



中國人民大學
RENMIN UNIVERSITY OF CHINA

香港科技大學
THE HONG KONG
UNIVERSITY OF SCIENCE
AND TECHNOLOGY

Persistent Data Sketching

Zhewei Wei¹, Ge Luo², Ke Yi², Xiaoyong Du², Ji-Rong Wen¹

¹ DEKE and School of Information, Renmin University of China

² Department of Computer Science and Engineering, The Hong Kong

University of Science and Technology

Contact: zhewei@ruc.edu.cn

Motivation and Problem Statement

Persistent Database:

- Microsoft Immortal DB, SNAP, Ganymed, Skippy and LIVE
- Allow queries on **past version** of the database
- Linear space: store all updates

(Ephemeral) Sketches:

- Count-Min sketch, AMS sketch
- Allow queries on current data
- Answer queries approximately with **sub-linear** space

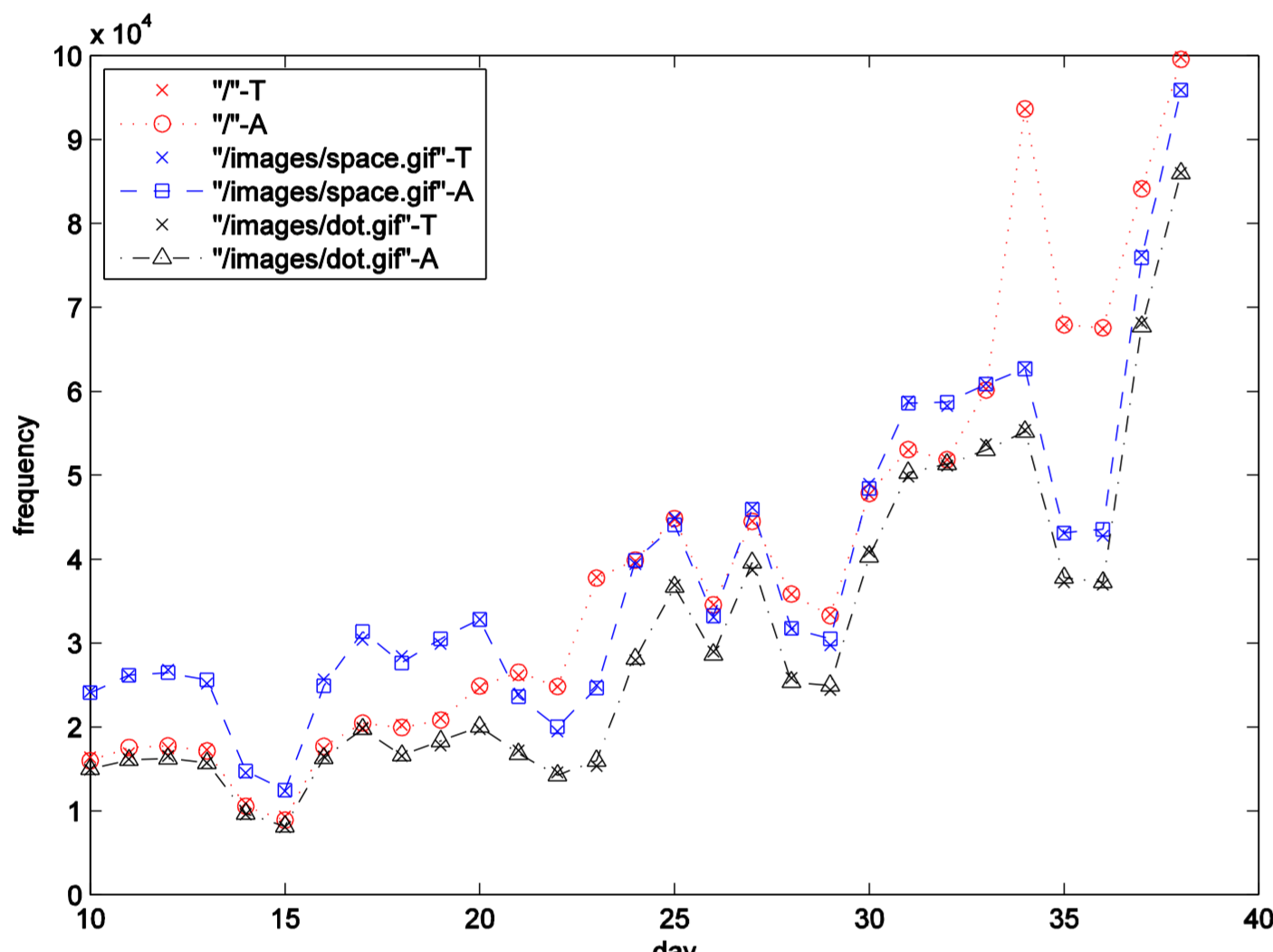
Persistent Sketches:

