

# Fast Memory-Efficient Anomaly Detection in Streaming Heterogenous Graphs



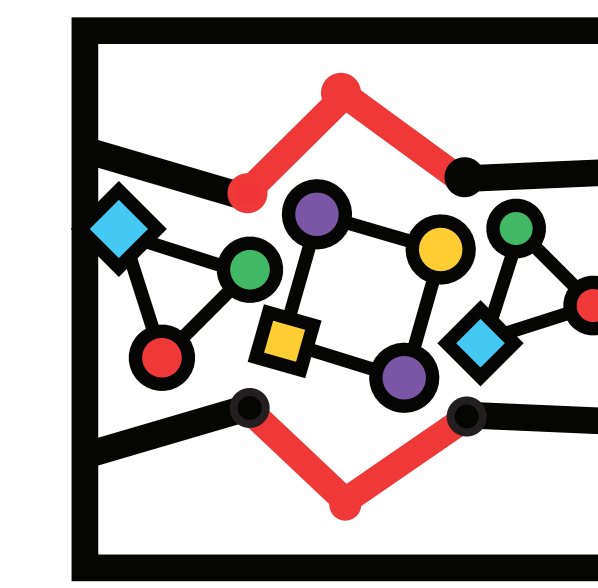
Emaad Manzoor\*



Sadegh M. Milajerdi

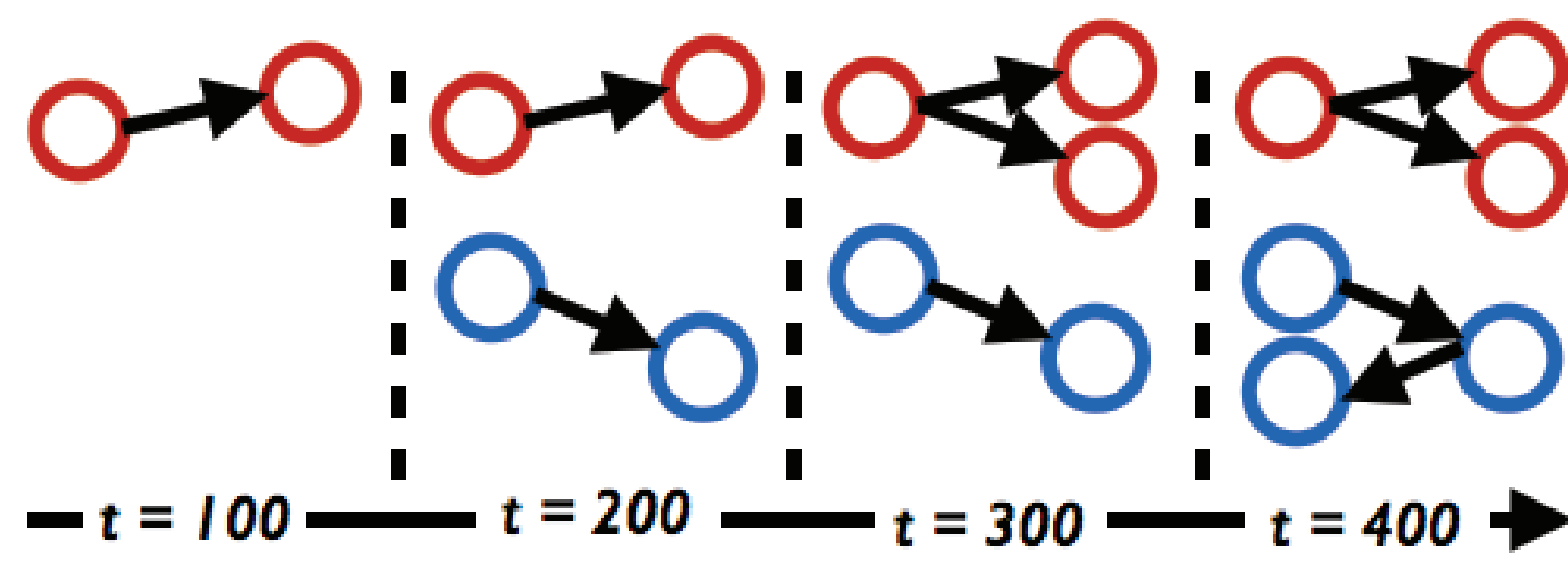


Leman Akoglu



Code and data at  
[bit.ly/streamspot](http://bit.ly/streamspot)  
[emanzoor@cs.stonybrook.edu](mailto:emanzoor@cs.stonybrook.edu)

