

CatchCore: Catching Hierarchical Dense Subtensor

Wenjie Feng, Shenghua Liu, Huawei Shen, and Xueqi Cheng

Institute of Computing Technology ICT, CAS University of Chinese Academy of Science

fengwenjie@ict.ac.cn, liu.shenghua@gmail.com, {shenhuawei, cxq}@ict.ac.cn



中国科学院
CHINESE ACADEMY OF SCIENCES

Summary

- Goal:** to detect the hierarchical dense subtensors (HDS-tensors)
- Previous Work:**
 - Showed that dense subtensors in real-world tensors signal anomalies or fraud
 - Existing related methods assume subtensors are exclusive, and detect flatly and separately
 - Real-world tensors present hierarchical structure properties
- Proposed Detection method:**
 - **Unified metric:** (g, h, ϕ) -entry-plenum density measures
 - **CatchCore:** gradient based algorithm detecting hierarchical dense subtensors
 - **Quality measure:** MDL principle to evaluate the detection quality and select optimal parameters
- Result:**
 - **Accurate:** CatchCore detect densest subtensor and HDS-tensors with perfect performance
 - **Effective:** our algorithm detect anomalies in various different applications including periodical network intrusion attacks, dense co-authorship research group pattern
 - **Flexible and Stable:** our algorithm adapts to many density metrics and is robust with parameters
 - **Provable Scalability:** CatchCore runs in (sub-) linear time with all aspects of tensor

Proposed Algorithm: CatchCore

- Hierarchical Dense Subtensors Requirements
 - **Density:** significant density difference $\rho(\mathcal{B}^k) \geq \eta \rho(\mathcal{B}^{k-1}) \quad \eta > 1$
 - **Structure:** multi-layer cores $\mathcal{B}^k \leq \mathcal{B}^{k-1}$, i.e., $\mathcal{B}_n^k \subseteq \mathcal{B}_n^{k-1}, \forall n \in [N]$
- Given:** \mathcal{R} : N -way tensor; η : density ratio; K : the max #hierarchies
- Find:** $r (\leq K)$ significant HDS tensors $\{\mathbf{X}^1, \mathbf{X}^2, \dots, \mathbf{X}^r\}$
- Densest subtensor (block) detection formulation**

$$\max_{\mathbf{X}: \{\mathbf{x}_1, \dots, \mathbf{x}_N\}} \mathcal{F}(\mathbf{X}) = (1+p)\mathcal{R} \bar{\times} \mathbf{X} - p \prod_{\mathbf{x}_n \in \mathbf{X}} \|\mathbf{x}_n\|_1$$

subject to $\mathbf{x}_n \in \{0, 1\}^{|\mathcal{R}_n|}, \forall n = 1, \dots, N$.
- HDS-tensors detection Alg.**
 - **Gradient based Optimization Algorithm**

$$\max_{\mathbf{X}^1, \dots, \mathbf{X}^K} \sum_{k=1}^K \mathcal{F}(\mathbf{X}^k)$$

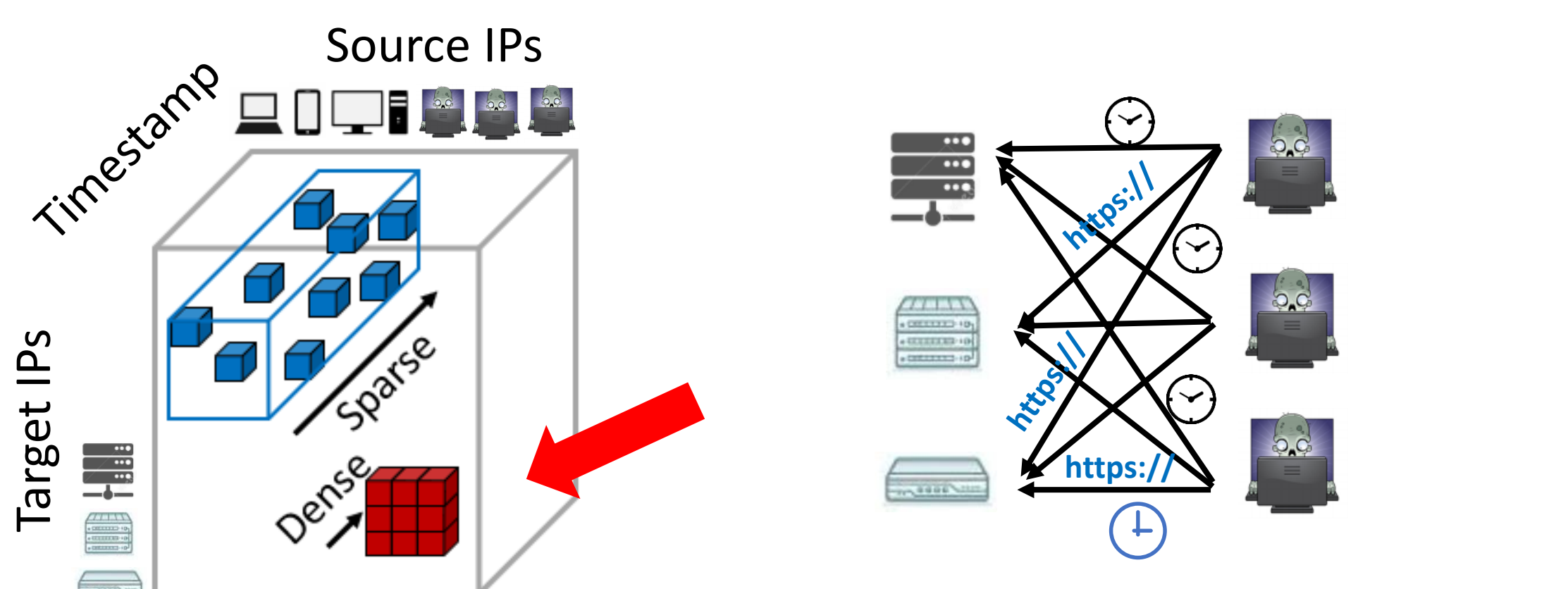
subject to $\rho_{\mathbf{X}^{h+1}} \geq \eta \rho_{\mathbf{X}^h}$
 - **Alternative Gauss-Seidel for updating**

$$\mathbf{X}_{(h+1, n, \cdot)} \leq \mathbf{X}_{(h, n, \cdot)} \leq \mathbf{X}_{(h-1, n, \cdot)}$$