- Allow queries on historical data
- Sub-linear space

Historical query (0, t)
Historical window query (s, t)

An Illustrating Example

URL	actual count	estimation
/	1138896	1138970
/images/space.gif	1117634	1120050
/images/dot.gif	880322	880765
/images/hm_nbg.jpg	818126	818586
/images/home_intro.anim.gif	799697	800323



Tracking top-k items

● 1998 World Cup web site access log:

(timestamp, IP address of the request, requested URL; size of the response, method).

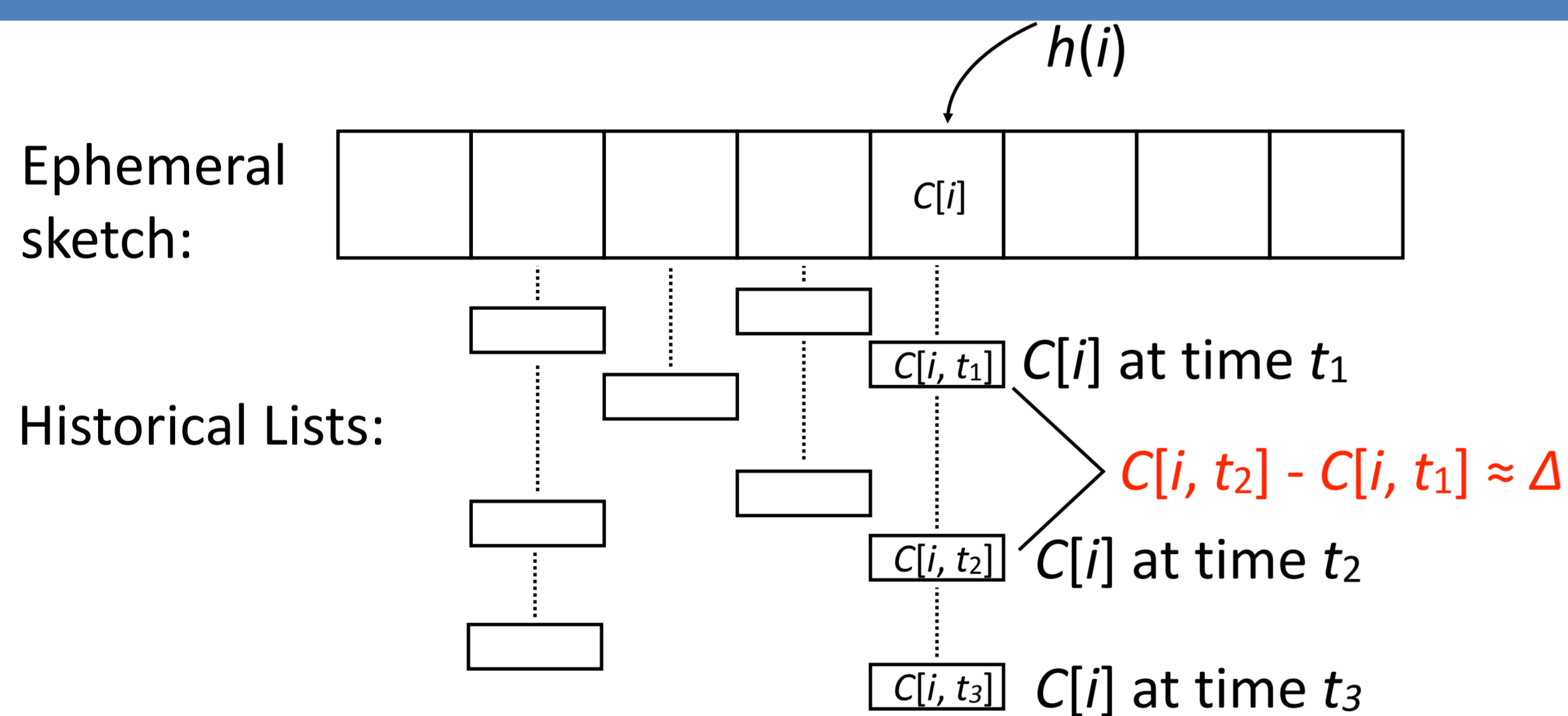
● Ephemeral sketch