**StreamSpot** tracks anomalous heterogenous graph objects as they evolve from a stream of typed edges.



## Challenges

**Bounded-space** graph representation

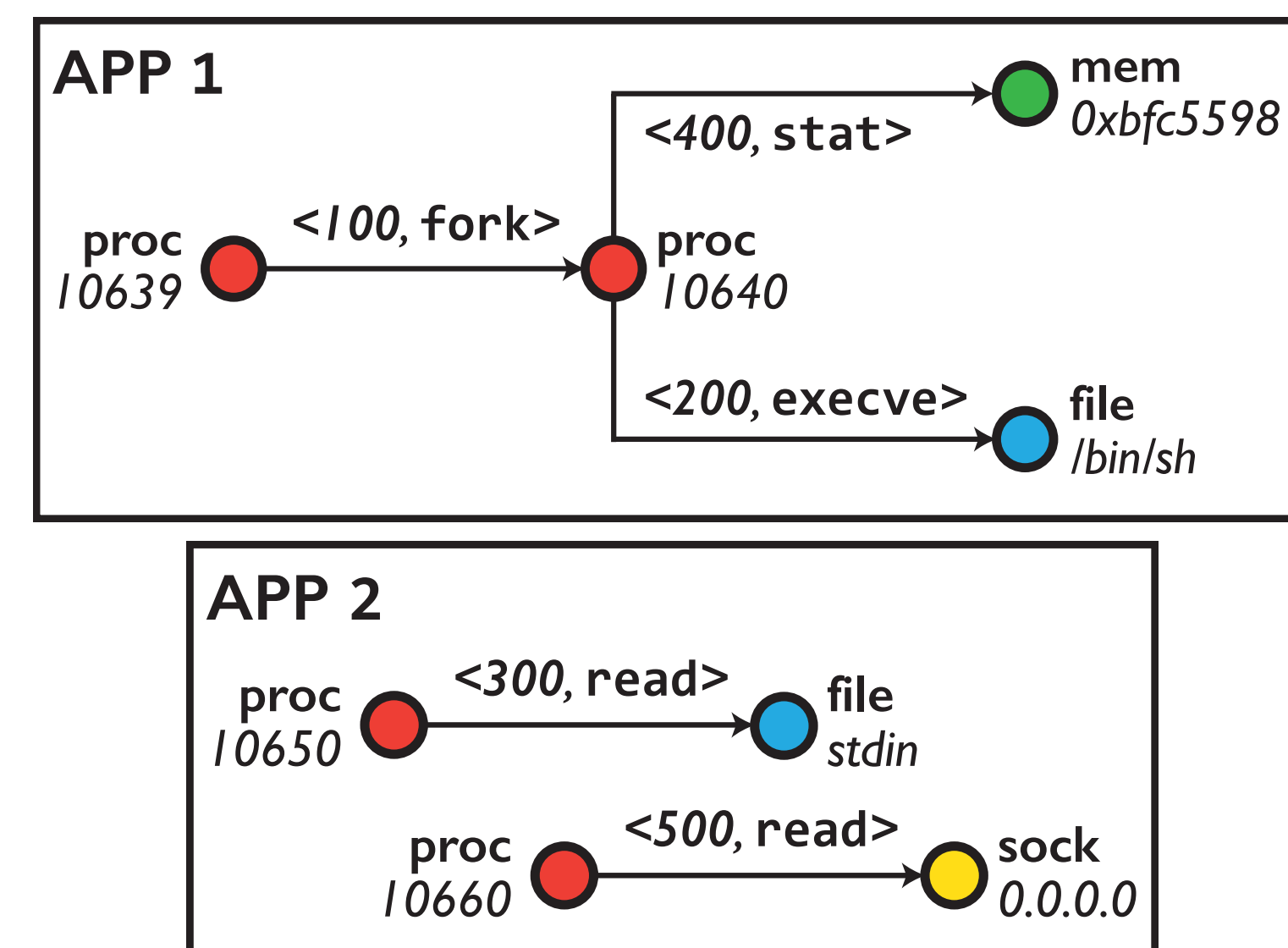
**Streaming** graph comparison

**Fast and constant-time** per-edge

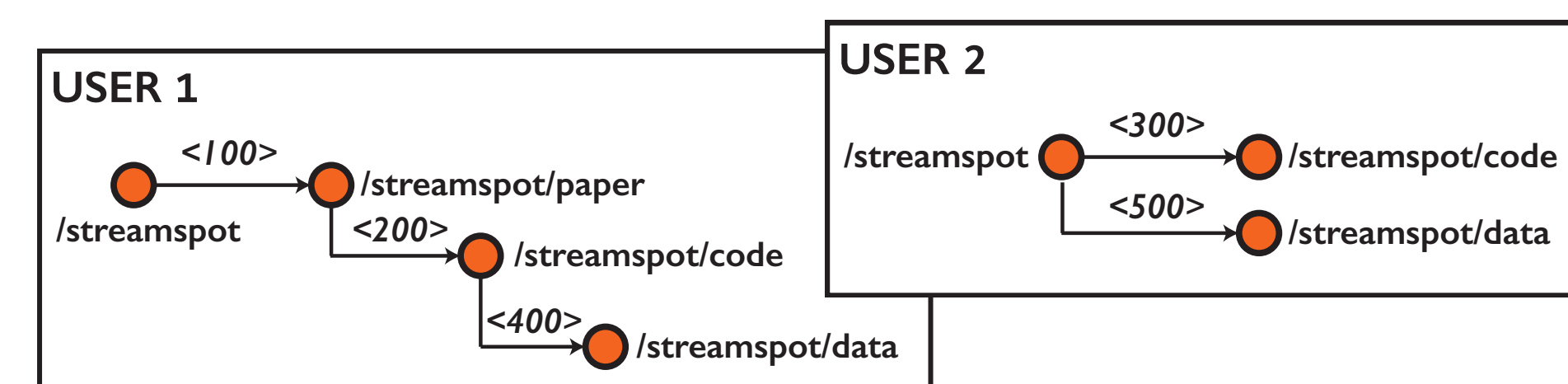
