnin, Lang, Liberty	✗	✗	$\log \log(1/\varepsilon)/\varepsilon$	Space Optimal
Still open...	✓	✓	$\log \log(1/\varepsilon)/\varepsilon$	

The basic buffer idea

Buffer of size k



The basic buffer idea

Stores k stream entries



The basic buffer idea

The buffer sorts k stream entries



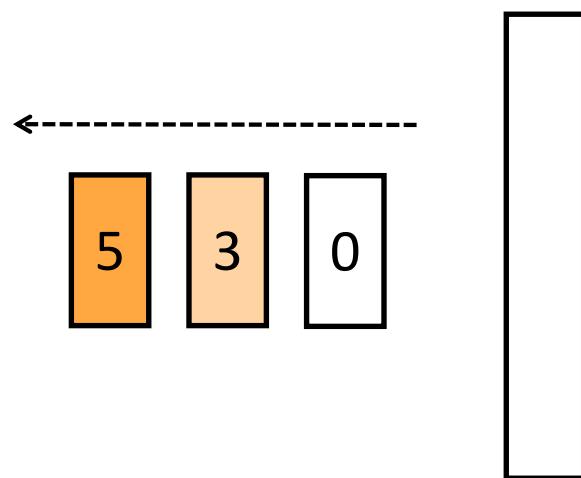