Track top-k URLs at the end of the stream

● Persistent sketch

Track how top-k frequencies change over time

Baseline Solution



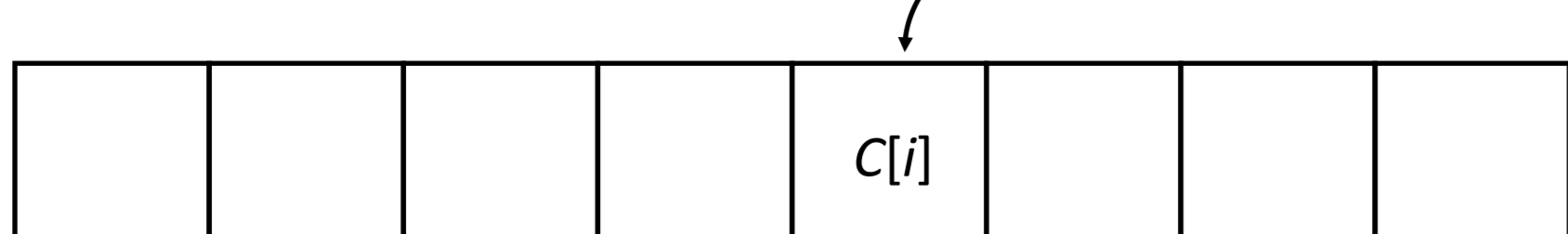
- Query for time t : find counter with timestamp closest to t
- Error: ϵn (ephemeral error) + Δ (persistent error)
- Space: proportional to $(1/\epsilon + m/\Delta)$

PLA-based Count-Min Persistent Sketch

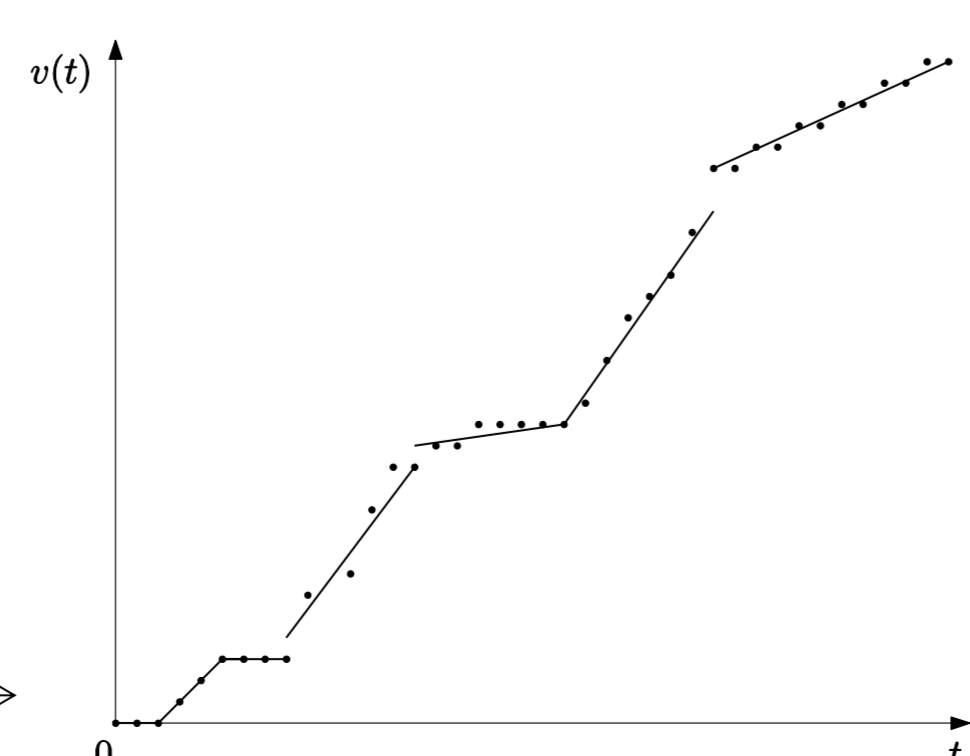
- Each counter is a discrete function according to timestamp
- (Count, time) points can be approximated with a line segment
- Approximate ratio: at most Δ

