



# EIGENPULSE: DETECTING SURGES IN LARGE STREAMING GRAPH WITH ROW AUGMENTATION

Jiabao Zhang<sup>+</sup>, Shenghua Liu<sup>+</sup>, Wenjian Yu<sup>\*</sup>, Wenjie Feng<sup>+</sup> and Xueqi Cheng<sup>+</sup>

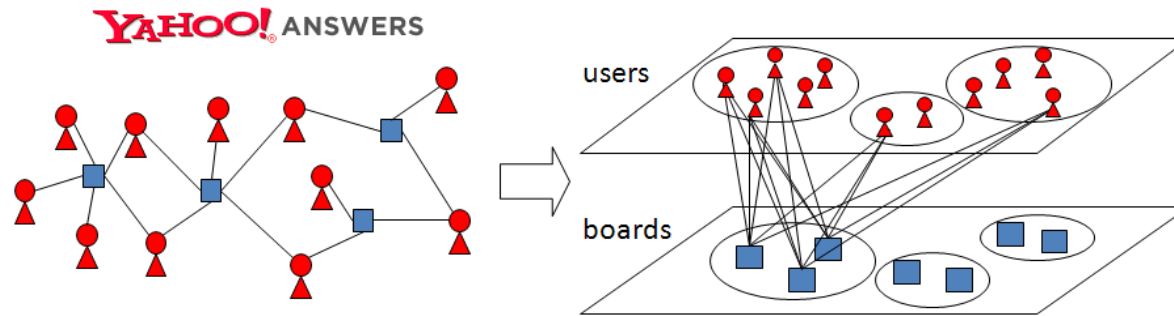
2020/4/20

<sup>\*</sup>Tsinghua University


<sup>+</sup>Institute of Computing Technology ICT, CAS

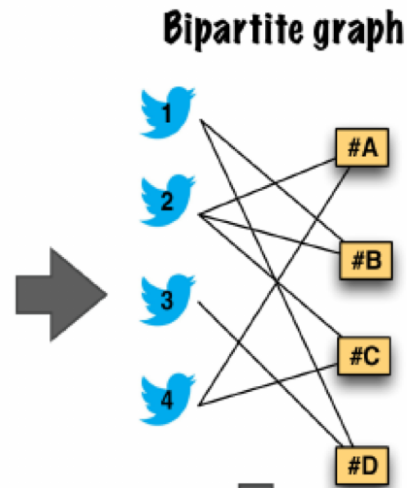
# Graphs are everywhere.

## ■ Yahoo answers

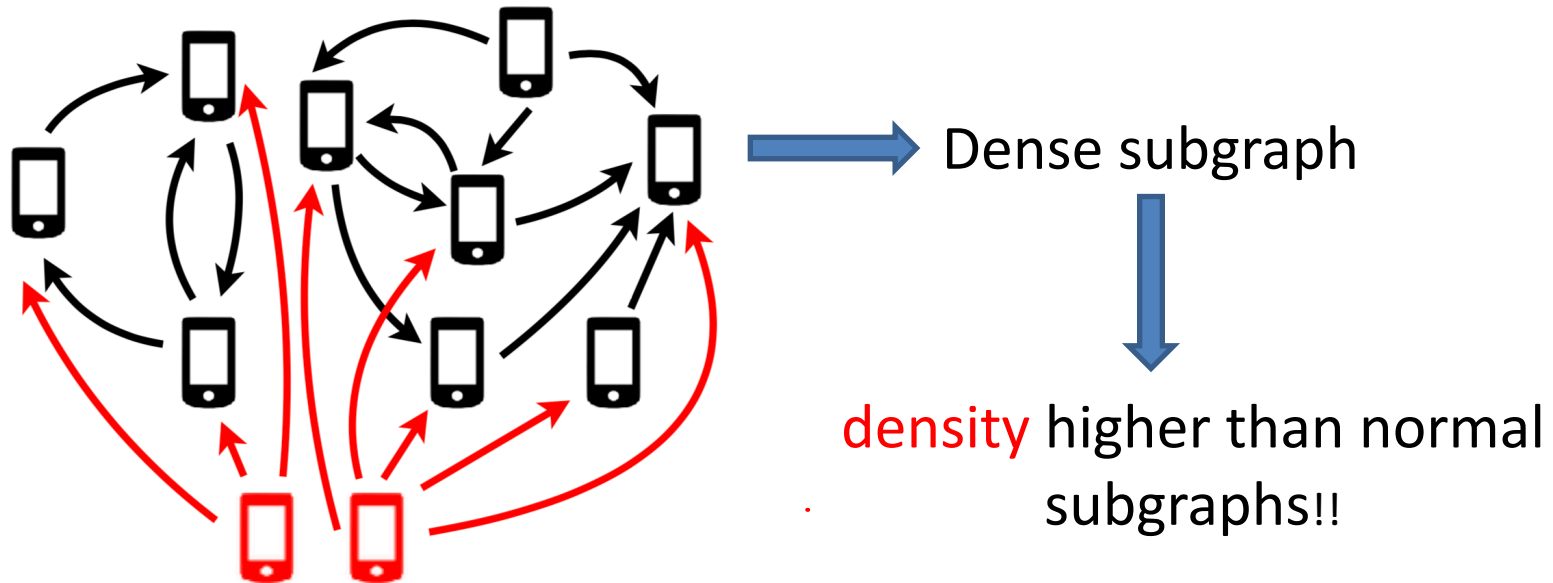


## ■ Twitter

	text1 #B, #F
	text2 #A, #B, #H
	text3 #D #G
	...



# Density usually indicates unusual events










Eswaran, D., Faloutsos, C., Guha, S., Mishra, N. Spotlight: Detecting anomalies in streaming graphs. In: SIGKDD. pp. 1378–1386. ACM (2018)

# Adjacency matrix for a graph

$G(U, V, E)$

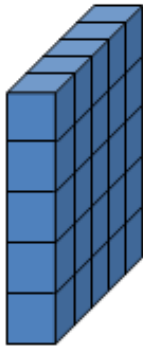


Adjacency Matrix

0			0		
		1			1
		1			0
				0	
1			1		
					
					
					
					

# Streaming graph

- graphs usually expand with time.



time

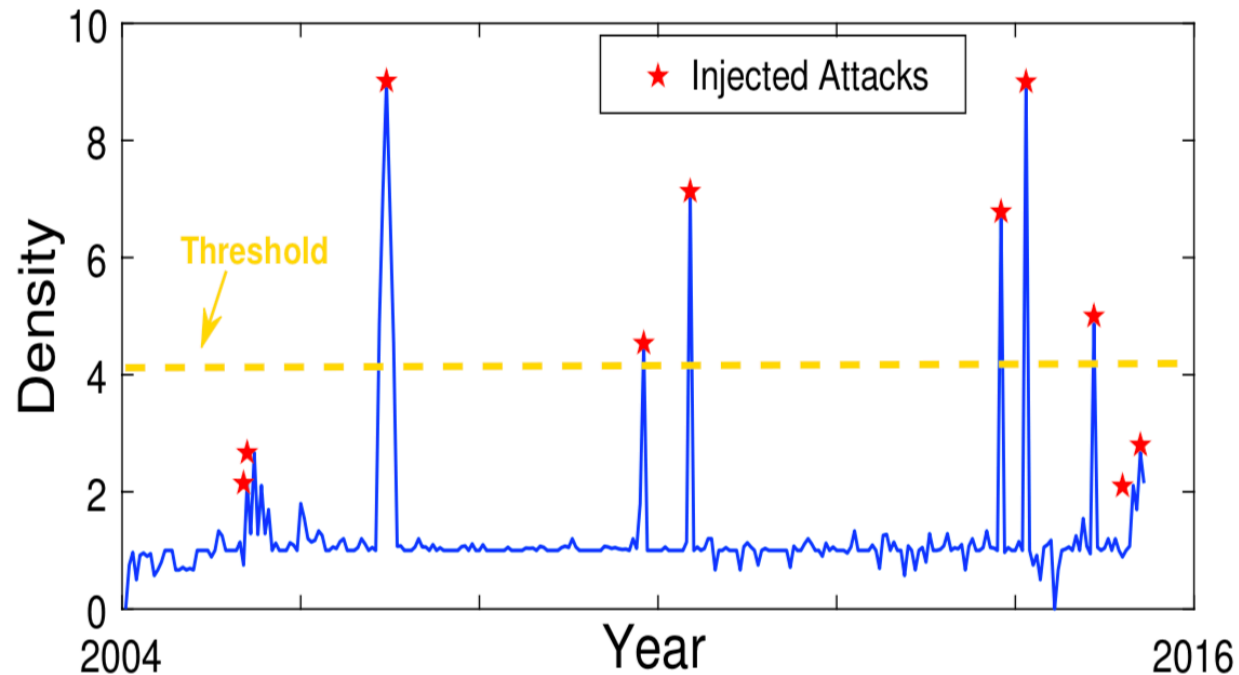
**How do we detect such **anomalies** in  
streaming graphs?**

