

KHyperLogLog

Estimating Reidentifiability and Joinability of Large Data at Scale

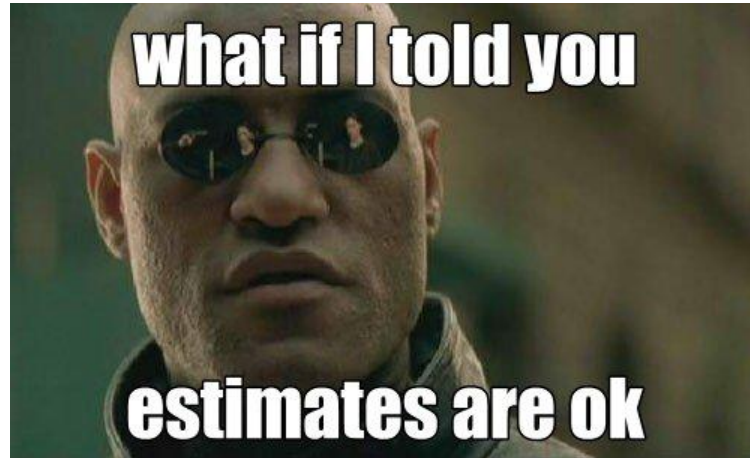
Pern Hui Chia¹, Damien Desfontaines^{1,2}, Irippuge Milinda Perera¹,
Daniel Simmons-Marengo¹, Chao Li¹,
Wei-Yen Day¹, Qiushi Wang¹, Miguel Guevara¹

¹ Google, ² ETH Zurich

Is data reidentifying?

Are data sets joinable?

Problem: Data is LARGE



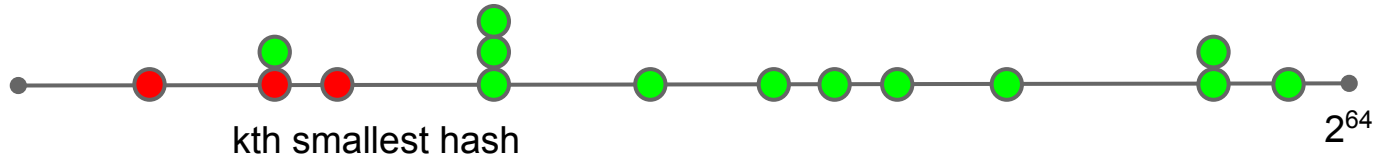
Background on **Approximate Counting**

Approximate Counting

- sublinear memory
 - compact data structures (sketches)
 - often bounded size (KBs)
- streaming algorithm
- easily parallelizable
- bounded error rate

K Min Values^[1]

- Use a uniform hash function
- Hash input values, keep **K** smallest hashes
- Extrapolate from the density of **K smallest hashes**

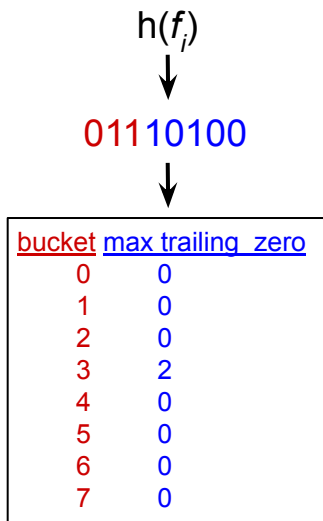


Error rate $\sim 1/\text{sqrt}(K)$

With $K=1024$, 64-bit hash: error rate $\sim 3\%$, size $\sim 8\text{KB}$

HyperLogLog_[2]

- Hash input values \Rightarrow keep max count of trailing 0's
- High number of trailing 0's \Rightarrow has seen a lot of unique values
- High variance \Rightarrow average over M counts