**Accurate** anomaly scoring

## Applications

Detecting **malicious software** from a stream of **system call logs**



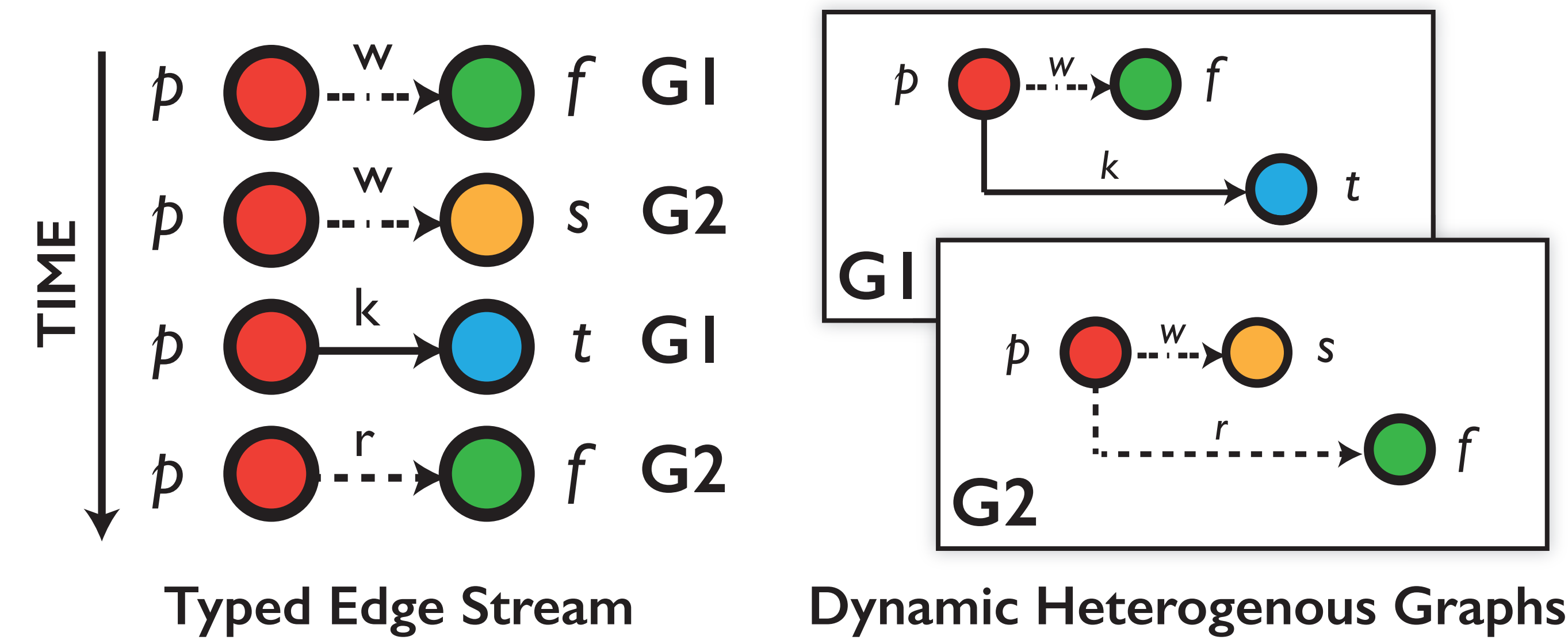
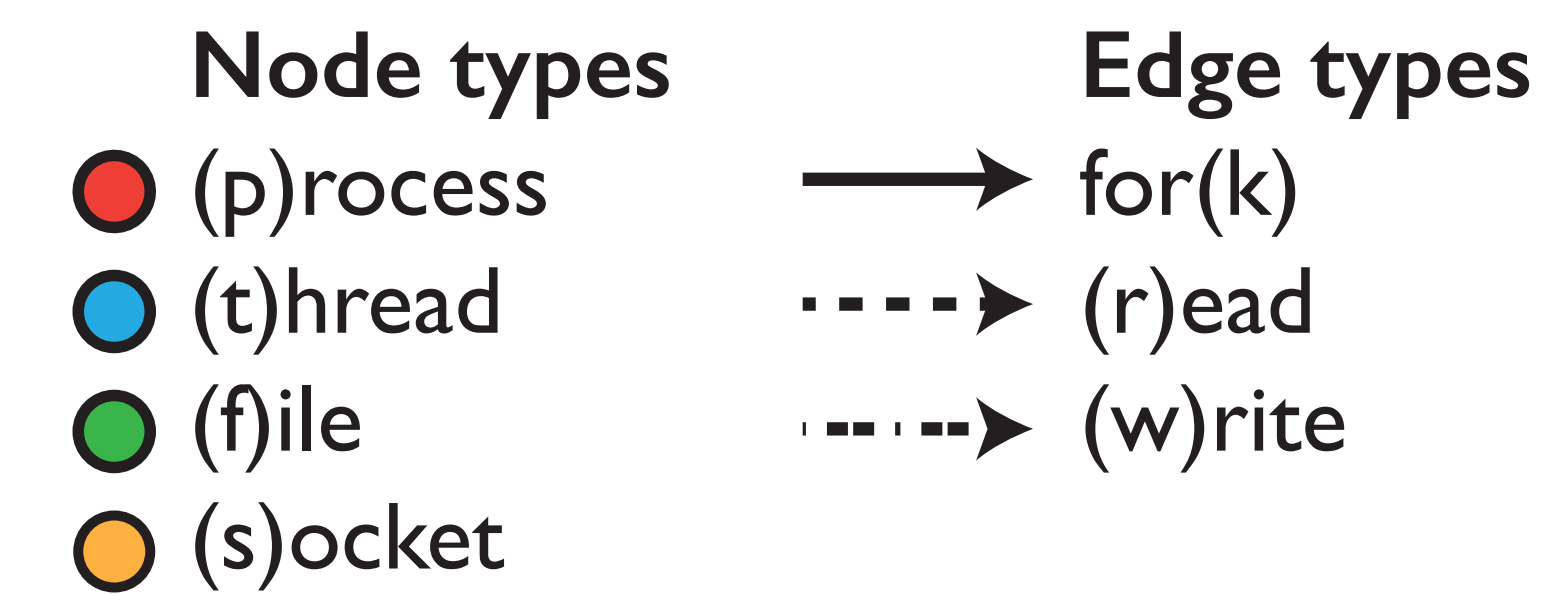
Detecting **malicious users** from a stream of **filesystem navigation** traces