How do we even characterize these **anomalies**?

.

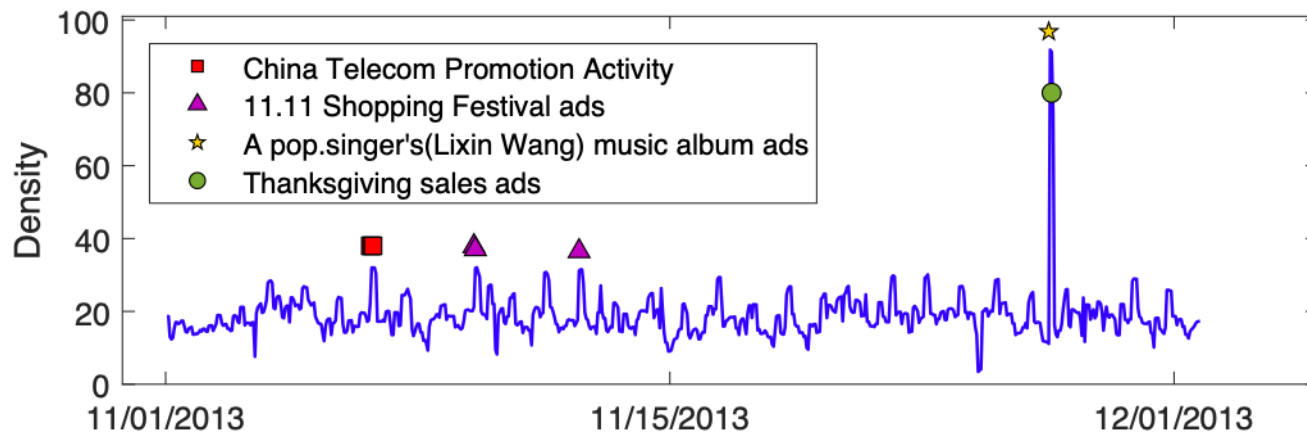
# EigenPulse: detect injection accurately

EigenPulse: “*graph heart-beat*”



# EigenPulse: detect anomalous surges on real data

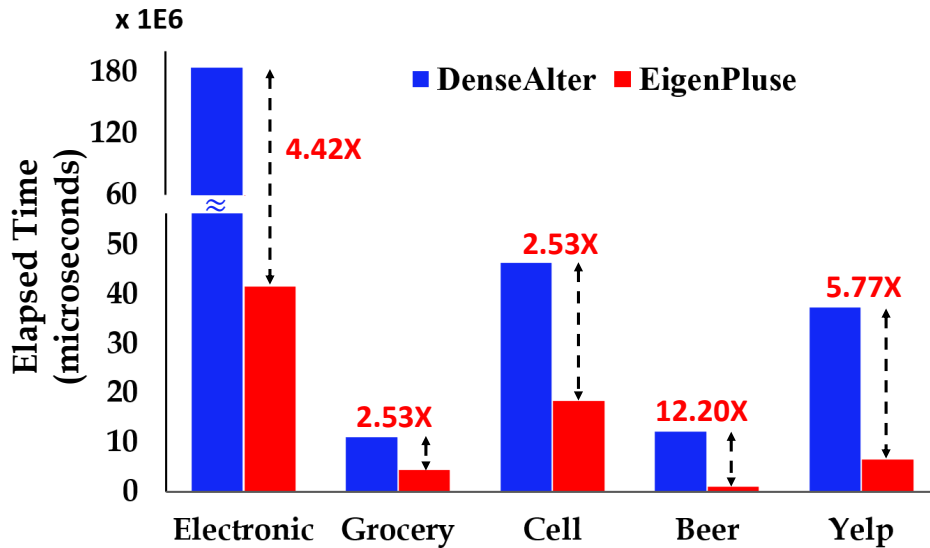
- Microblog: Sina Weibo, Nov. 2013
  - node size: 2.74M x 8.08M, # of edges: 50.06M



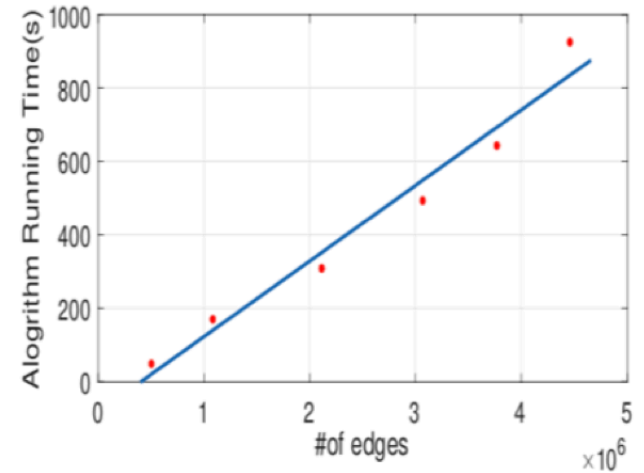
EigenPulse: “graph heart-beat”



# EigenPulse: run fast and near-linearly



EigenPulse outperforms DenseAlert by more than 2.53x.