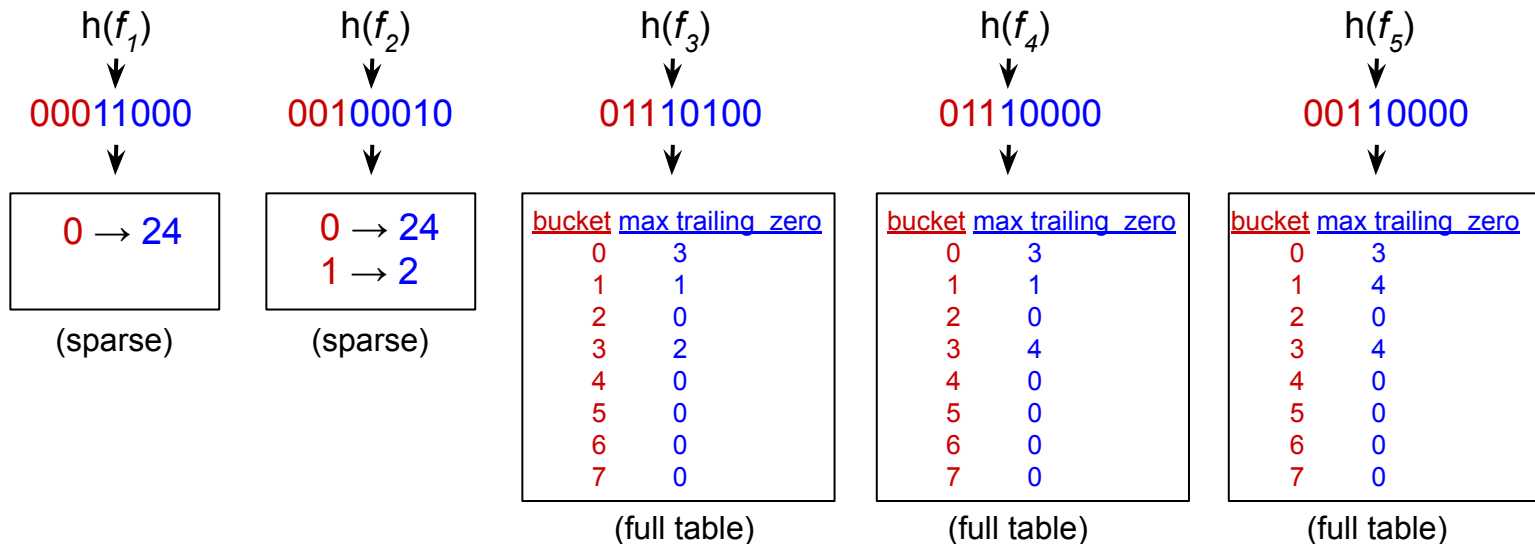
Error rate $\sim 1 / \text{sqrt}(M)$

With $M=1024$, 64-bit hash:

error rate $\sim 3\%$, size $\sim 1\text{KB}$

HyperLogLog++ ^[3]

- Sparse representation for low cardinalities
⇒ better estimates + less memory



But cardinality estimate is not enough

Introducing **KHyperLogLog**
to estimate reidentifiability and joinability

KHyperLogLog (KHLL)

