

# An Improved Data Stream Summary: The Count-Min Sketch and its Applications

Graham Cormode, DIMACS  
graham@dimacs.rutgers.edu

S. Muthukrishnan, Rutgers  
muthu@cs.rutgers.edu

# Data Streams

- Data is growing fast — faster than our ability to store or compute on it.
- Information in Networks (phones, internet)  
Scientific data readings (satellites, sensor networks)  
Databases (financial transactions, etc.)
- One approach: take one pass over data, summarize for later querying (for some class of queries):  
the **data stream model**

# Data Stream Model

- Data stream represents a high-dimensional vector  $a$ , initially all zero: for  $1 \leq i \leq U$  .  $a[i] = 0$
- $n$  items in the stream:  $t$ 'th update is  $(i(t), c(t))$ , meaning  $a[i(t)]$  is updated to  $a[i] + c(t)$ .
- $c$  may be negative in some cases,  $a[i]$  may or may not be allowed to be negative  
(here, assume non-negative; general case in paper)

# Sketches

"Sketches" are a class of data stream summaries

- Typically, formed by linear projections of source data with appropriate (pseudo)random vectors
- Introduced by Alon Matias & Szegedy in 1996 for estimating  $F_2$  (later:  $L_2$  norm, inner products)
- Also: Indyk '00 for  $L_1, L_p$  norms  
Flajolet-Martin '83 for  $F_0$  (distinct items)  
Charikar, Chen, Farach-Colton for point estimates

# Limitations of Sketches

So why do we need new sketches?

- Space dependency is  $1/\epsilon^2$  for  $1 + \epsilon$  approximations: unusable for even reasonable values of  $\epsilon < 1\%$ . (for some problems  $1/\epsilon^2$  is a lower bound)
- Update time often slow (linear in space), doesn't scale to network line speeds
- Independence and randomness requirements sometimes excessive or unclear
- Sometimes limited to one application

# CM Sketch

Count-Min Sketch sets out to solve all these problems.

Gives simple, fast solutions for:

- Point Estimation (Estimate  $a[i]$ )
- Range Sums (Estimate  $\sum_{i=j}^k a[i]$ )
- Inner Products (Estimate  $\sum_i a[i] * b[i]$ )

Applications to

- Heavy Hitters (with departures)
- Dynamic Quantile Maintenance

# Point Estimation

Point Estimation: given  $i$  return an estimate of  $a[i]$ .

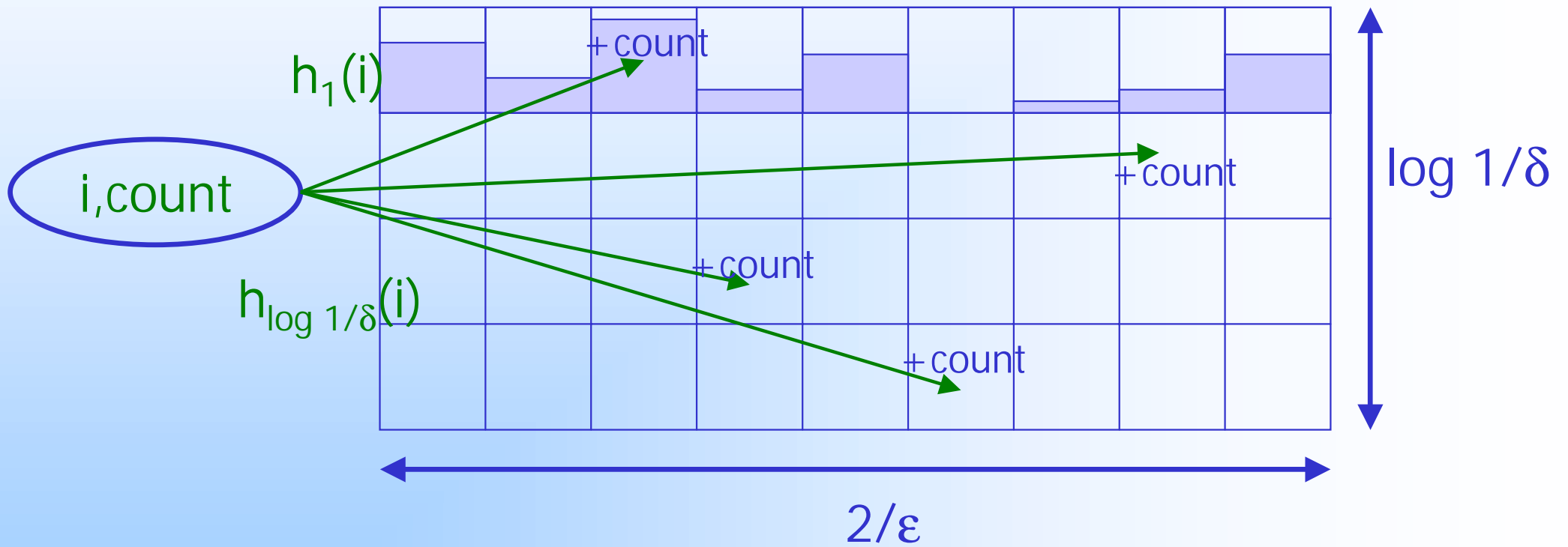
Set  $N = \sum c(t) = \|a\|_1$

Replace the vector  $a$  with small sketch which approximates all  $a[i]$  upto  $\epsilon N$  with probability  $1-\delta$