The basic buffer idea

Deletes every other item

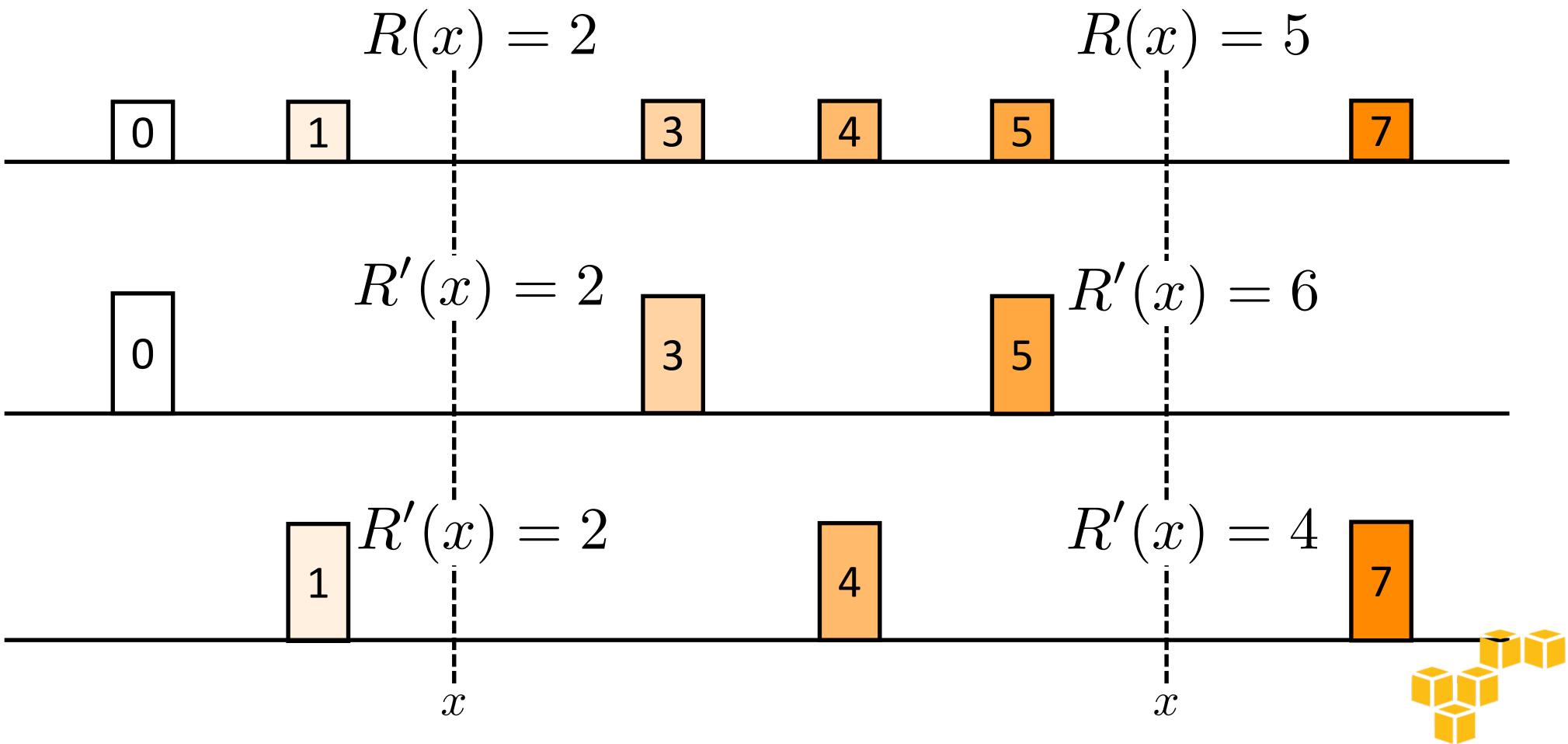


The basic buffer idea

And outputs the rest
with double the weight

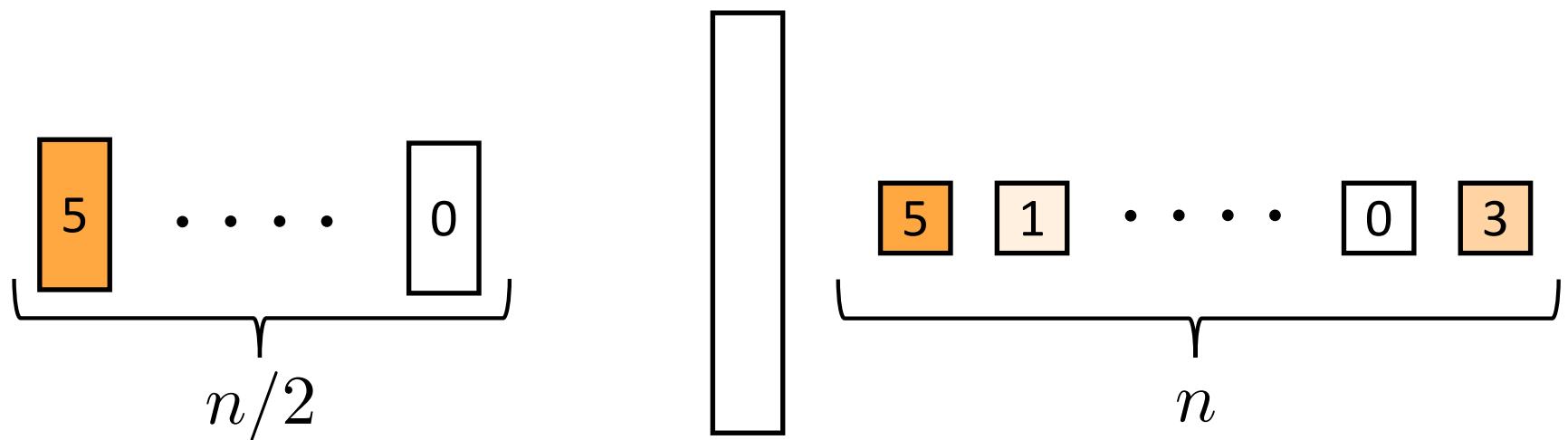


The basic buffer idea



The basic buffer idea

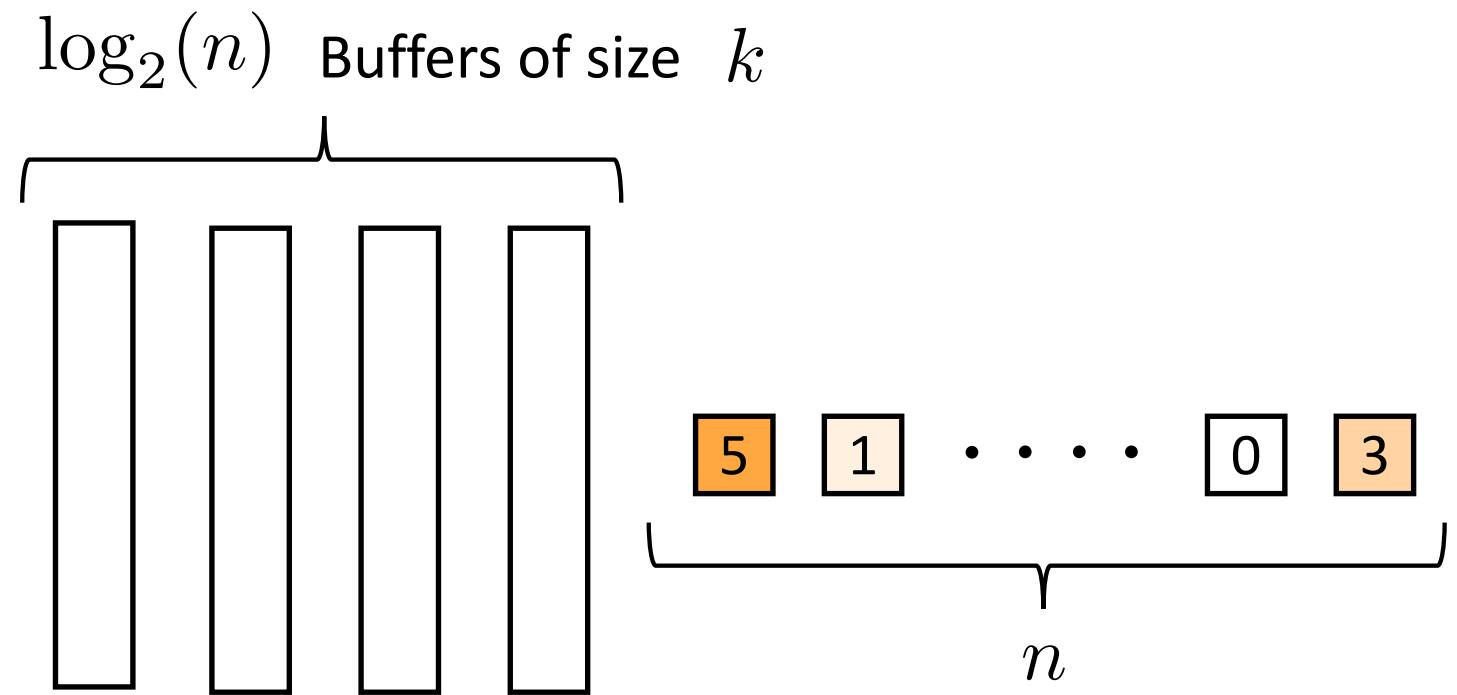
Repeat n/k time until
the end of the stream



$$|R'(x) - R(x)| < n/k$$



Manku-Rajagopalan-Lindsay (MRL) sketch



$$|R'(x) - R(x)| \leq n \log_2(n)/k$$



Manku-Rajagopalan-Lindsay (MRL) sketch

If we set $k = \log_2(n)/\varepsilon$

We get $|R'(x) - R(x)| \leq \varepsilon n$

And we maintain only $\log_2^2(n)/\varepsilon$ items from the stream!