Thanks to support from



## Typed event stream as dynamic heterogenous graphs



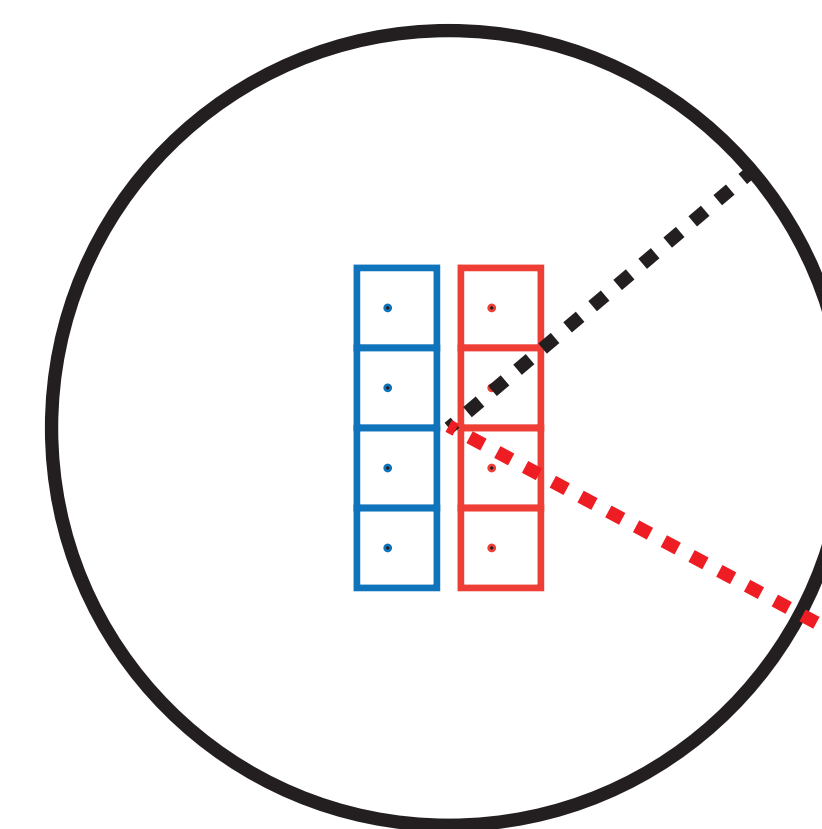
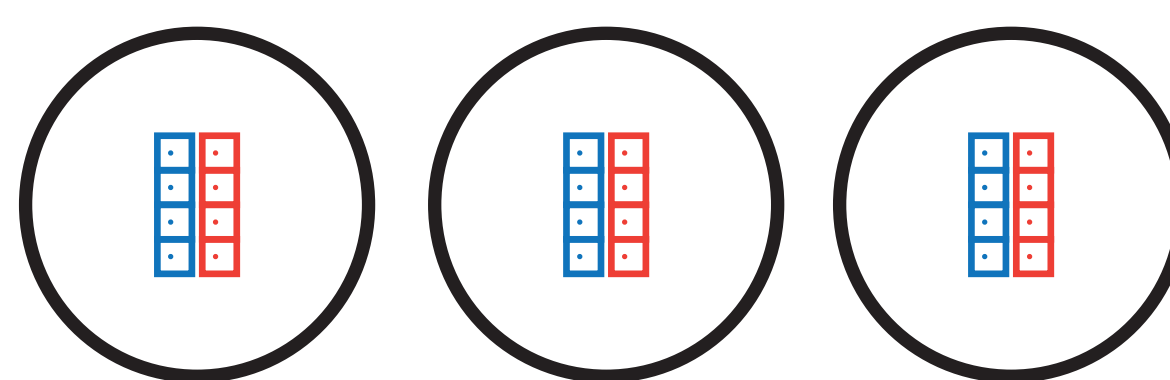
## Clustering-based anomaly detection