run near-linearly in # of edges

# Outline

- Problem
- Related works
- Our model
- Our algorithm
- Experiments
- Conclusion

# Problem

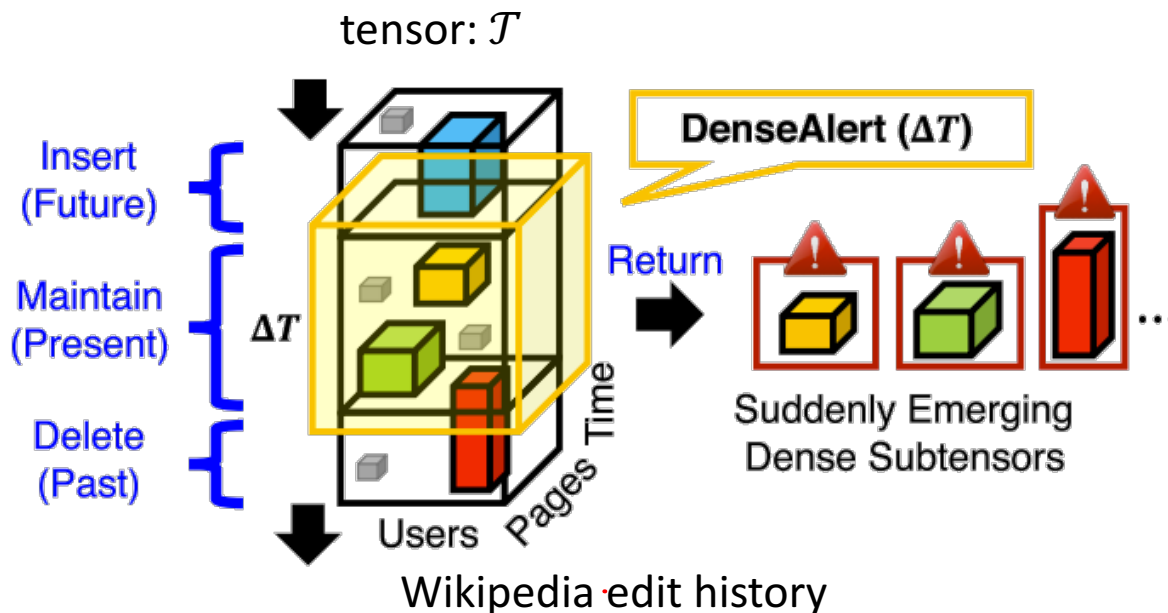
## ■ Given:

- a stream of triplets (*user, object, timestamp*).

## ■ Find:

- at *each time step*, a group of users and objects who have *densest* edges
- *detect* suspicious surges of density

# Related work: DenseAlert



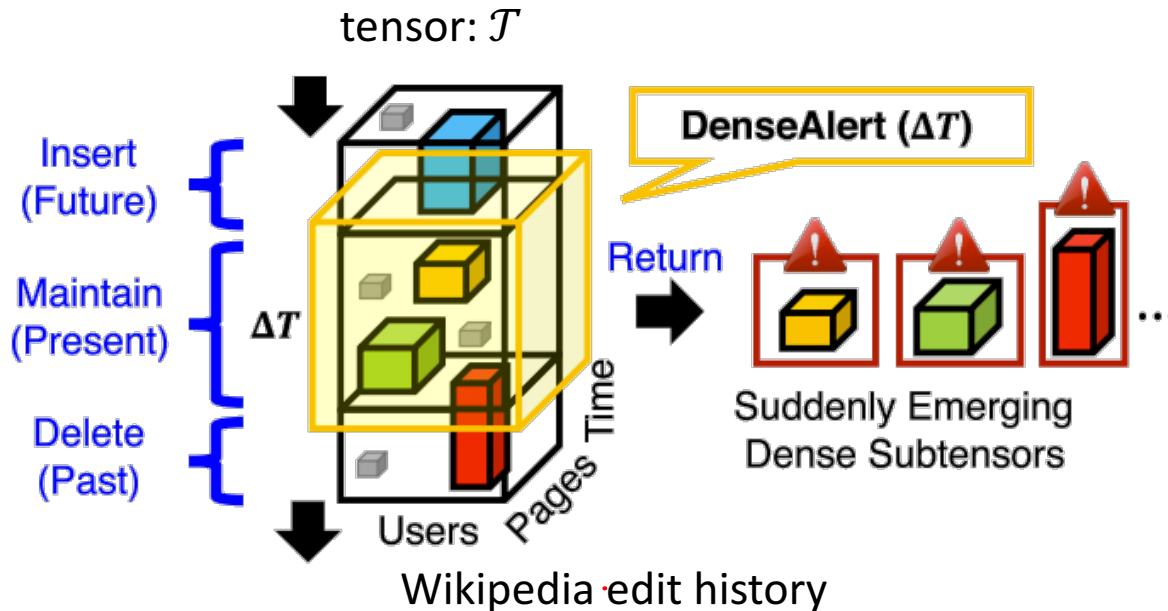
**Given:** a **stream** of changes, e.g. adding/removing edges in tensor  $\mathcal{T}$ .

- **Maintain:** a subtensor  $\mathcal{T}(S)$ , where  $S$  is set of slice indices.

- **to maximize:** density  $\rho(\mathcal{T}(S))$

$$\rho(\mathcal{T}(S)) = \frac{\text{sum}(\mathcal{T}(S))}{|S|}$$

# Related work: DenseAlert

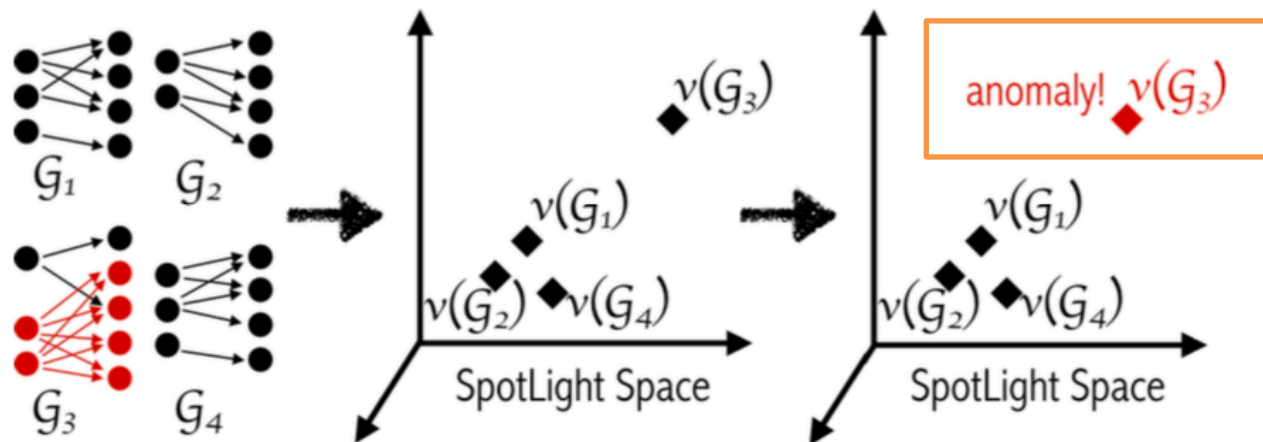


- Need an updating operation for *every single adding or removing* edge.
- Running a bit *slow*.

Shin, K., Hooi, B., Kim, J., Faloutsos, C.: Densealert: Incremental dense-subtensor detection in tensor streams. In: KDD. ACM (2017)

# Related work: SpotLight

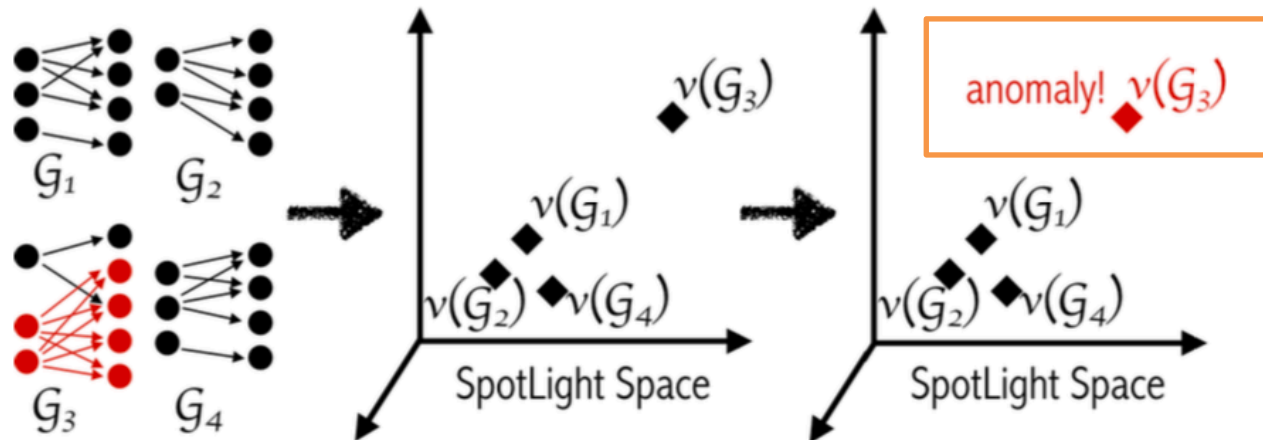
- Given a sequence of *snapshots* of dynamic graphs
  - extract a K-dimensional *sketch*  $v(G)$  for every snapshot.
  - *detect* the anomalous snapshot in sketch space