Greenwald-Khanna (GK) sketch

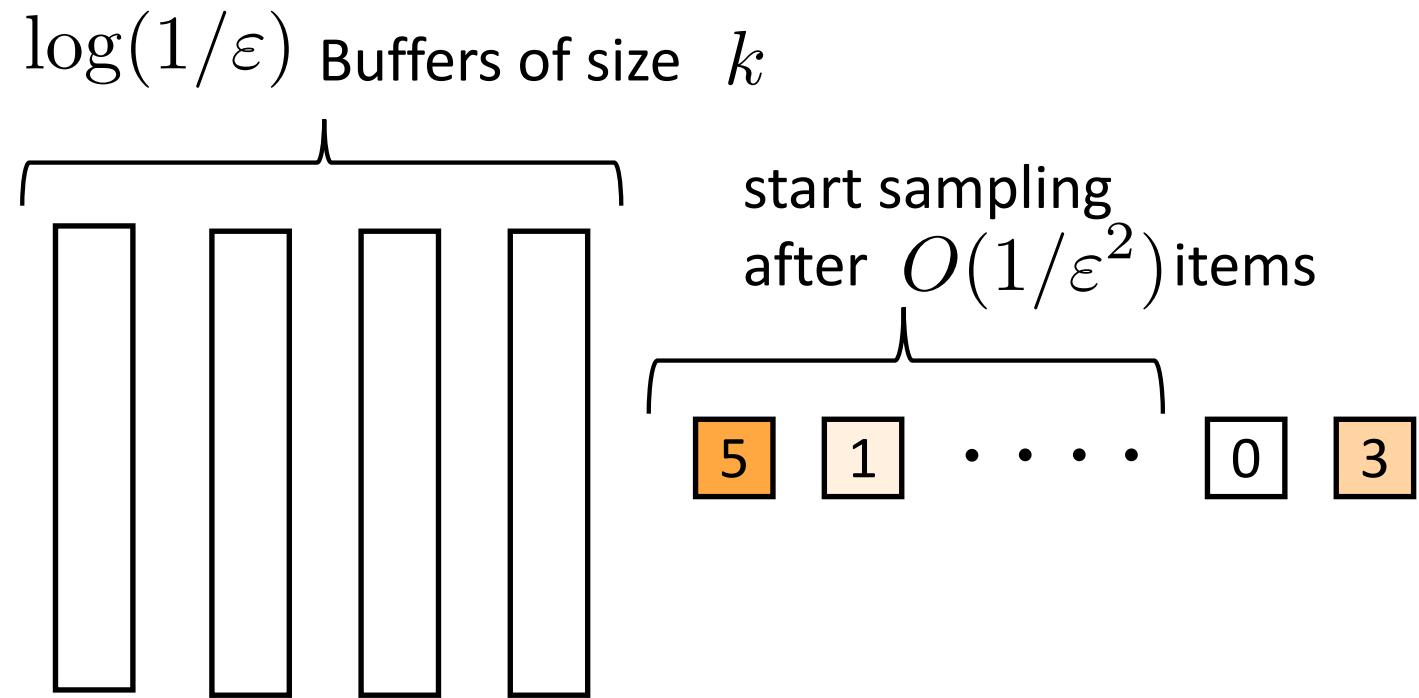
Uses a completely different construction

It gets $|R'(x) - R(x)| \leq \varepsilon n$

And maintains only $O(\log(n)/\varepsilon)$ items from the stream!



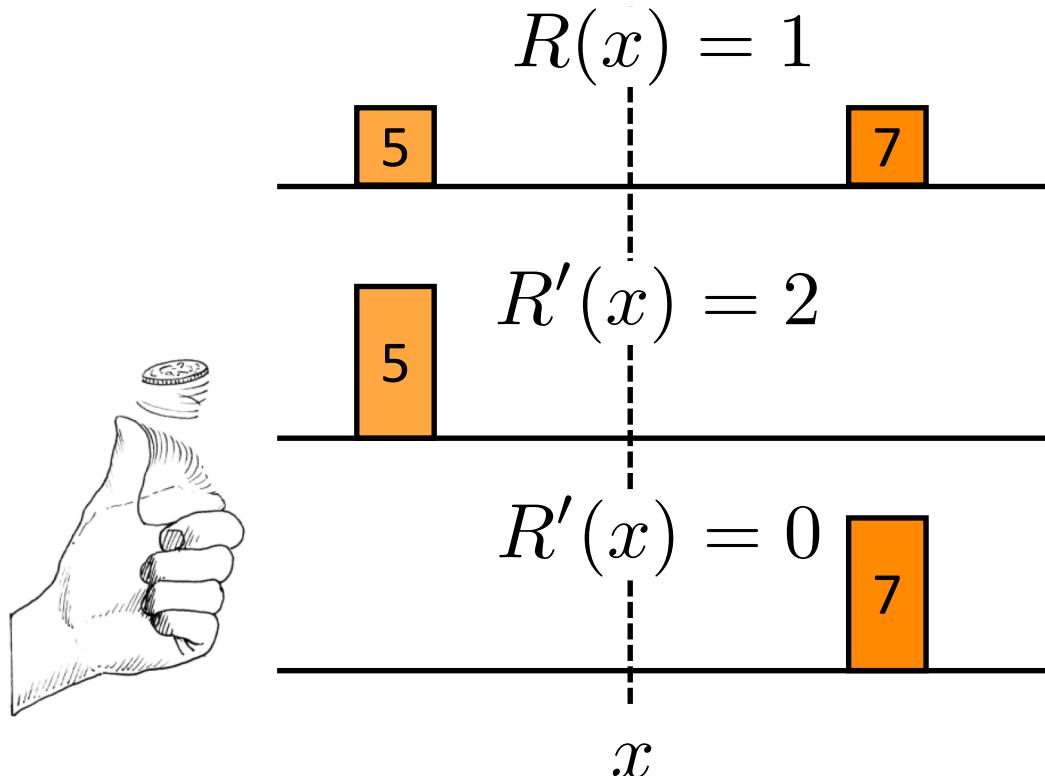
Agarwal, Cormode, Huang, Phillips, Wei, Yi (1)



Reduces space usage to $\log^2(1/\varepsilon)/\varepsilon$ items from the stream.