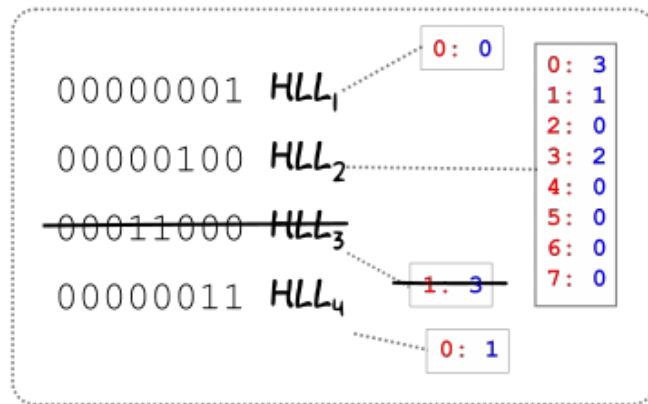
Two level counting

- 1st level → K Min Values
- 2nd level → HyperLogLog++_{HalfByte}

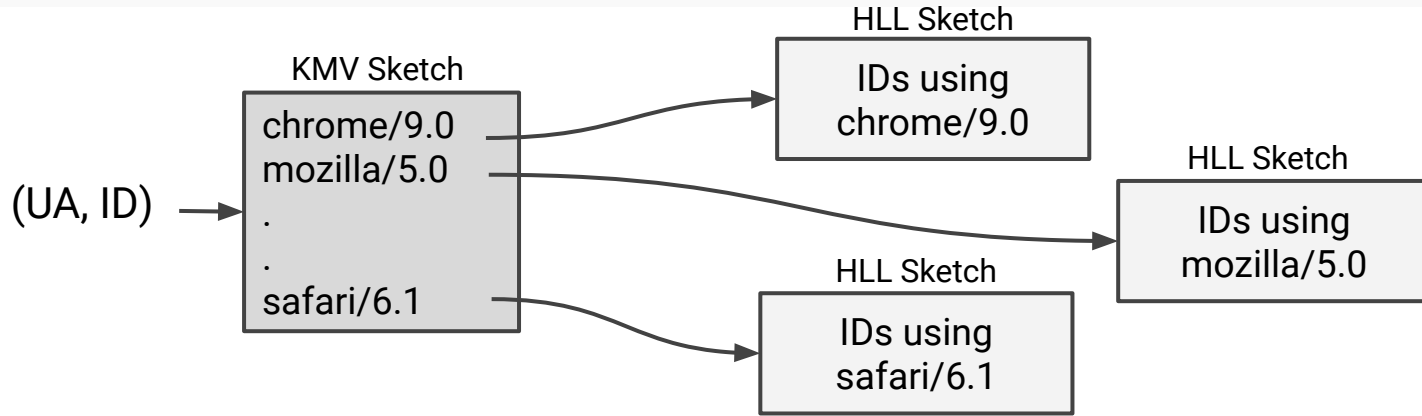
Ex: estimate uniqueness of User Agent (UA) strings

h (UA)	h (ID)
00000001	00000011
00000100	00011000
00000100	00100010
00000100	01110100
00011000	00101000
00000011	00011010