Eswaran, D., Faloutsos, C., Guha, S., Mishra, N. Spotlight: Detecting anomalies in streaming graphs. In: SIGKDD. pp. 1378–1386. ACM (2018)

# Related work: SpotLight

- *Only* spot anomalous snapshots at time windows
- *Not* detect the exactly suspicious groups (subgraphs)



Eswaran, D., Faloutsos, C., Guha, S., Mishra, N. Spotlight: Detecting anomalies in streaming graphs. In: SIGKDD. pp. 1378–1386. ACM (2018)

# Related works: summary of comparisons

	Fraudar	HoloScope D-Cube M-Zoom	DenseAlert	SpotLight	EigenPulse
temporal information		✓	✓	✓	✓
streaming graphs			✓	✓	✓
suspicious groups	✓	✓	✓		✓
theoretical analysis	✓	✓	✓	✓	✓
scalability	✓	✓	✓	✓	✓



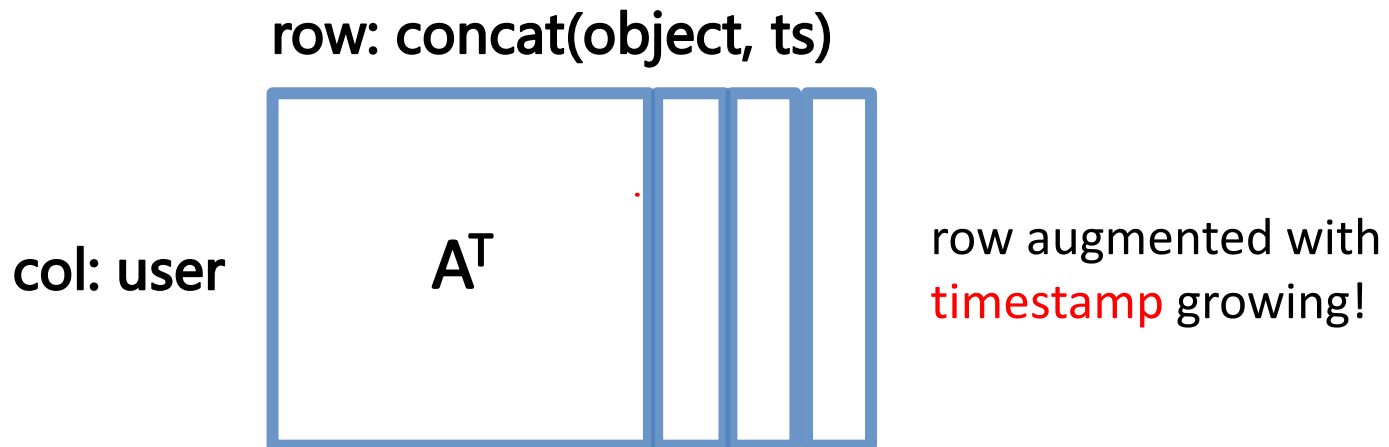
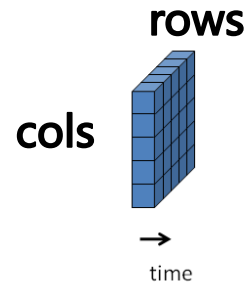
# Outline

- Problem
- Related works
- Our model
- Our algorithm
- Experiments
- Conclusion

# Model: Row-Augmented Matrix

## ■ Row-Augmented Matrix (RAM)

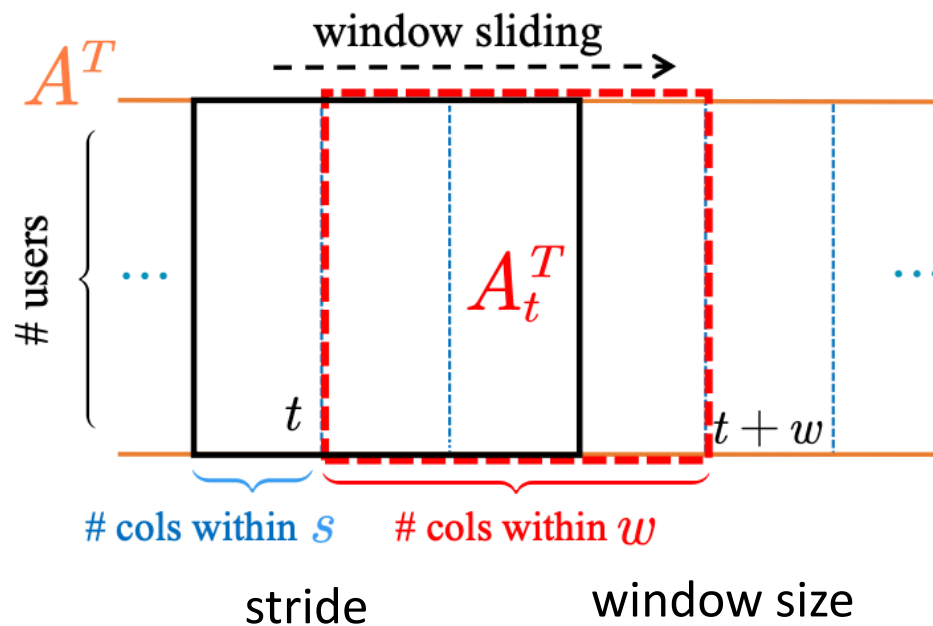
- rows are augmented with same cols
- *concatenate*



- equivalently, unfolding dynamic tensors (removing empty rows)
- *e.g. same restaurant changes overtime, upgrading of a product, etc.*