Piece-wise linear approximation (PLA):

Ephemeral sketch:



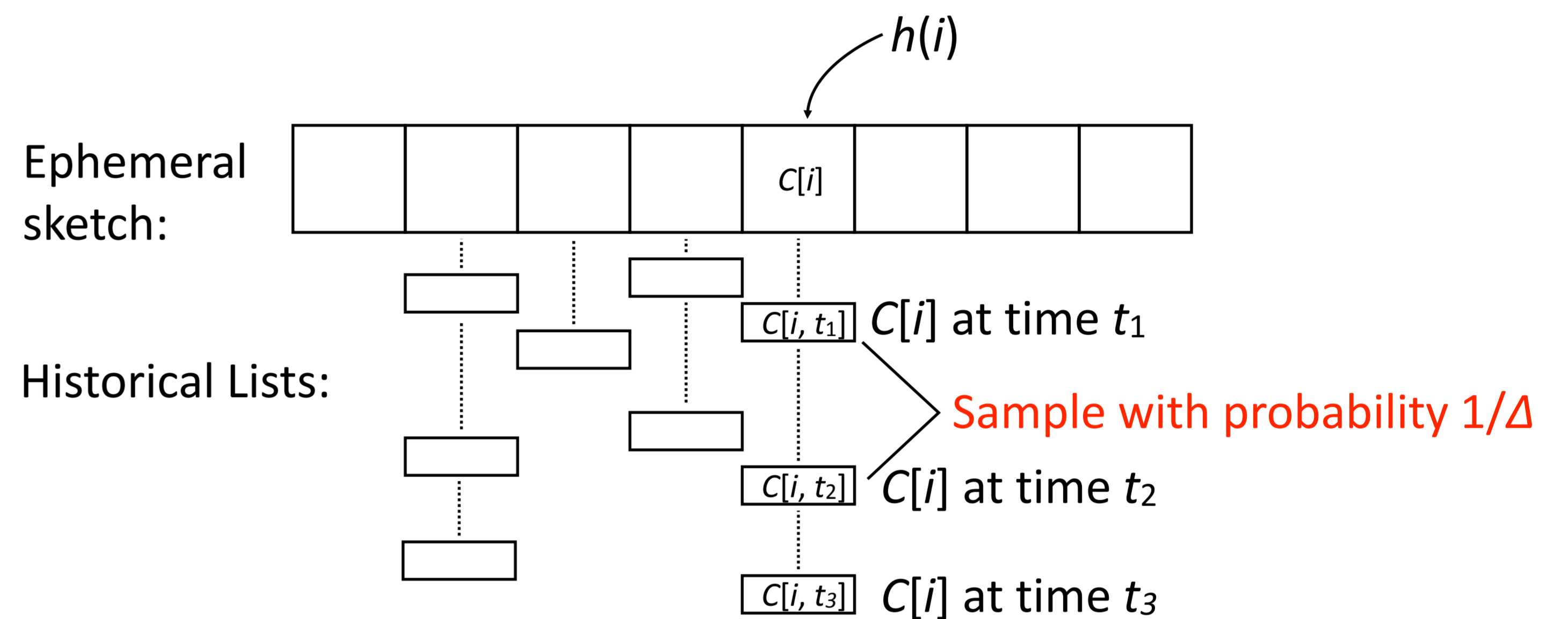
PLA generators:



- Query for time t : return the value of the PLA function at t
- Error: ϵn (ephemeral error) + Δ (persistent error)
- Space: proportional to $(1/\epsilon + m/\Delta^2)$ in **random stream model**
 - Random walk along the piece-wise linear function
 - Expected time to escape the linear function by Δ : Δ^2

Sampling-based AMS Persistent Sketch

- Estimating (self) join size: $\sum_{i \in [n]} (C[i] + \text{error of } \Delta)^2$
- **Bias** will amplify error significantly



- Given a query time t
 - Set $C[i]$ to be $C[i, t_k] + \Delta - 1$ if $C[i, t_k]$ is the last sampled counts precedes t , and 0 if there is no sample preceding t .
 - **Unbiased estimator** for $C[i]$ at time t
- Error: ϵF_2 (ephemeral error) + $(\Delta/\epsilon)^2$ (persistent error)
- Space: proportional to $(1/\epsilon + m/\Delta)$