Ingredients:

- Universal hash fns  $h_1 \dots h_{\log_2 1/\delta} \{1..U\} \rightarrow \{1..2/\epsilon\}$
- Array of counters  $CM[1..2/\epsilon, 1..\log_2 1/\delta]$

# Update Algorithm



Count-Min Sketch



# Approximation

Approximate  $\hat{a}[i] = \min_j \text{CM}[h_j(i), j]$

Analysis: In  $j$ 'th row,  $\text{CM}[h_j(i), j] = a[i] + X_{i,j}$

$$X_{i,j} = \sum a[k] \mid h_j(i) = h_j(k)$$

$$\begin{aligned} E(X_{i,j}) &= \sum a[k] * \text{Pr}[h_j(i) = h_j(k)] \\ &\leq \text{Pr}[h_j(i) = h_j(k)] * \sum a[k] \\ &= \epsilon N / 2 \text{ by pairwise independence of } h \end{aligned}$$

# Analysis

$$\begin{aligned}\Pr[X_{i,j} \geq \varepsilon N] &= \Pr[X_{i,j} \geq 2E(X_{i,j})] \\ &\leq 1/2 \text{ by Markov inequality}\end{aligned}$$

$$\begin{aligned}\text{Hence, } \Pr[\hat{a}[i] \geq a[i] + \varepsilon N] &= \Pr[\forall j. X_{i,j} > \varepsilon N] \\ &\leq 1/2^{\log 1/\delta} = \delta\end{aligned}$$

Final result:

with certainty  $a[i] \leq \hat{a}[i]$  and

with probability at least  $1-\delta$ ,  $\hat{a}[i] < a[i] + \varepsilon N$

# Inner Products

- Want to estimate  $\sum a[i] * b[i]$
- Estimate with  $\min_j \sum_i CM(a)[i] * CM(b)[i]$
- Error is  $\epsilon \|a\|_1 \|b\|_1$  , similar Markov proof.
- Result from AMS96: Error  $\epsilon \|a\|_2 \|b\|_2$  with space  $1/\epsilon^2 \log 1/\delta$ .
- Which is better? Depends on distribution of  $a, b$

# Applications of CM Sketch

Heavy Hitters

Dynamic Quantiles

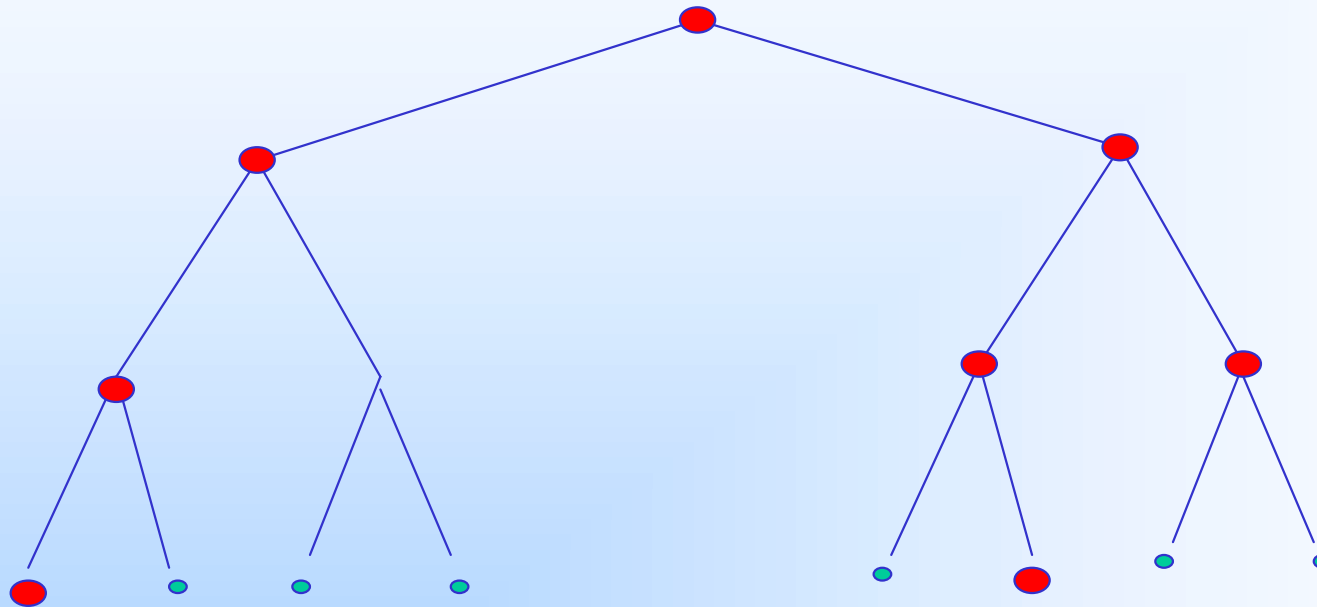
# Heavy Hitters

- See a sequence of items arriving (and departing?). Given  $\phi$ , find all items occurring more than  $\phi N$  times.
- That is, find  $i$  for which  $a[i] > \phi N$
- CCFC: Solve the arrivals only problem by remembering the largest estimated counts (in a heap) as items arrive, update sketch.
- Here: find all heavy hitters with certainty, prob  $1 - \delta$  of outputting an item with  $a[i] < (\phi - \epsilon)N$