K Bootstrap Clusters

From benign training graphs

Cluster Centroid

Average graph of the cluster

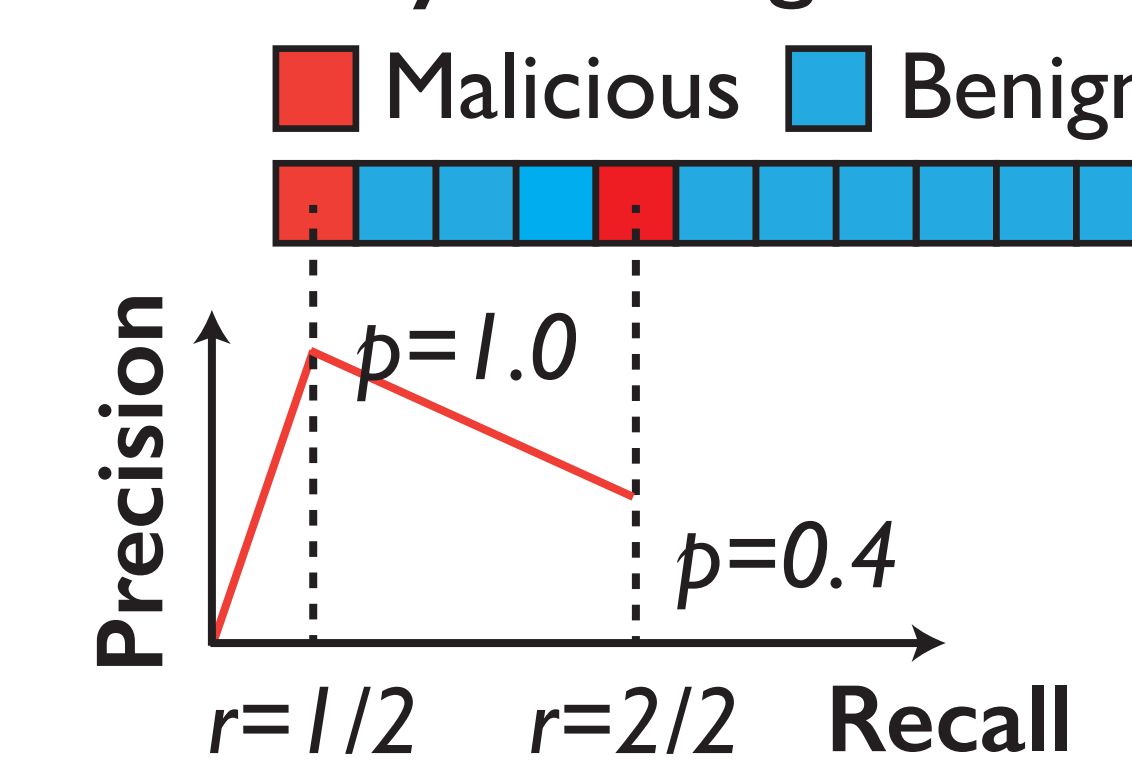


Cluster Radius

3 standard-deviations from the mean graph to centroid distance

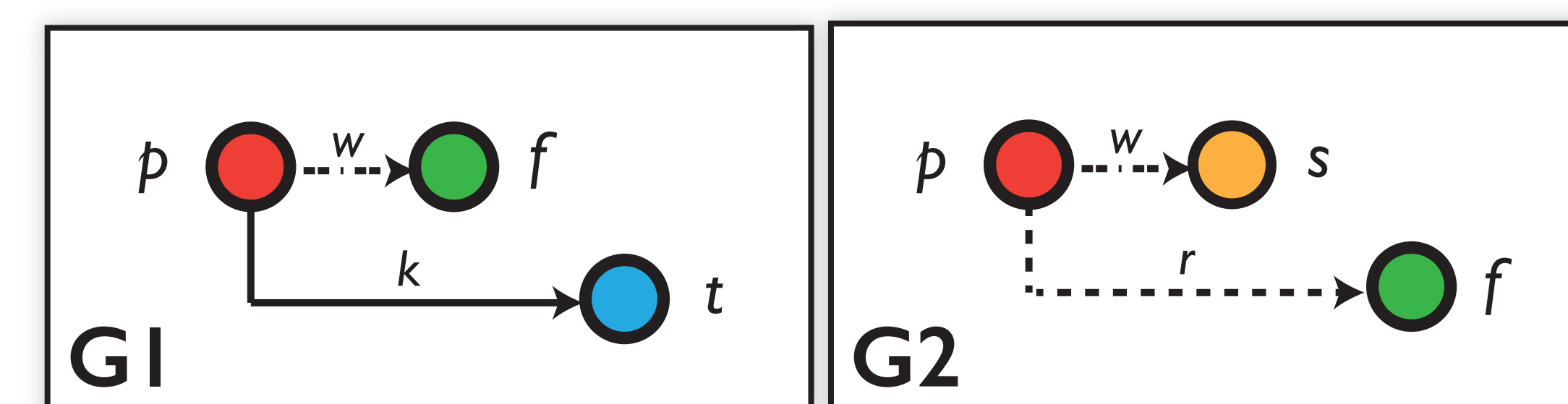
Anomaly Score  
Distance to nearest cluster centroid