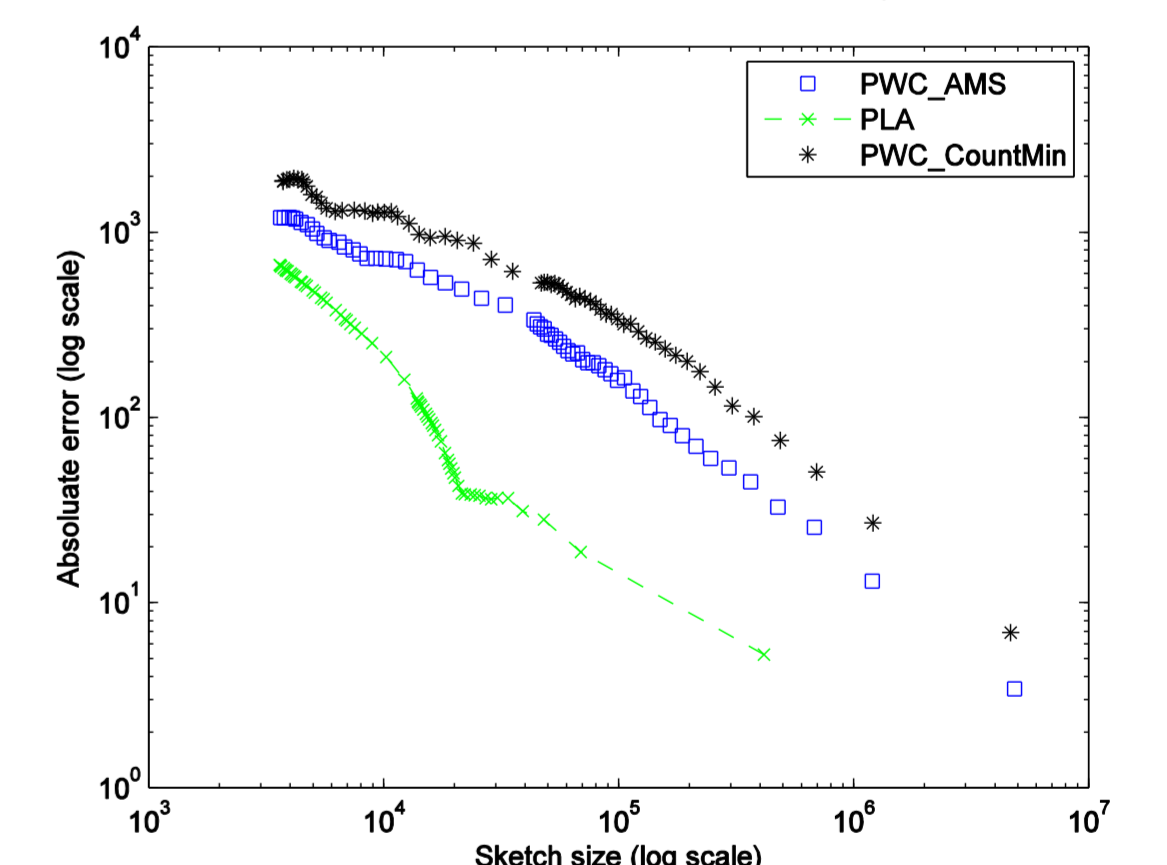
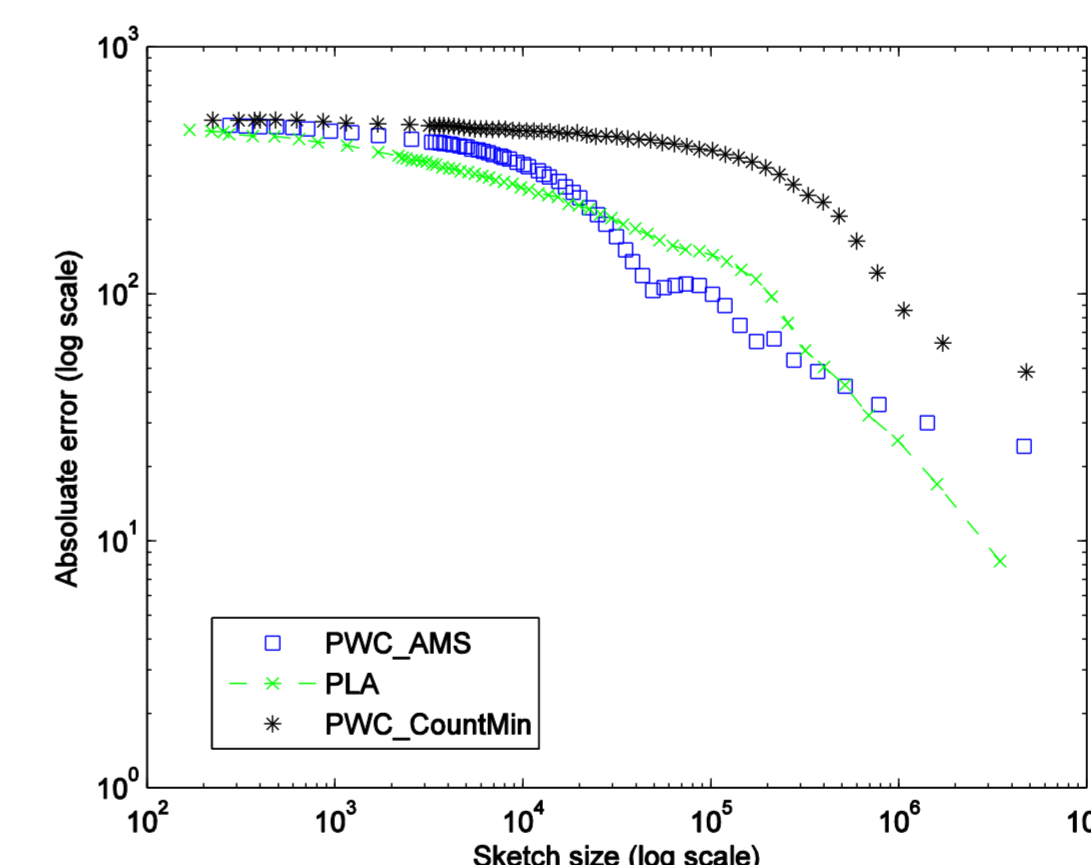
Experimental Results