# Model: Row-Augmented Matrix

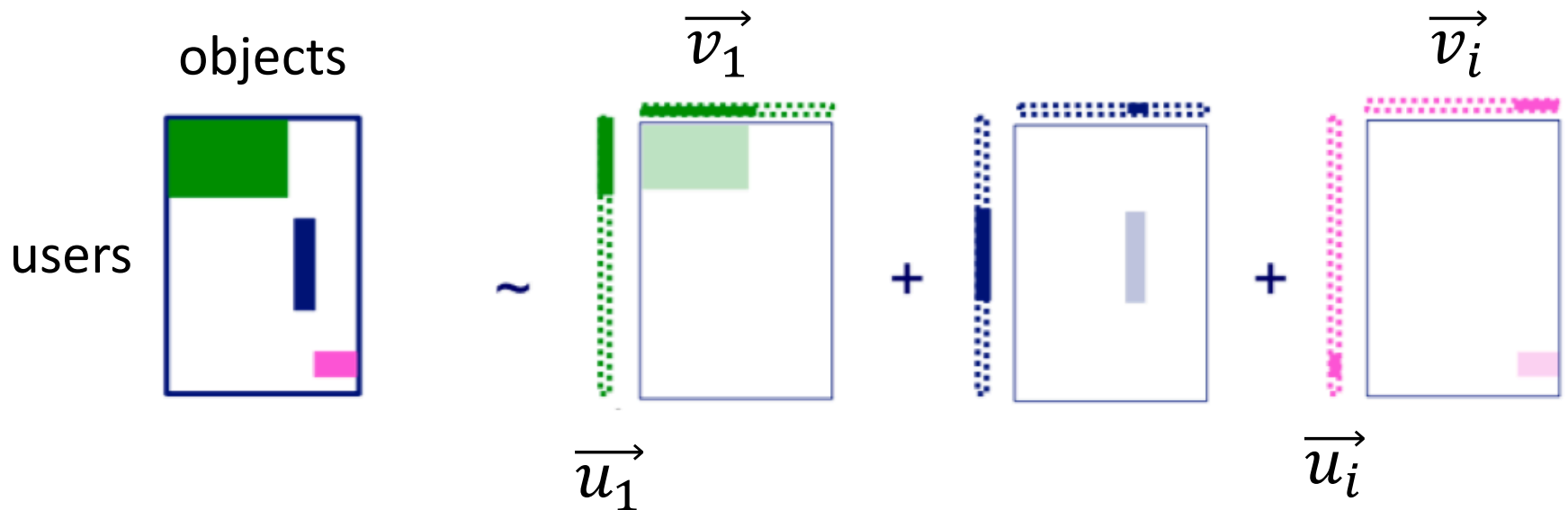
- RAM with sliding window



# Crash intro to SVD

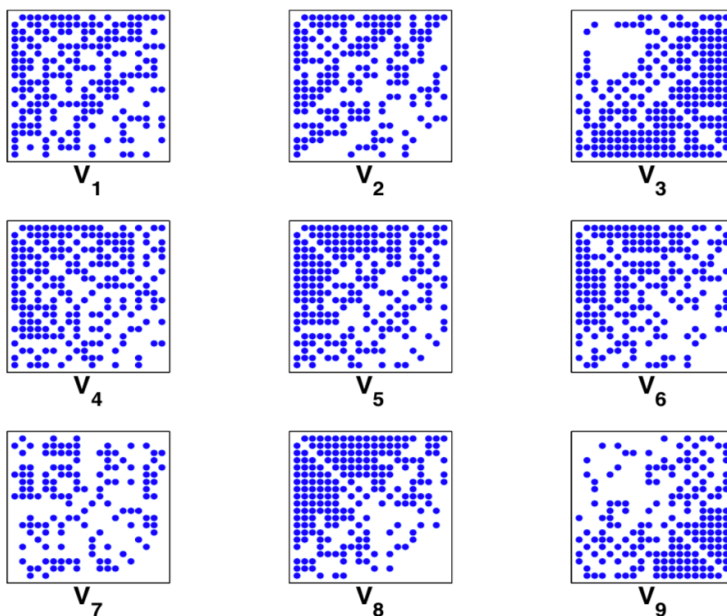
- (SVD) matrix factorization: finds blocks

$$A = U\Sigma V^T$$



# Properties of Singular Vectors

- Find **dense** groups of users by SVD
  - 20 nodes with **the highest magnitude projection** along the first 9 singular vectors

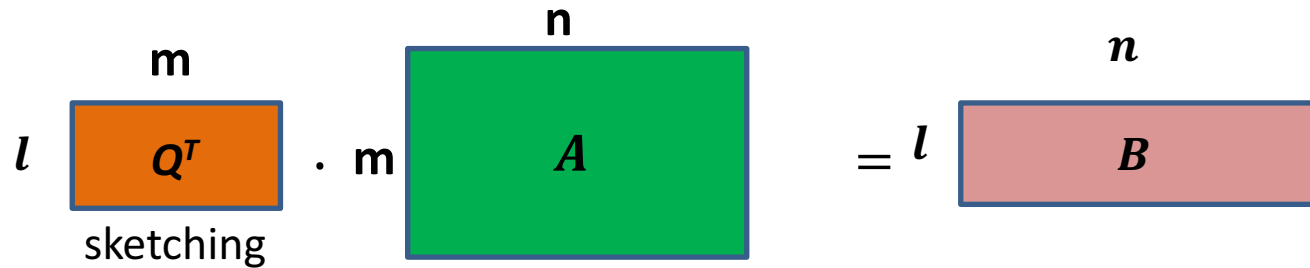


inducing sub-graphs  
contain **near-cliques**.

# Outline

- Problem
- Related works
- Our model
- **Our algorithm**
- Experiments
- Conclusion

# QB approximation



- (1)  $\Omega = \text{randn}(n, k + s)$
- (2)  $Q = \text{orth}(A\Omega)$
- (3)  $B = Q^T A$

$$\begin{aligned}
 A: m \times n, l &= k + s \\
 \Omega: n \times l \\
 Q: m \times l \\
 B: l \times n
 \end{aligned}$$

$$\left. \begin{aligned}
 B &= \tilde{U} \tilde{\Sigma} \tilde{V}^T \\
 A \approx QB &= Q \tilde{U} \tilde{\Sigma} \tilde{V}^T
 \end{aligned} \right\} \longrightarrow U = Q \tilde{U}, \Sigma = \tilde{\Sigma}, V = \tilde{V}$$

# AugSVD: incrementally building Q

■ generate matrices  $Q, B$  by  $G, H$

```

6: repeat
7:   Read rows  $\mathbf{a}$  for next stride  $s$  in augmented  $A$ 
8:    $\mathbf{g} = \mathbf{a}\Omega$ ;  $\mathbf{h} = \mathbf{a}^T \mathbf{g}$ 
9:    $glist.enqueue(\mathbf{g})$ ;  $hlist.enqueue(\mathbf{h})$ 
10: until the elements in  $glist$  corresponds to a window  $w$ 
11: for all  $\mathbf{g}$  in  $glist$ ,  $\mathbf{h}$  in  $hlist$  do
12:    $\mathbf{G} = [\mathbf{G}, \mathbf{g}]$ ;  $\mathbf{H} = \mathbf{H} + \mathbf{h}$ 
13: end for

```

$\mathbf{Q} = []$ ;  $\mathbf{B} = []$

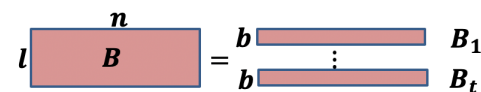
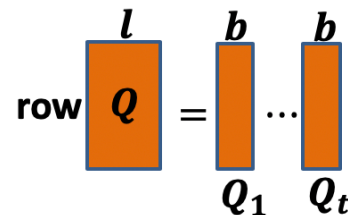
```

15: for  $i = 1, 2, \dots, t$  do
16:    $\Omega_i = \Omega(:, (i-1)b + 1 : ib)$ ;  $\mathbf{Y}_i = \mathbf{G}(:, (i-1)b + 1 : ib) - \mathbf{Q}(\mathbf{B}\Omega_i)$ 
17:    $[\mathbf{Q}_i, \mathbf{R}_i] = qr(\mathbf{Y}_i)$ 
18:    $[\mathbf{Q}_i, \tilde{\mathbf{R}}_i] = qr(\mathbf{Q}_i - \mathbf{Q}(\mathbf{Q}^T \mathbf{Q}_i))$ 
19:    $\mathbf{R}_i = \tilde{\mathbf{R}}_i \mathbf{R}_i$ 
20:    $\mathbf{B}_i = \mathbf{R}_i^{-T} (\mathbf{H}(:, (i-1)b + 1 : ib))^T - \mathbf{Y}_i^T \mathbf{Q} \mathbf{B} - \Omega_i^T \mathbf{B}^T \mathbf{B}$ 
21:    $\mathbf{Q} = [\mathbf{Q}, \mathbf{Q}_i]$ ;  $\mathbf{B} = [\mathbf{B}^T, \mathbf{B}_i^T]^T$ 
22: end for

