STREAMING ALGORITHMS Space! -> the key resource today. Stream of . 011012, ---, 014, --- data one datapoint arrives at each time t Soul: Maintain Some function of 01, is seen up to time t Ex: ais are integers. Maintain the Sum.

Ideas maintain a "Sum" variable. Add as at time step + to the current sum. Alg: Sum < 0
attimet,
Sum < sum + at

Let's analyze how good this algorithm is.

Key resource: Space!

Assume: each ai is an integer in [1,-1,N]
Then each ai con be written in
[loy2N] bits.

How large a space do I need to Stove "Sum". ?

All integers are [1,--,N] red los N/ bits. At time T, Sum & T.N [log_TIV] = [log_T]+[log_N] Space Usage & O(log_T + (og_N) "Maximum" Space usage? Ex 2: alg: max < 0
at time t
if max < ac, max < at [log_N] else discard at Median ???? But thenk about

Since we need log_N space to even write down an element in the stream, O(logN) is the "best" case space A trivial upper bound would be D(T. logN) Costore the entire stream!) Ourgoal: fead algorithms that do

2 Thogh 2 uifact are abset o Ollog T + log N)

"Heavy Hitters" Fix E-Bruen 0,----; at, ∈ ∑ At all times t, maintain a list of all elements that have occurred >Et times. Note: do not reed a lest of Size $\geq \left| \frac{1}{2} \right| - 1$. Why? Exercise: prove that there cannot be more than [/2]-I heavy
Key Principles hitters. 1) Approxumation is skay & important. 2) Hashing/randomnens is crucial-

Heavy Hitters Def: Count (e) = {ie{1,...,t}} Def: Elemente & Elemente & Jis an Element & Taphalet of stream Element & Lis an Count (e) > E.t. Relaxation of the guevantee talse positives as the relaxation false negatives X approximation we choose" So, if there is a heavy hitter, it will be output- But all outputs need not be heavy hetters.

Finding a Majority element (= \frac{1}{2}-heavy hitter)

9: memory = 1 memory = 1 Counter=0 When element at arrives if (counter=0)9 set memory=at, counter=1 if $a_b = = memory,$ Counter = Counter+1 else Counter (- Counter - 1 discard at Creturn the element in the memory)

Analysis: Claim 1: # 2-heavy hitters = 1. at all times T. Pf: if not, there are 2 distinct elements that each occur > 1.7 times. But since total stream length is $\leq T$. Contradiction. Claim 2: At all times T>0, memory Contains the majority element if it exists. It: (Note: we make no claims if there is no Majority element) Key: When we discard an element at from the memory, we also throw away arother (different) element as well. Decrementing the counter = throwing away one capy of element in memory So each time we discard the true majority we must have thrown away another different

If majority element is not in the memory of them # elements thrown away is > 7/2+ 7/2 > T. Contradiction Space Usage: O (log2N + log2T) counter memory = sigle element Lots now go figure out E-heavy hitters.

Let k= [=]-1 T[1,-..k] = I $((1,--k)=0^k$ When element at avvives if I je & 1,--, k3s.t. at=T[j], then C[j]++. else if some Counter (C[j]==0, then T[j] < at CLJJ < 1 else: decrementall counters by I (and discord ax)-Exè check Whythis algo is some as before of E=k.

est ce) = { C[j] if c=T[j] t o o/w Lemma 1° our estimate & the true count are somewhat close to the true count (e) - est(e) < to the t Covollary: Therefore if Count (ce) > Et then est (e) > 0 Proofs the Sum of the Counters Carnot be more than t Wherever I decrement g it goes down by K+1-Therefore no countr'is dec by more

than te. District Elements Input stream: a1, az,..., at Goal: Keepa count of number of dishact elements seen so fari

Not have to show that exact count requires memory or T. Ourgoal: Reepa est-count s.t.

for some constant (>0. "Constant factor approx"

It turns out the every deterministic algorithm provably heads SL(T) memory even for the approx. guerantee But, there is a simple vandomized algorithm! Idea: Hashing. the min h (ai) min = L at teme teme t, Compute h(ax) Switch if h(at) is timier than

·M 08 Analysis of Idealized Algorithm Lennal: Let k be the number of district elements.
The probability that the min hash Value is larger than $\frac{C}{K}$ is at most 1/00 Cfor large anough Constant (>0) C/R O/K

Chance that all K numbers are assigned in [C 1] $13 \leq (1-C)^{k}$