Anomaly Ranking



## Dynamic heterogenous graphs as vectors by shingling

OkBFT Ordered k-hop Breadth-first Traversal



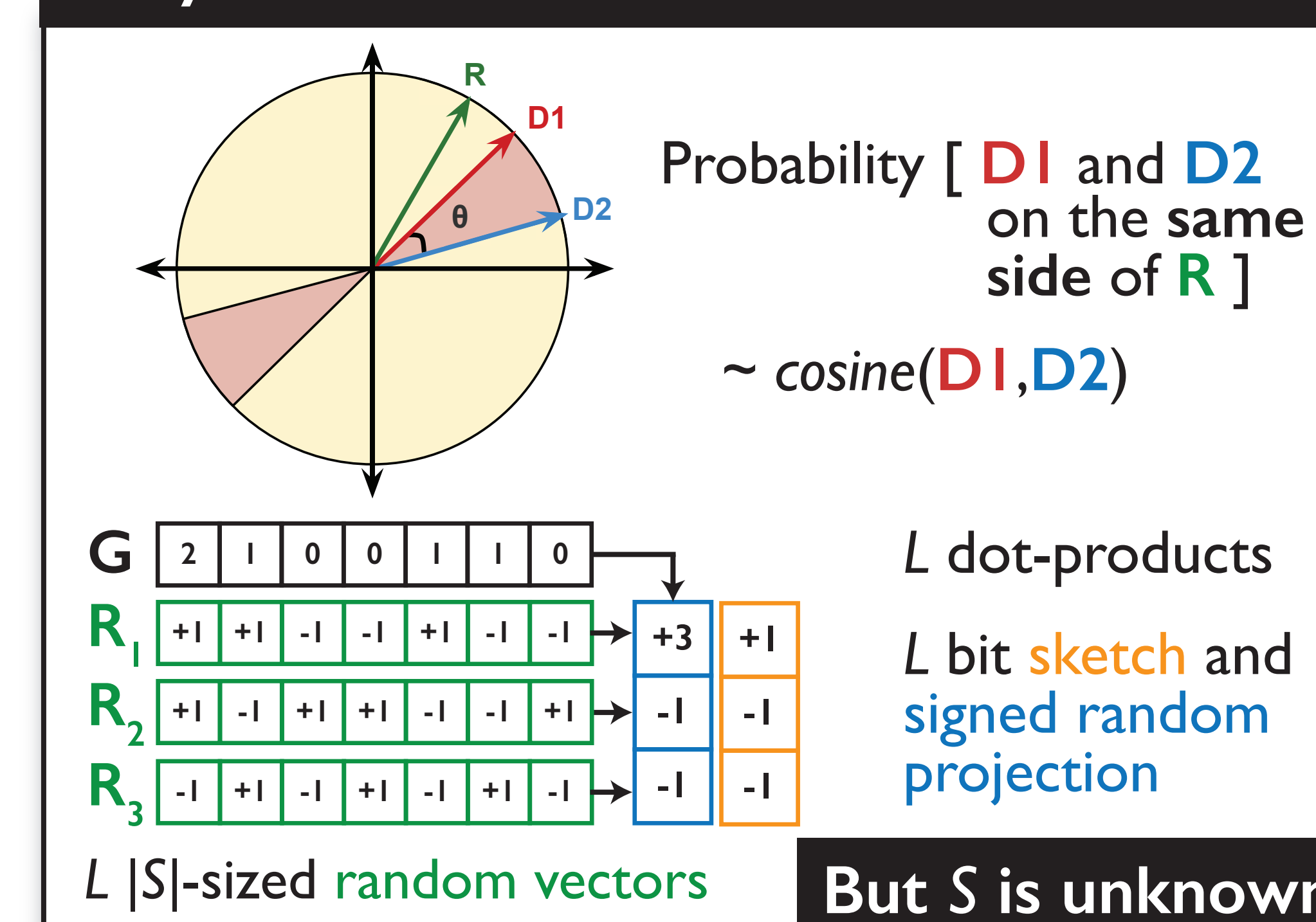
Shingle Vectors

G1	1	0
G2	0	1
	pwfkt	pwsrf

Shingle Universe S  
**Large!**  
**Unknown!**

## Streaming pairwise graph similarity with STREAMHASH

Why does SIMHASH fail?



Why cache when you can **hash**?