- 7,000,000 requests from the 1998 World Cup web site access log
- Built sketches on two attributes

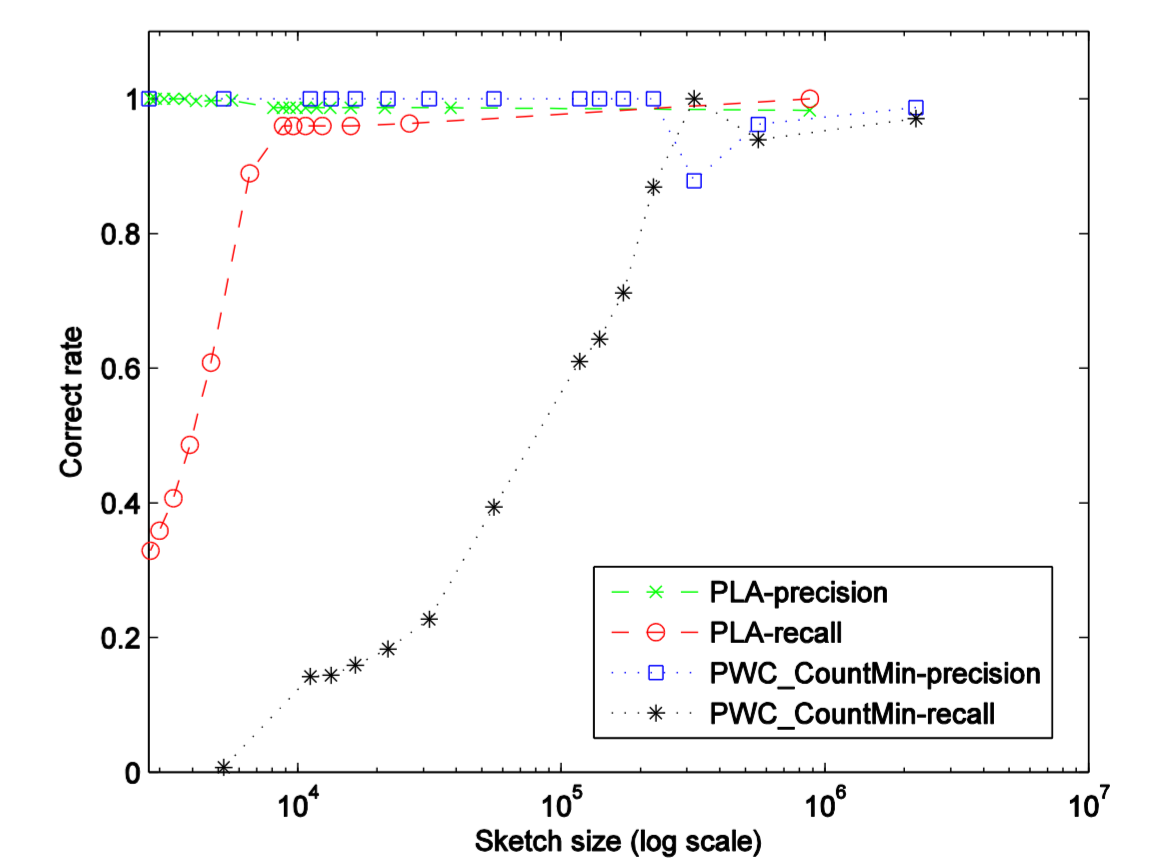