Agarwal, Cormode, Huang, Phillips, Wei, Yi (2)



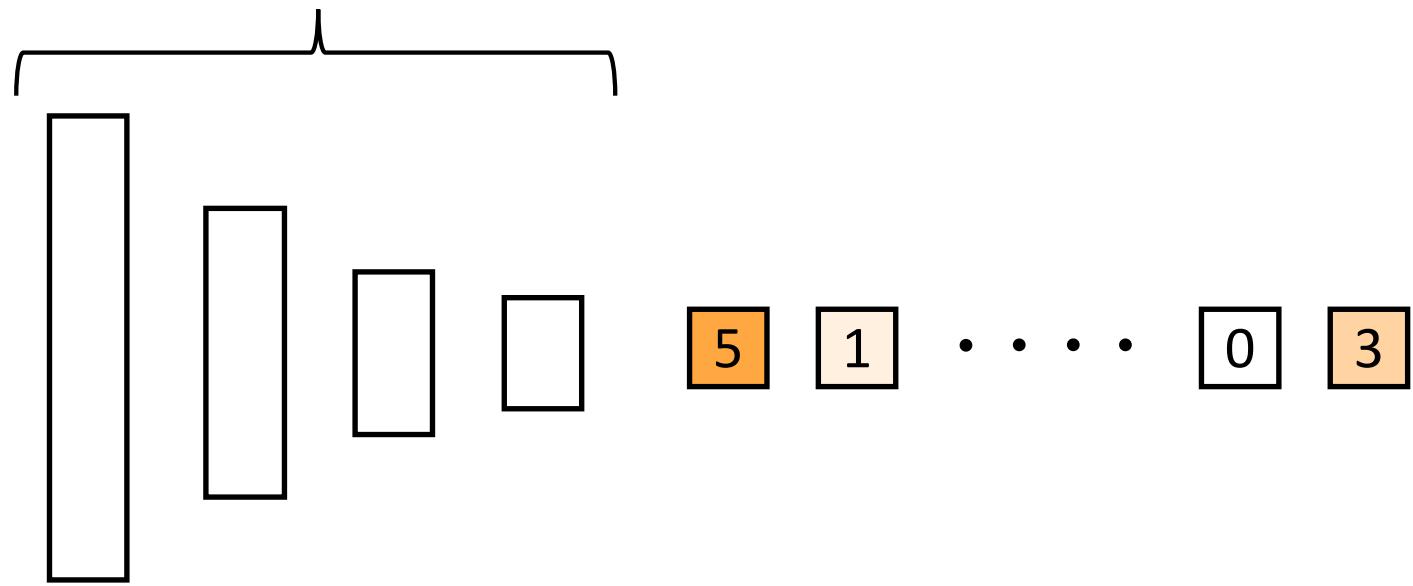
$R'(x)$ is a random variable now and
 $E[R'(x)] = R(x)$

Reduces space usage to $\log^{3/2}(1/\varepsilon)/\varepsilon$ items from the stream.



Lang, Karnin, Liberty (1)