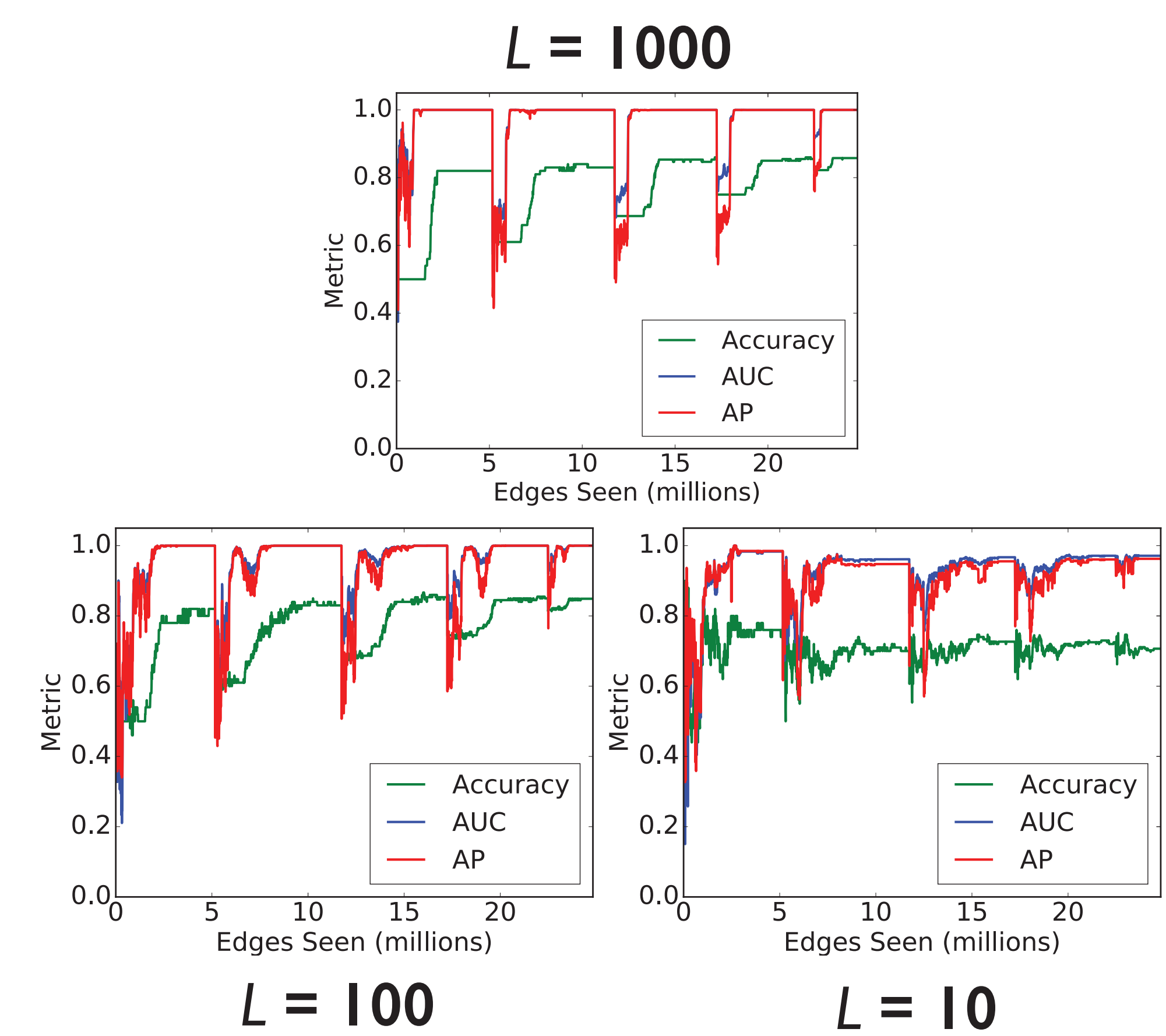
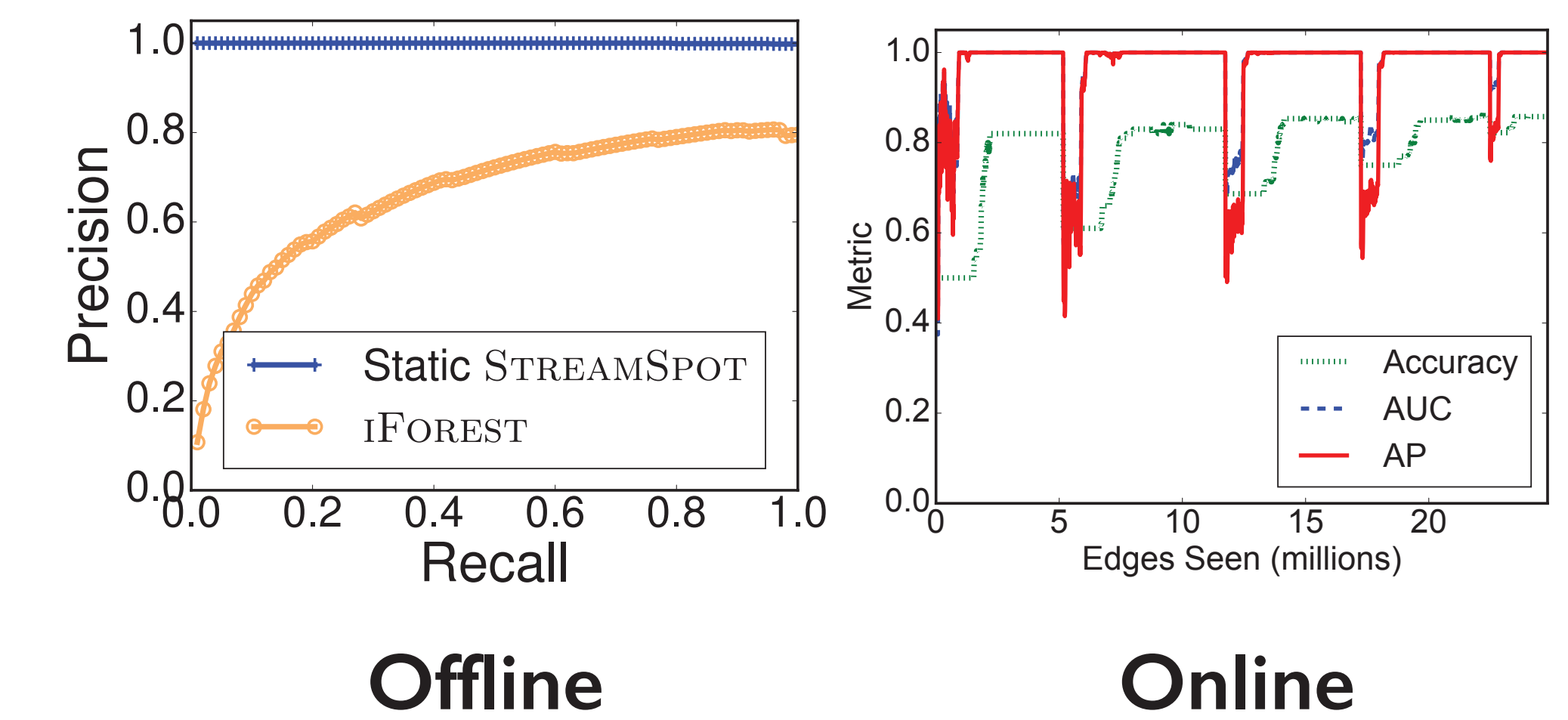
$h_i(s_1) \ h_i(s_2) \ h_i(s_{|S|})$   
 $h_i(s): s \rightarrow \{+1, -1\}$   
Store  $h_1 \dots h_L$  in constant space