# Solutions with Departures

- When items depart (eg deletions in a database relation), finding heavy hitters is more difficult.
- Items from the past may become heavy, following a deletion, so need to be able to recover item labels.
- Impose a (binary) tree structure on the universe, nodes correspond to **sum of counts** of leaves.
- Keep a sketch for nodes in each level and search the tree for frequent items with divide and conquer.

# Search Structure



Find all items with count  $> \phi N$  by divide and conquer  
(play off update and search time by changing degree)

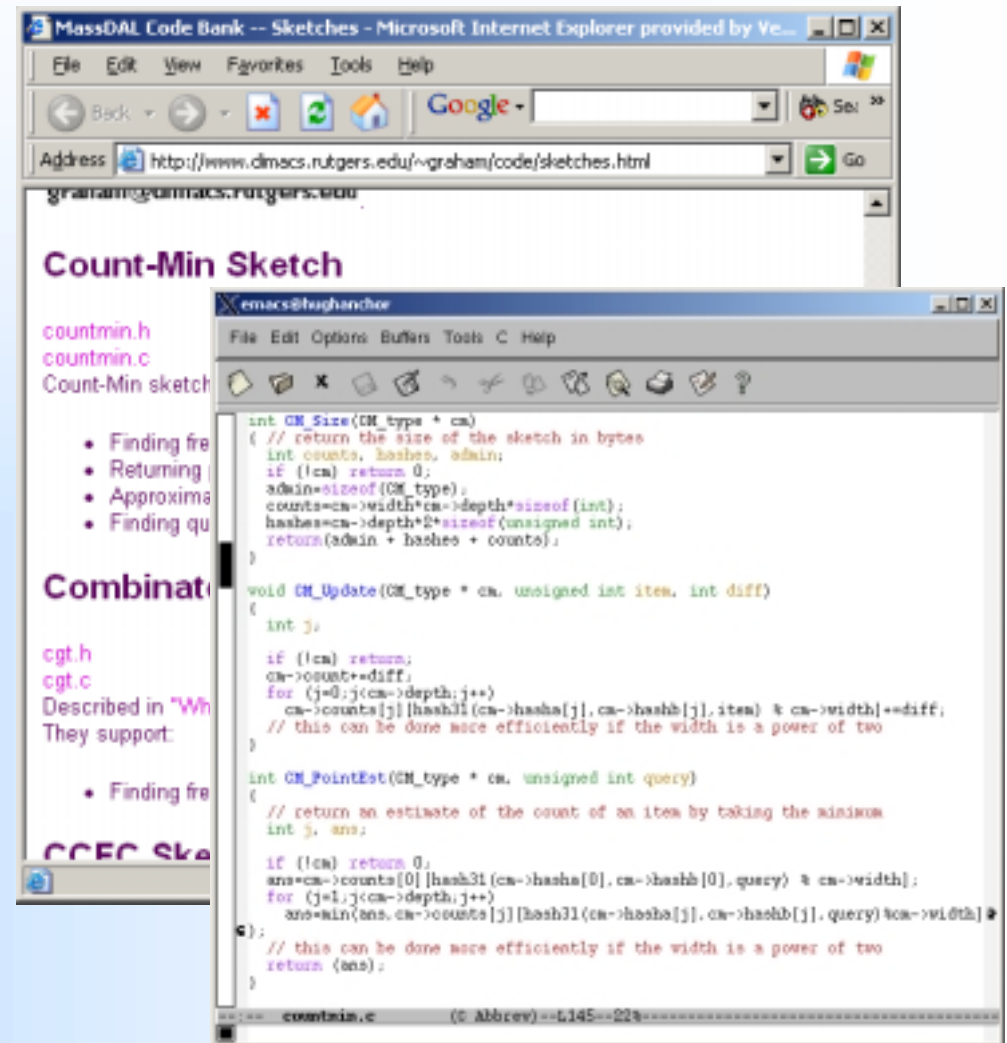
# Quantiles

- Result of GKMS02: find quantiles with range sums
- Eg Median: binary search for  $r$  so  $R(1,r) = N/2$
- Can generalize for arbitrary quantiles
- CM sketches improve space from  $O(1/\epsilon^2)$  to  $O(1/\epsilon)$
- Time is  $O(\log U \log 1/\delta)$  from  $O(1/\epsilon^2 \log^2 U \log 1/\delta)$



# Implementations

- Sketches running in AT&T Research's Gigascope network stream processing system, at 2.4Gbs
- Code for CM sketch is publicly available



<http://www.cs.rutgers.edu/~muthu/massdal-code-index.html>