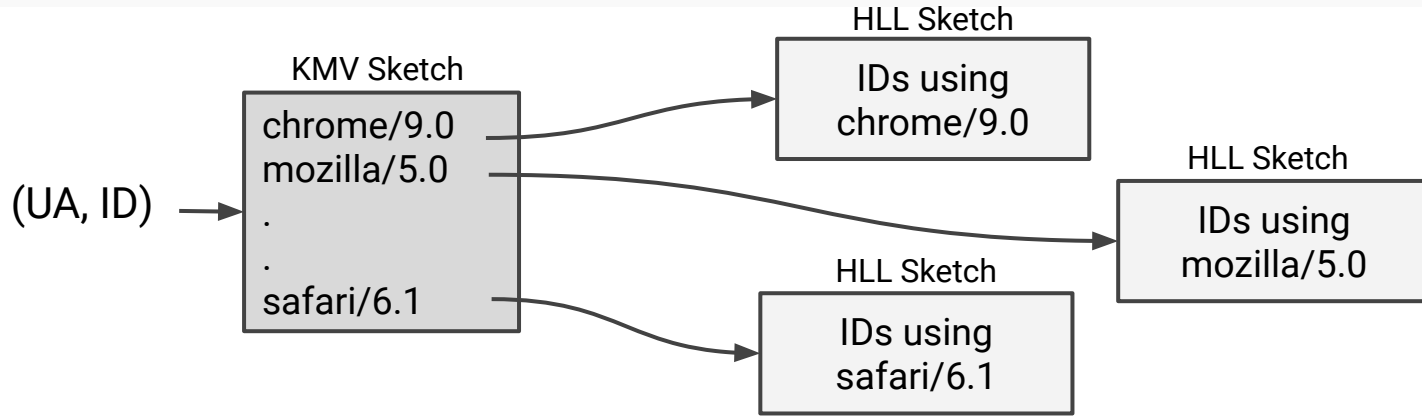
KHYPERLOGLOG, K=3, M=8



Uniqueness Distribution

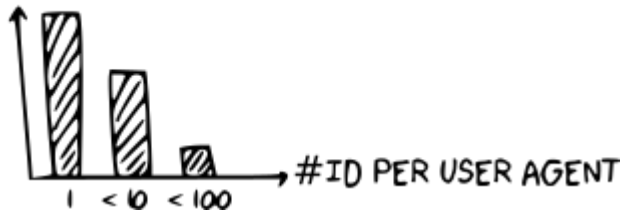


Uniqueness Distribution

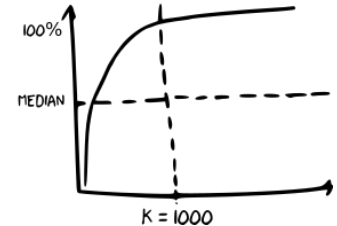
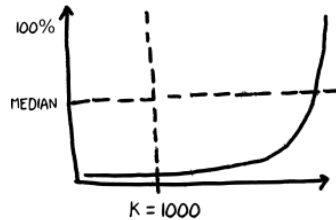