**Incremental** sketch construction

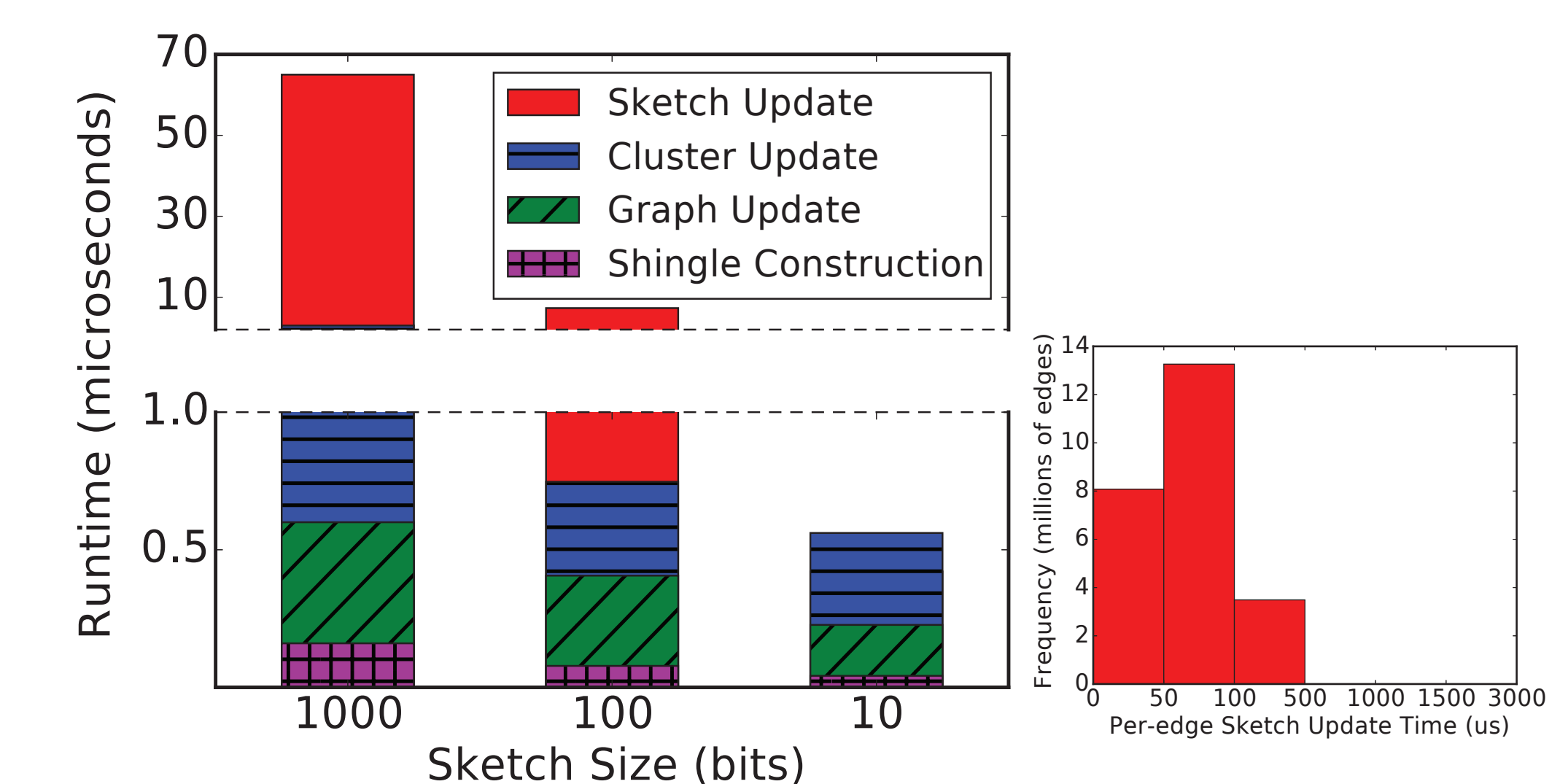
Add/remove shingle **s** in constant time

Old projection vector  $\pm$  StreamHash update vector = New projection vector

## Accuracy

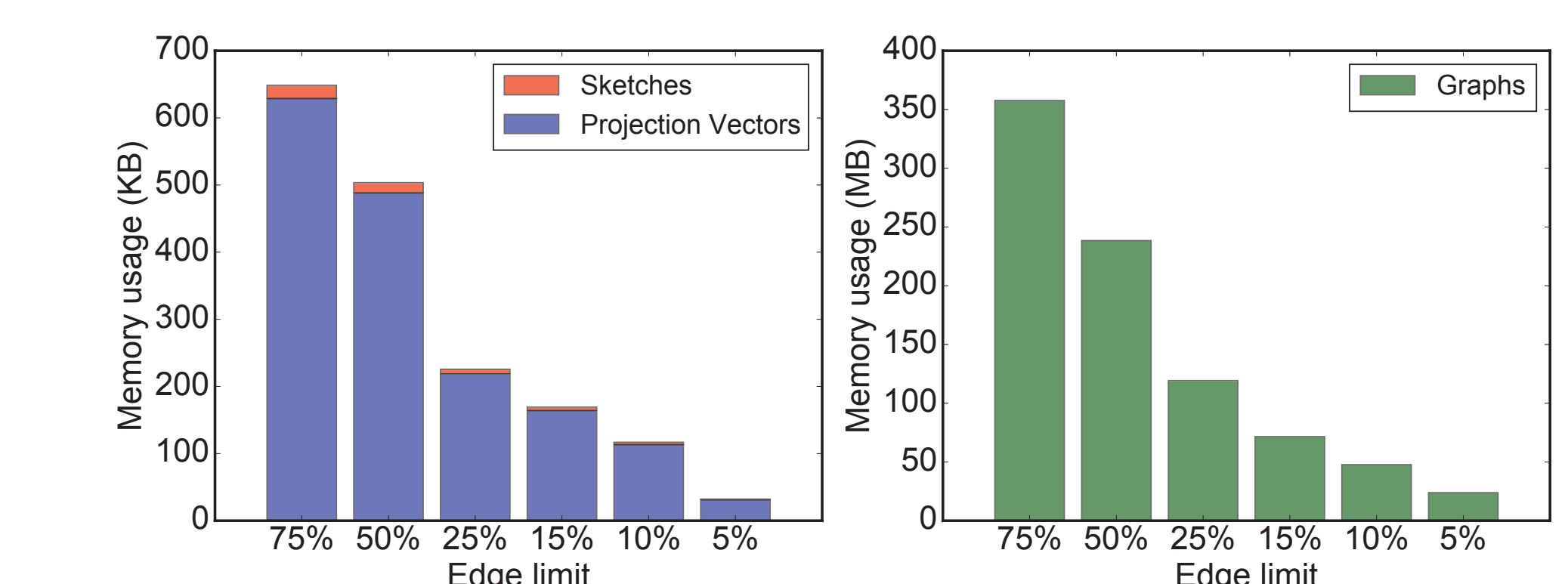


## Time



$L = 100 \rightarrow 100,000$  edges/s

## Space



$\sim 240\text{MB}$  with  $N = 12.5\text{M}$ ,  $L = 1000$  bits