Exponentially shrinking buffers

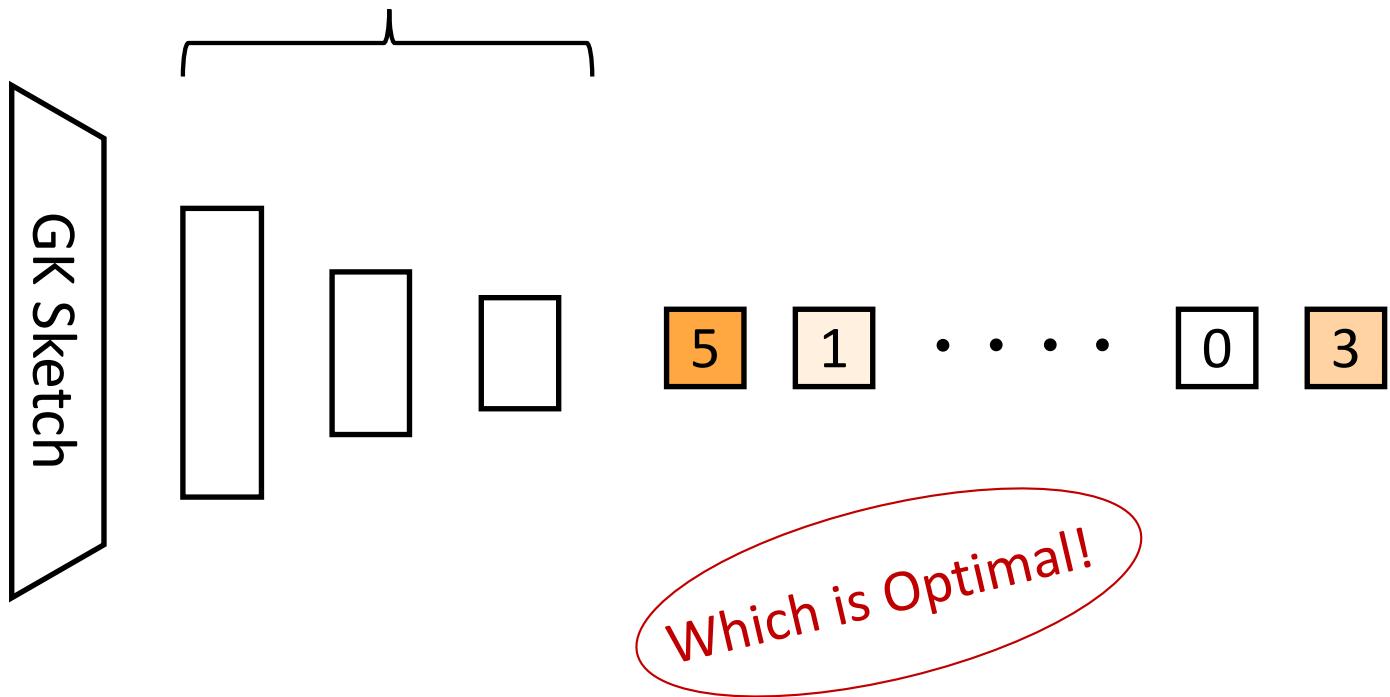


Reduces space usage to $\sqrt{\log(1/\varepsilon)}/\varepsilon$ items from the stream.



Lang, Karnin, Liberty (2)