```

$\mathbf{Y}_i: m \times b, \mathbf{Q}_i: m \times b$

$\mathbf{G}: m \times l$   
 $\mathbf{H}: n \times l$





# AugSVD: incrementally building Q

## ■ generate matrices $Q, B$ by $G, H$

```
6: repeat
7:   Read rows  $\mathbf{a}$  for next stride  $s$  in augmented  $A$ 
8:    $\mathbf{g} = \mathbf{a}\Omega$ ;  $\mathbf{h} = \mathbf{a}^T \mathbf{g}$ 
9:    $glist.enqueue(\mathbf{g})$ ;  $hlist.enqueue(\mathbf{h})$ 
10: until the elements in  $glist$  corresponds to a window  $w$ 
11: for all  $\mathbf{g}$  in  $glist$ ,  $\mathbf{h}$  in  $hlist$  do
12:    $\mathbf{G} = [\mathbf{G}, \mathbf{g}]$ ;  $\mathbf{H} = \mathbf{H} + \mathbf{h}$ 
13: end for
```

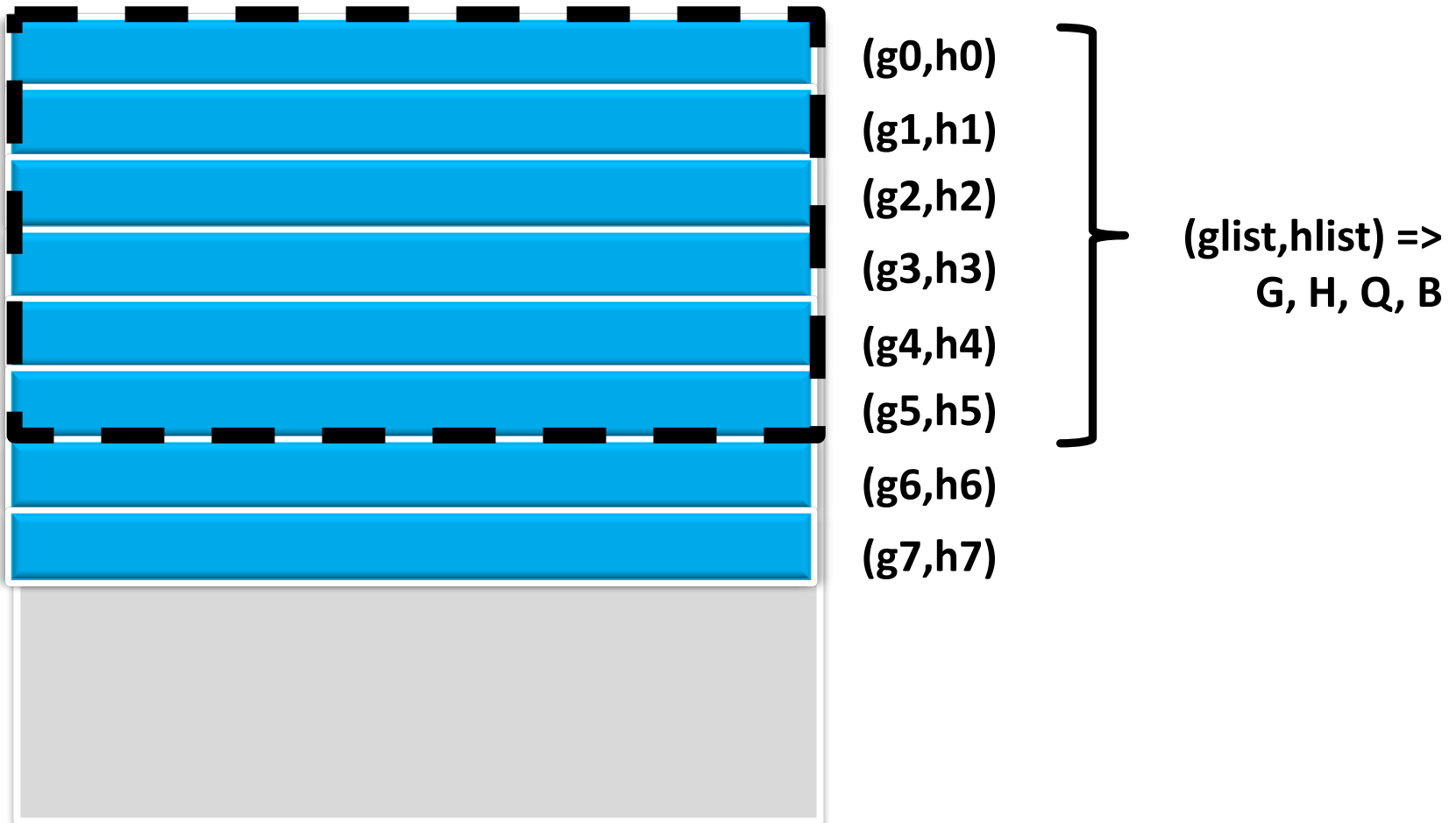
$\mathbf{G}: m \times l$   
 $\mathbf{H}: n \times l$

## ■ with Q and B

```
23:  $[\tilde{\mathbf{U}}, \mathbf{S}, \mathbf{V}] = svd(\mathbf{B})$ 
24:  $\mathbf{U} = \mathbf{Q}\tilde{\mathbf{U}}$ 
25:  $\mathbf{U} = \mathbf{U}(:, 1:k)$ ;  $\mathbf{V} = \mathbf{V}(:, 1:k)$ ;  $\mathbf{S} = \mathbf{S}(1:k, 1:k)$ 
```

# AugSVD

Combine *Sliding Window*, Change matrices  $G, H$  generation.



# Theoretical analysis

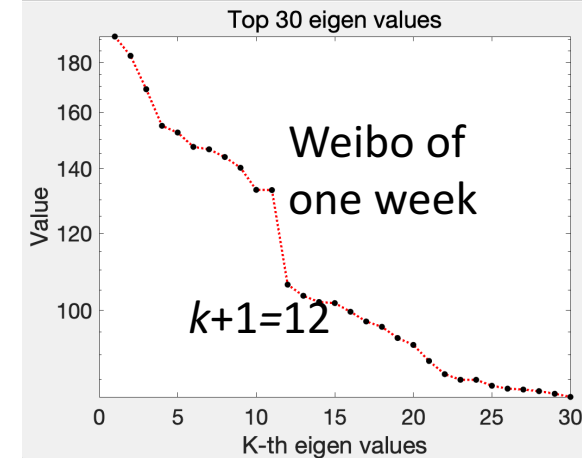
## ■ Theorem:

- Let approx error rate be  $\varepsilon_i = (\sigma_i - \hat{\sigma}_i)/\sigma_i$ , then

$$|\varepsilon_i| \lesssim 2 \frac{\sigma_{k+1}}{\sigma_i} + \frac{e\sqrt{2k+1}}{k} \left( \sum_{j=k+1}^{\min(m,n)} \left( \frac{\sigma_j}{\sigma_i} \right)^2 \right)^{1/2}, \quad i = 1, \dots, k$$

less than approximately.  $(k+1)$ -th original singular value  $k$  is the truncated length

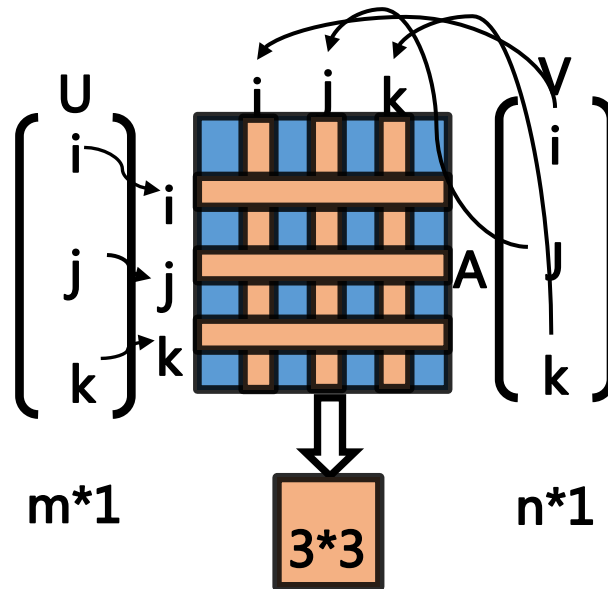
- Error is small when  $\sigma$  is *highly skewed*



# EigenPulse

- At every time stride,
  - Choose dense blocks based on the first several singular vectors.