From a KHLL sketch, we can estimate

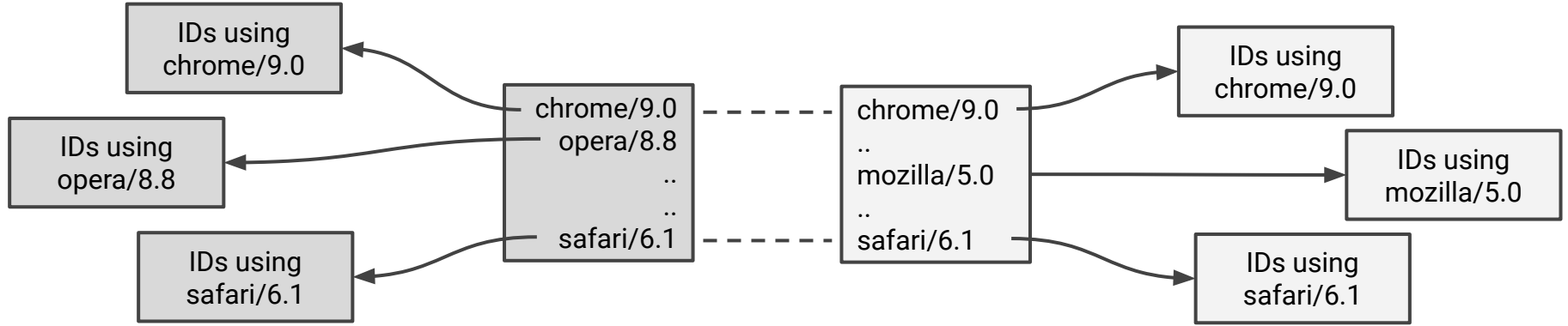
% USER AGENT



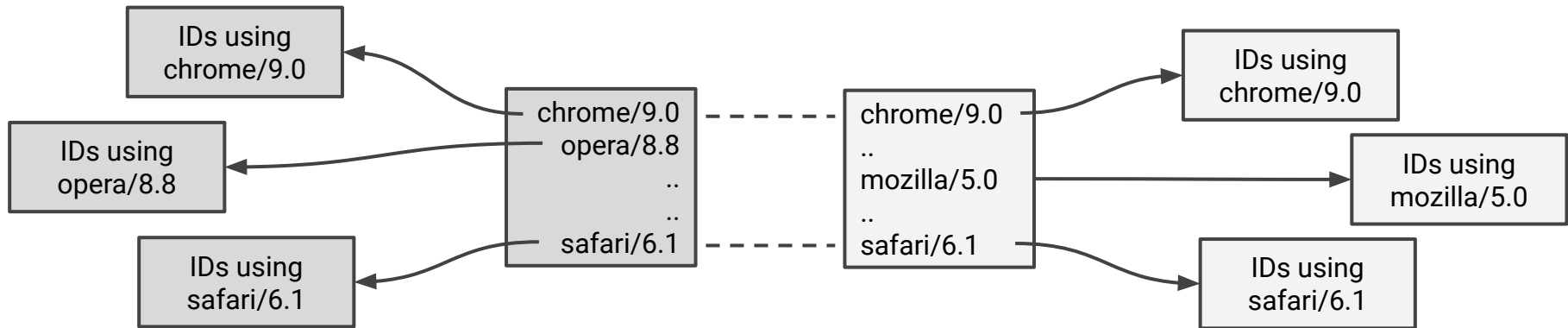
% VALUES NOT K-ANONYMOUS



Containment & Joinability



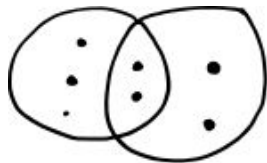
Containment & Joinability



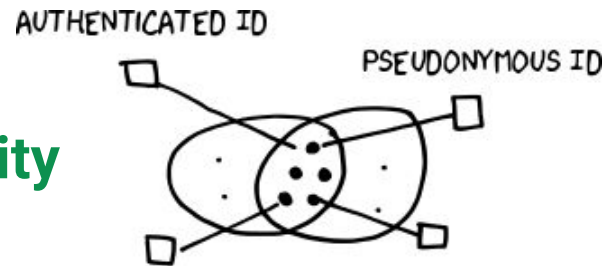
Given 2 KHLL sketches, we can estimate

Containment

$$= |A \cap B| / |A|$$



Joinability



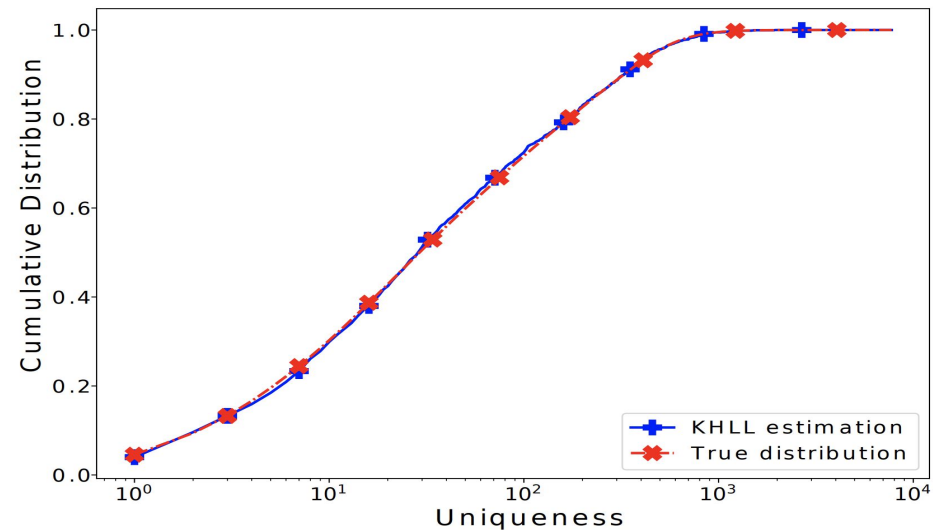
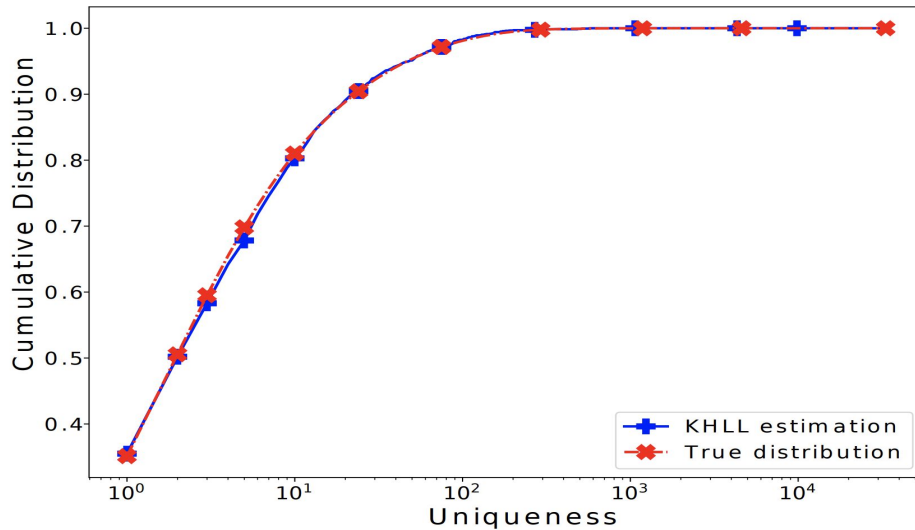
How does **KHyperLogLog** perform?