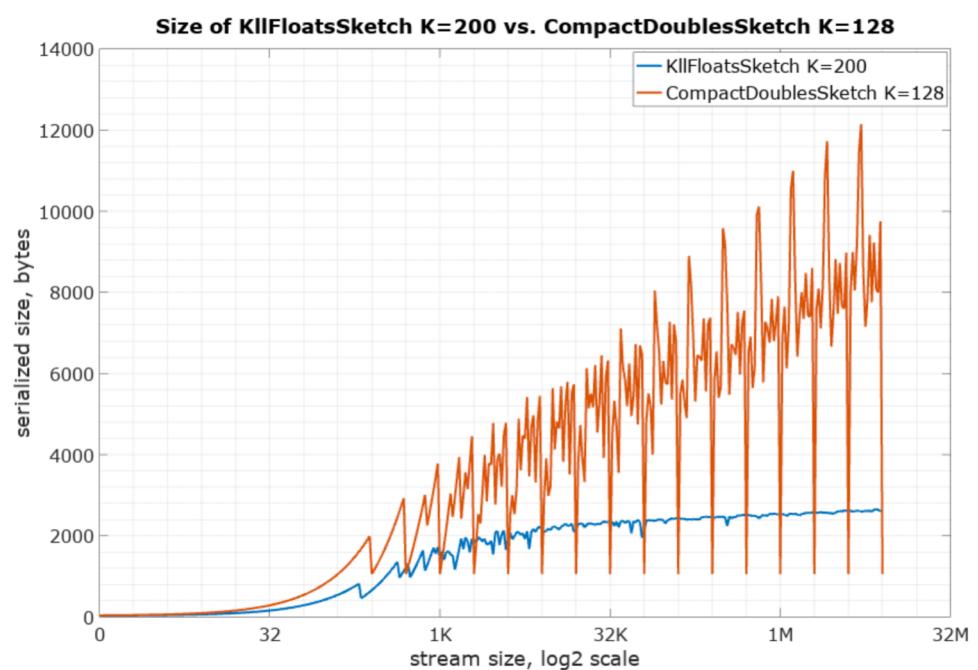
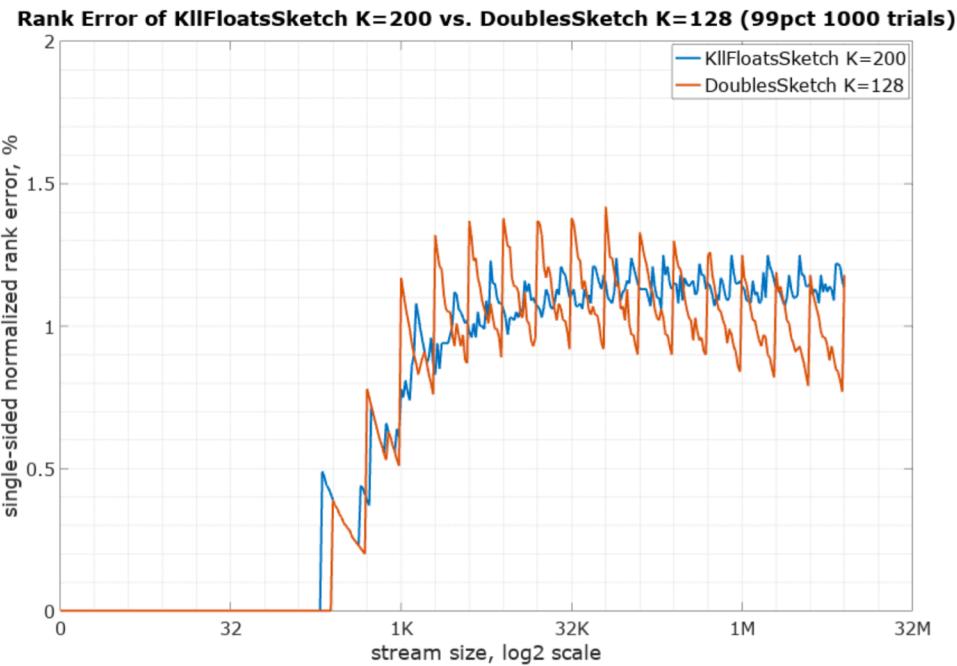
Exponentially decreasing buffer sizes



Reduces space usage to $\log \log(1/\varepsilon)/\varepsilon$ items from the stream.