✓ above average  $\tau_u = \frac{1}{\sqrt{m_t}}; \tau_v = \frac{1}{\sqrt{n}}$



# EigenPulse

- At every time stride,
  - Choose dense blocks based on the first several singular vectors.
  - [Optional] dense block detection in small selected blocks
    - ✓ use Fraudar and HoloScope (HS- $\alpha$ )
  - Calculate density

$$D_t(\text{rowset}, \text{colset}) = \frac{\sum_{i \in \text{rowset}} \sum_{j \in \text{colset}} \mathbf{A}_t(i, j)}{|\text{rowset}| + |\text{colset}|}$$

- plotting EigenPulse, and detecting anomalies.

# Outline

- Problem
- Related works
- Our model
- Our algorithm
- **Experiments**
- **Conclusion**

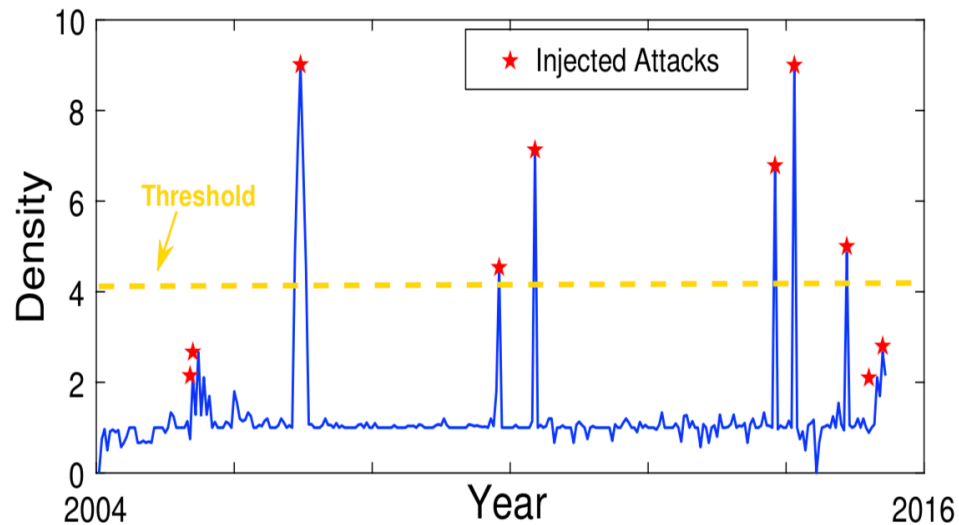
# Data statistics

**Table 1.** Datasets Statistic Information

Name	nodes	edges	span time
BeerAdvocate	26.5K × 50.8K	1.08M	Jan 08 - Nov 11
Yelp	686K × 85.3K	2.68M	Oct 04 - Jul 16
Amazon Phone & Acc	2.26M × 329K	3.45M	Jan 07 - Jul 14
Amazon Electronics	4.20M × 476K	7.82M	Dec 98 - Jul 14
Amazon Grocery	763K × 165K	1.29M	Jan 07 - Jul 14
Sina Weibo	2.74M × 8.08M	50.06M	Sep 01 - Dec 01

# EigenPluse: detect injection accurately and instantly

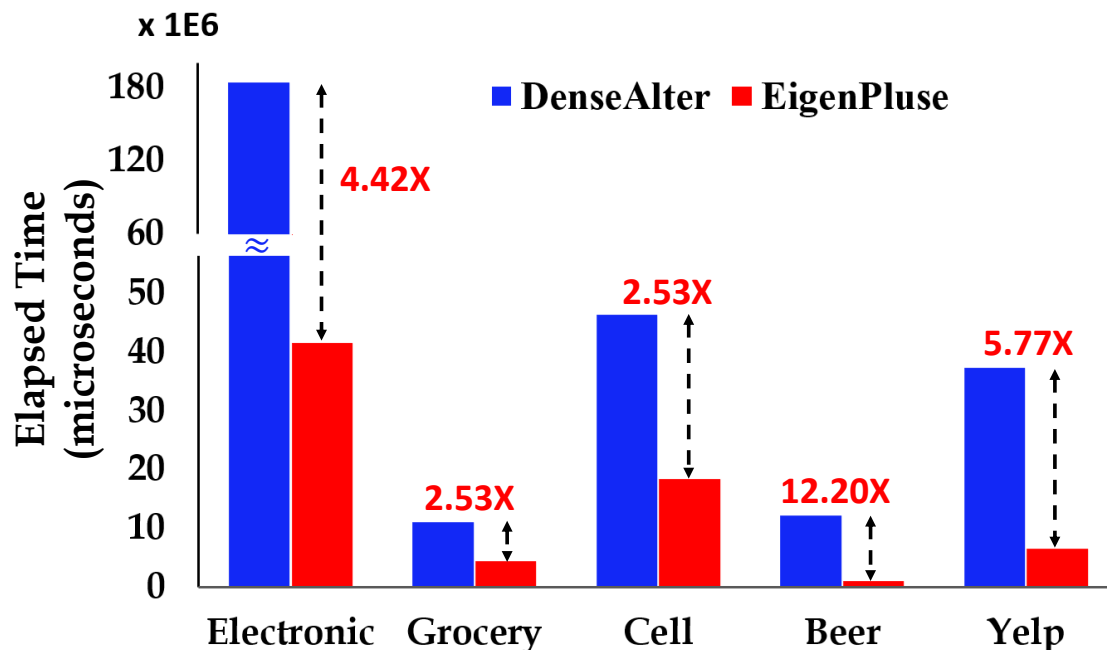
- Injected 10 dense blocks for *Yelp* dataset