Efficiency

- KHLL vs exact counting
 - CPU performance 1 or 2 orders of magnitude faster
 - Bounded memory usage
 - With $K=1024$, $M=512$, size $\sim 256\text{KB}$
 - Scale to large data sets

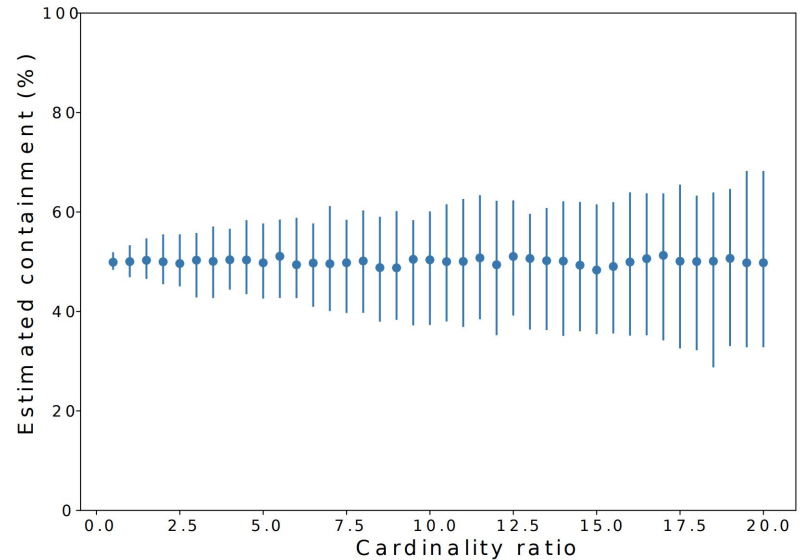
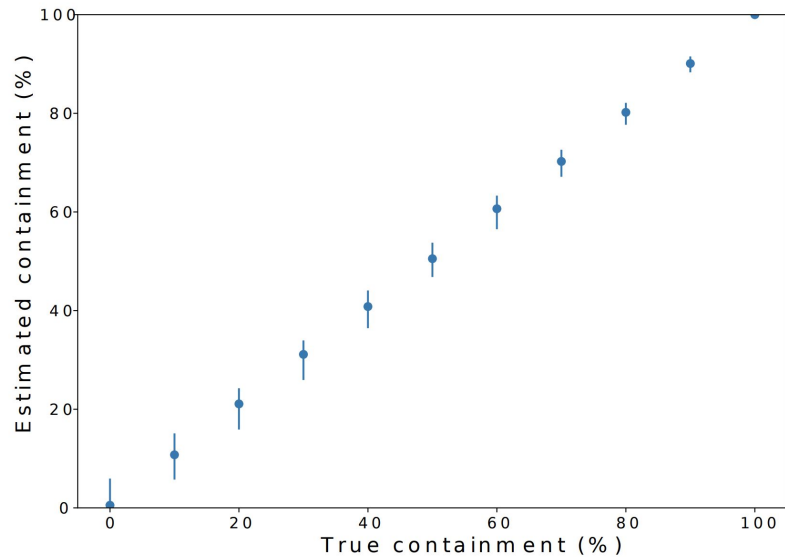
Accuracy: Uniqueness Distribution

- Prepared test distributions using Netflix & US census data
- Validated accuracy of estimated distributions



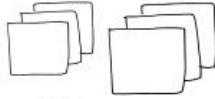
Accuracy: Containment

- Sets of equal size: $\pm 5\%$ error 90% of time
- But, problematic for sets of highly unequal sizes