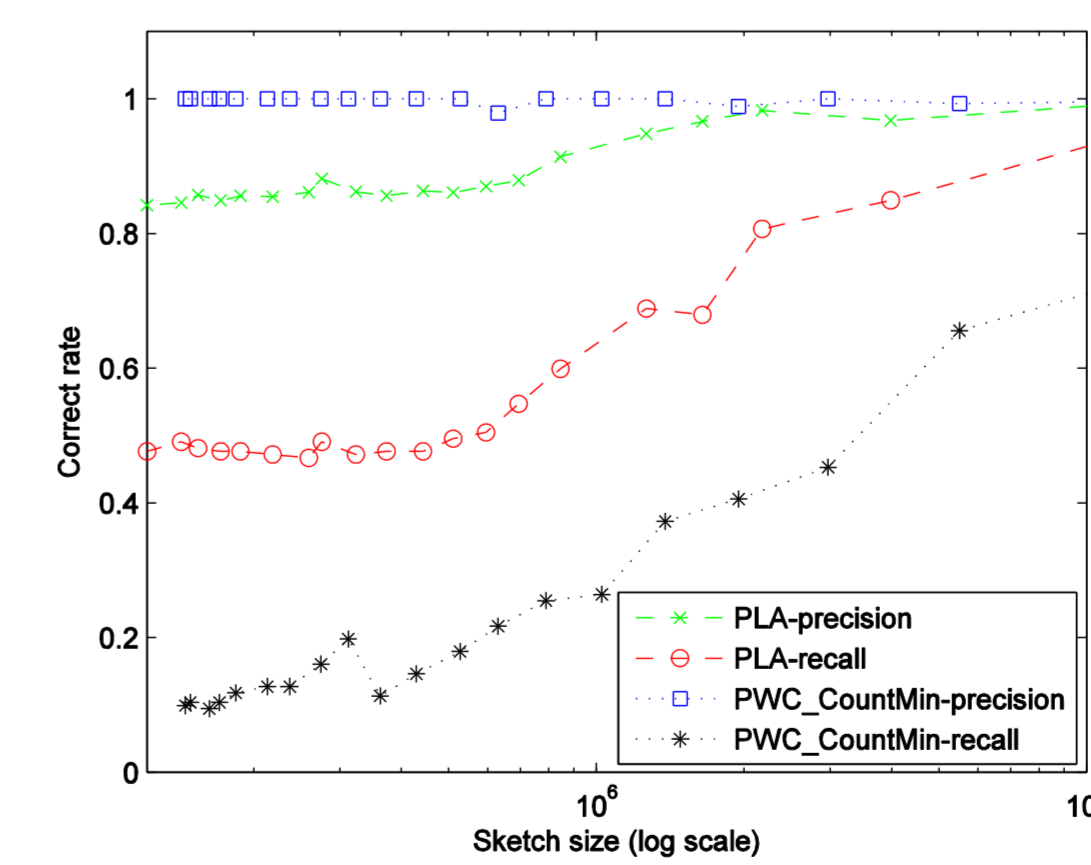
Attributes: Requested URL

IP address of the request

Point query



Heavy hitters query



Self join size query