*threshold = 99.7 percentile*

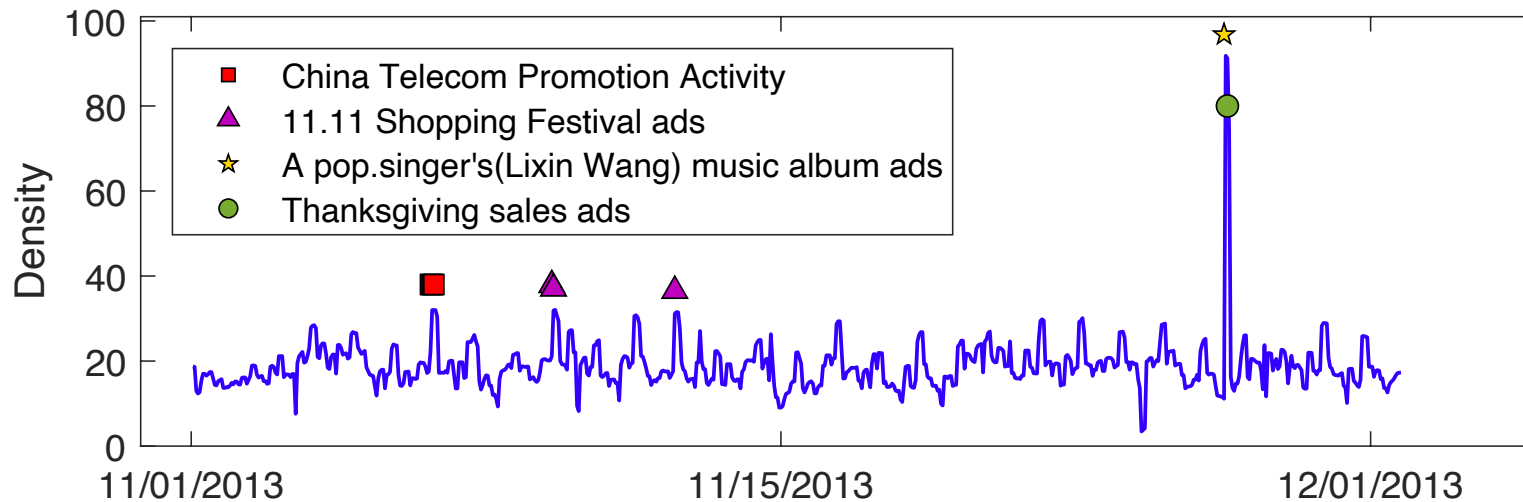


# EigenPluse: run faster than state-of-art methods



EigenPulse achieves more than  $2.53 \times$  speed up.

# EigenPluse: detect anomalous surges on Microblog data



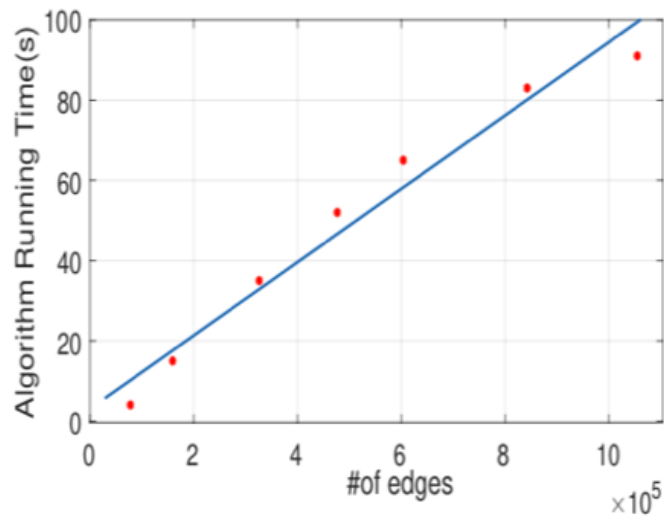
set  $w = 2h$  and  $s = 1h$  on Sina Weibo data

# Detected Blocks in Sina Weibo

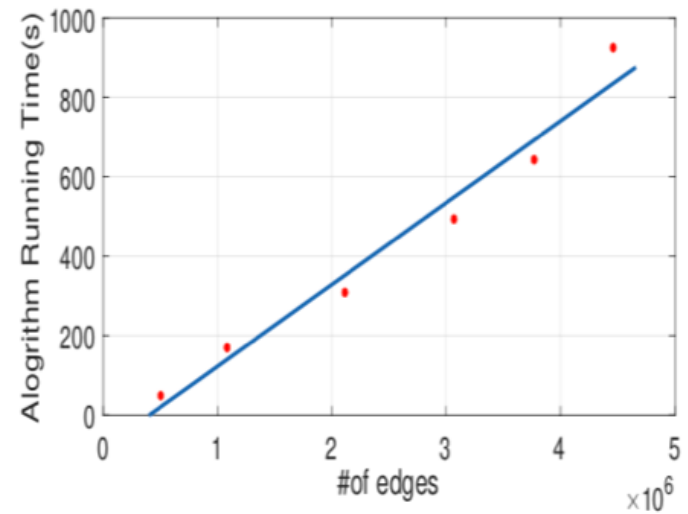
Message Topic	Size	Time range	#Edges	
China Telecom Promotion Activity	39 × 57	6:00~ 8:00, Nov 7	2,004	
	78 × 58	7:00~9:00, Nov 7	4,051	
	151 × 119	8:00~10:00, Nov 7	8,295	
11.11 Shopping Festival ads	201 × 139	6:00~8:00, Nov 10	7,012	
	196 × 111	7:00~9:00, Nov 10	9,668	≈ 100 <i>rt/user</i>
	126 × 93	8:00~10:00, Nov 13	638	
A pop. singer's (Lixin Wang) music album ads.	7 × 8	22:00~24:00, Nov 26	953	≈ 140 <i>rt/user</i>
Thanksgiving sale ads	26 × 36	23:00, Nov 26~1:00, Nov 27	629	
	43 × 34	1:00~3:00, Nov 27	263	

**7 users × 8 messages, 953 edges in 2 hours means every user retweeted **more than once per minute.****

# EigenPulse: Run linear with # of edges



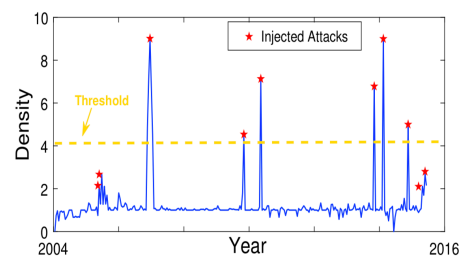
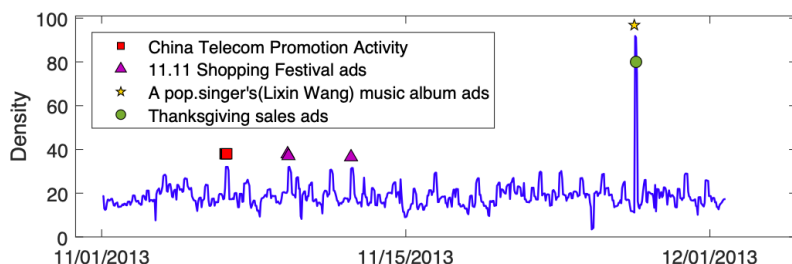
(a) **BeerAdvocate dataset**



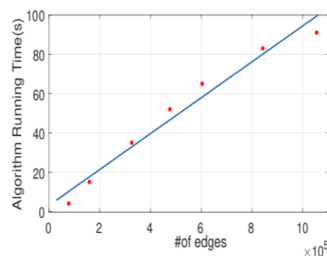
(b) **Amazon Electronic dataset**

# Conclusion

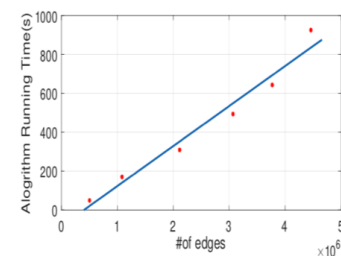
- Detect dense blocks given a streaming graph in form of triplet (*user, object, timestamp*)
- Robust and effective
  - theoretical robust approximation to batch SVD.



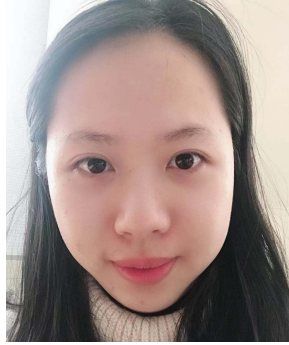
## ■ Scalable



(a) BeerAdvocate dataset



(b) Amazon Electronic dataset



Questions and Answers

**THANK YOU!**