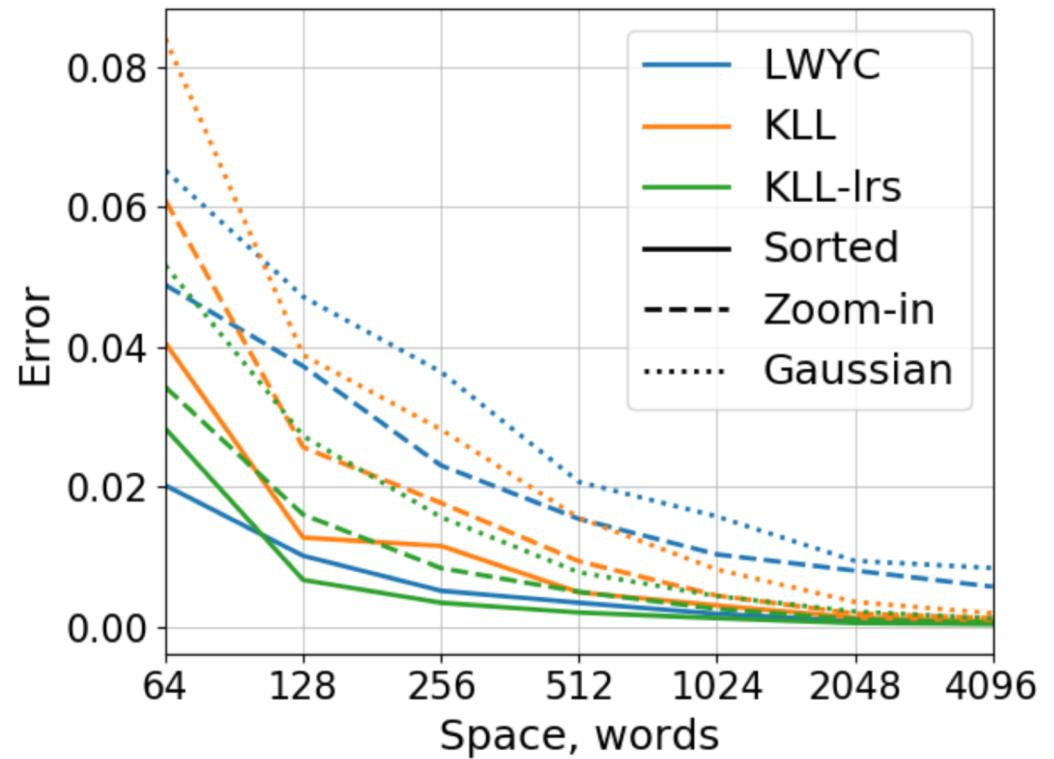
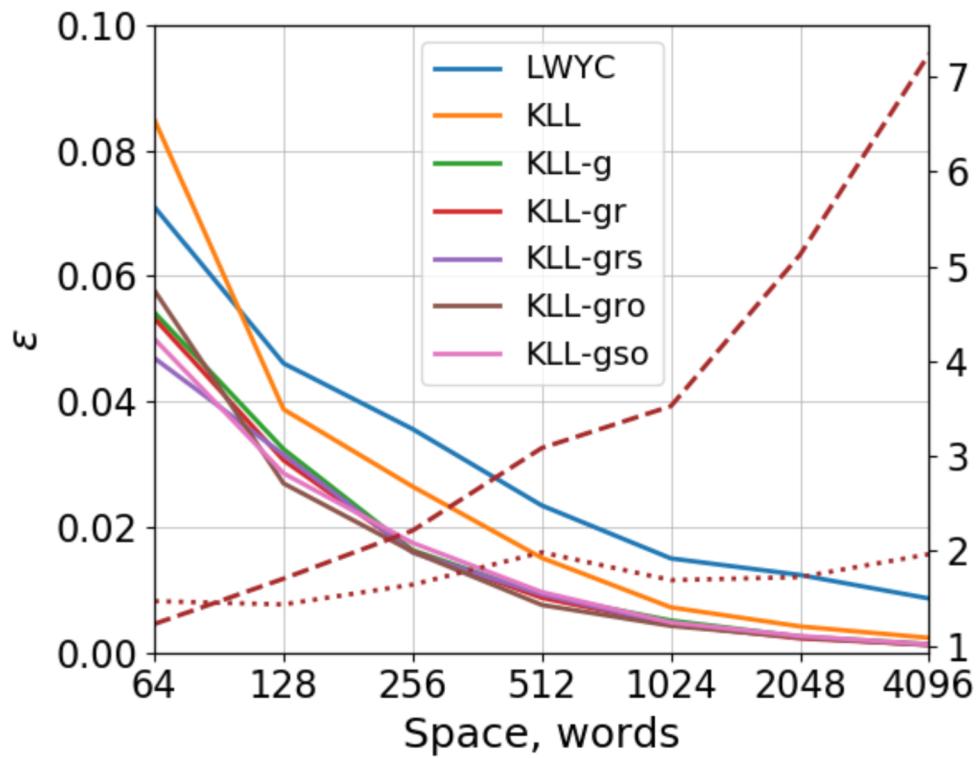


Experimental Results



More experiments: <https://datasketches.github.io/docs/Quantiles/KLSSketch.html>
<https://datasketches.github.io/docs/Quantiles/KIISketchVsTDigest.html>

Even Newer Experimental Results



Ivking, Lang, Karnin, Liberty, Braverman, Streaming quantiles algorithms with small space and update time

What else can we do in this model?

Items

(words, IP-addresses, events, clicks,...)

- Counting distinct elements
- Item frequencies
- Approximate Quantiles
- Moment and entropy estimation
- Approximate set operations
- Sampling

Matrices

(text corpora, recommendations, ...)

- Covariance estimation matrix
- Low rank approximation
- Sparsification

Vectors

(text documents, images, example features,...)

- Dimensionality reduction
- Clustering (k-means, k-median,...)
- Linear Regression
- Machine learning (some of it at least)
- Density Estimation / Anomaly detection

Graphs*

(social networks, communications, ...)

- Connectivity
- Cut Sparsification
- Weighted Matching



Thank you!