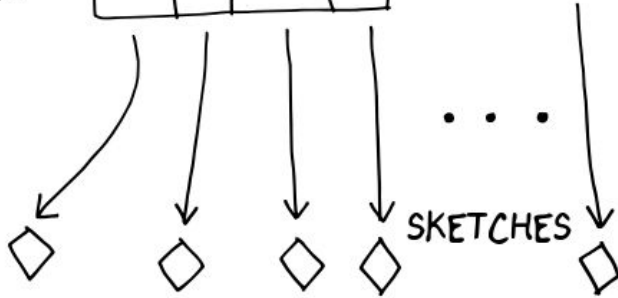


SKETCH
ALL FIELDS
ALONG
WITH IDS

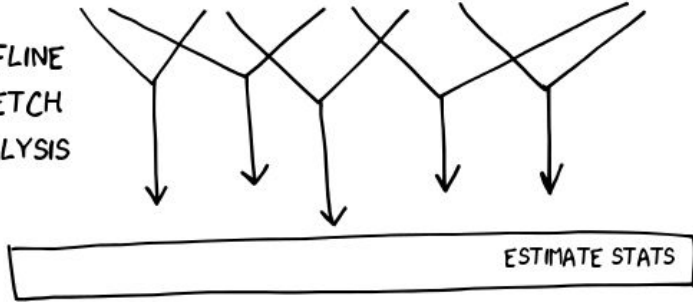
UA	IP	...	ID



--	--	--



OFFLINE
SKETCH
ANALYSIS



Practical Use of KHLL

- Read data once (petabytes)
- Produce sketches
- Offline
 - Estimate reidentifiability
 - Estimate joinability

Summary

- KHLL for estimating reidentiability & joinability at scale
- Approximation errors possible, but great for
 - Understanding data e.g., for evaluating data strategies
 - Regression testing esp. detecting data joinability

Summary

- KHLL for estimating reidentifiability & joinability at scale
- Approximation errors possible, but great for
 - Understanding data e.g., for evaluating data strategies
 - Regression testing esp. detecting data joinability
- Future work
 - Memory efficient but still CPU intensive
 - Applications in privacy enhancing techniques?

Thanks! Questions about **KHyperLogLog**?

References

1. **On Synopses for Distinct-Value Estimation Under Multiset Operations.**
— K. Beyer, P. J. Haas, B. Reinwald, Y. Sismanis, and R. Gemulla.
2. **HyperLogLog: the analysis of a near-optimal cardinality estimation algorithm.** — P. Flajolet, E. Fusy, O. Gandouet and F. Meunier.
3. **HyperLogLog in Practice: Algorithmic Engineering of a State of The Art Cardinality Estimation Algorithm.** — S. Heule, M. Nunkesser, A. Hall.

HLL++_{half-byte} VS HLL++

HLL++

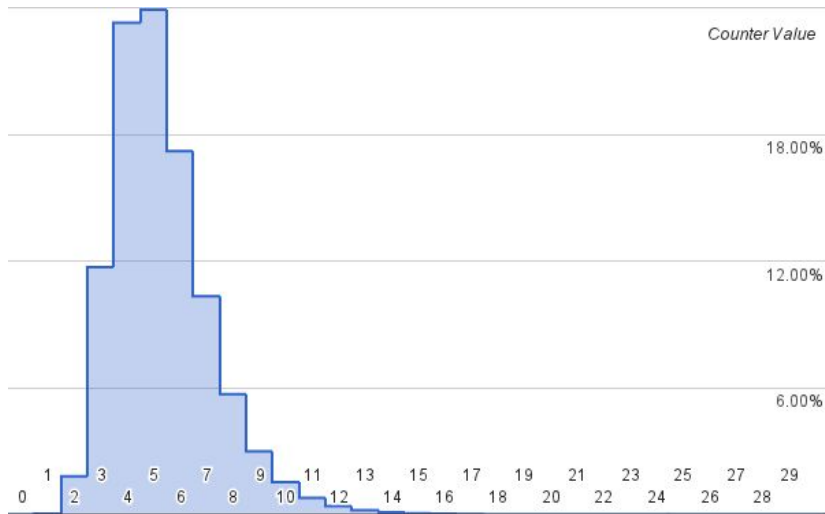
- Compressed sparse list
- Bias correction (between linear and HLL counting)

HLL++_{half-byte} -- readability | cpu efficiency | accuracy

- Sparse list
- Half-byte counters

Half-byte Counters in HLL++

- Observation: counters tend to clump together



<u>bucket</u>	<u>max trailing zero</u>
0	4
1	8
2	12
3	7
4	5
5	5
6	7
7	6



<u>bucket</u>	<u>max trailing zero</u>
0	0
1	4
2	8
3	3
4	1
5	1
6	3
7	2

offset 4