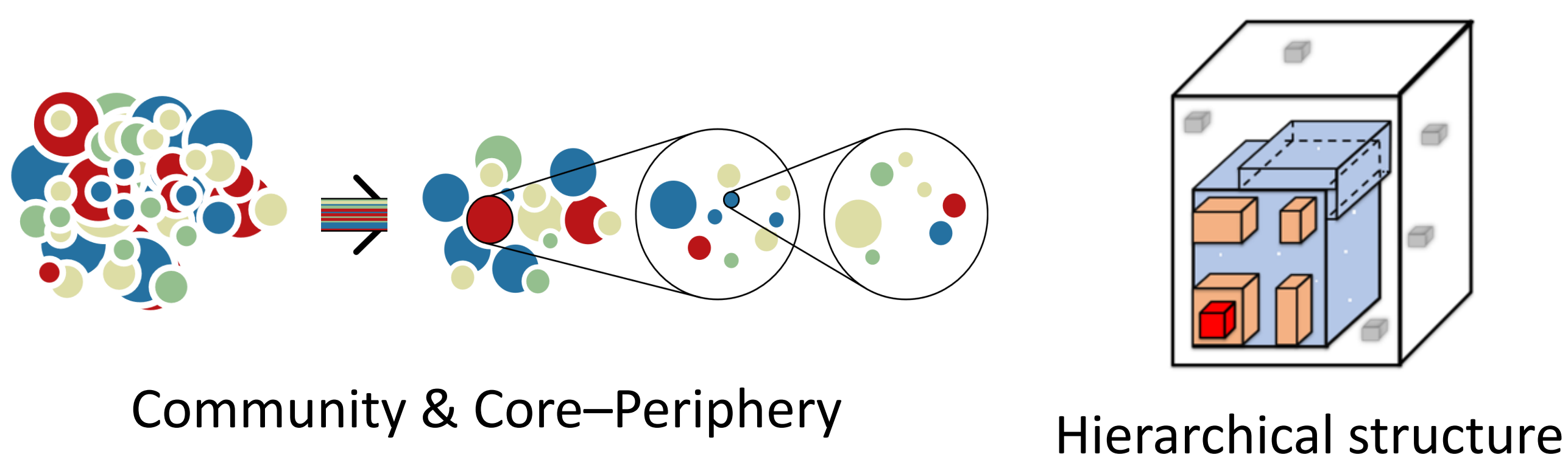
$\forall h = 1, \dots, K; \forall n = 1, \dots, N$
- Using MDL to evaluate the quality of detection result of given parameter configuration: $M^* = \operatorname{argmin} L(M) + L(D|M)$

Motivation

- Dense subtensor in various multi-aspect data (tensor) e.g., the TCP dump for network intrusion, etc.



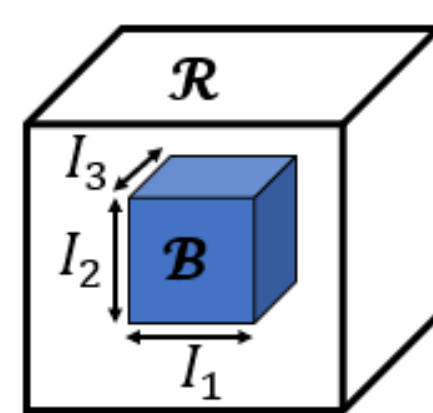
- Hierarchical structures in real-world data



How can we detect hierarchical structures for multi-aspect data?

Terminology

- A subtensor \mathcal{B} in a tensor \mathcal{R}
 - Sub-tensor inclusion: $\mathcal{B} \leq \mathcal{R}$
 - Mass: $M_{\mathcal{B}}$, Volume: $V_{\mathcal{B}}$, Cardinalities: $D_{\mathcal{B}}$, Density: $\rho_{\mathcal{B}}$
- Indicator vectors paradigm
 - $\mathbf{X}_{\mathcal{B}} = \{\mathbf{x}_n \in \{0, 1\}^{|\mathcal{R}_n|}; \forall n \in [N]\}$
 - $M_{\mathcal{B}} = \mathcal{R} \bar{\times} \mathbf{X}_{\mathcal{B}} = \mathcal{R} \times_1 \mathbf{x}_1 \cdots \times_N \mathbf{x}_N$
 - \times_n : n -mode tensor-vector product
- (g, h, ϕ) -entry-plenum density measures
 - $M_{\mathcal{X}}, S_{\mathcal{X}}$: mass and size (vol. /card.) of subtensor \mathcal{X}
 - g, h : two increasing functions, ϕ : constant factor



Date	2	0	4	6
Sep-8	1	7	5	1
Sep-7	1	7	5	1
Alex	1	7	5	1
Chris	0	4	3	3
Dora	2	1	0	1
	A	B	C	

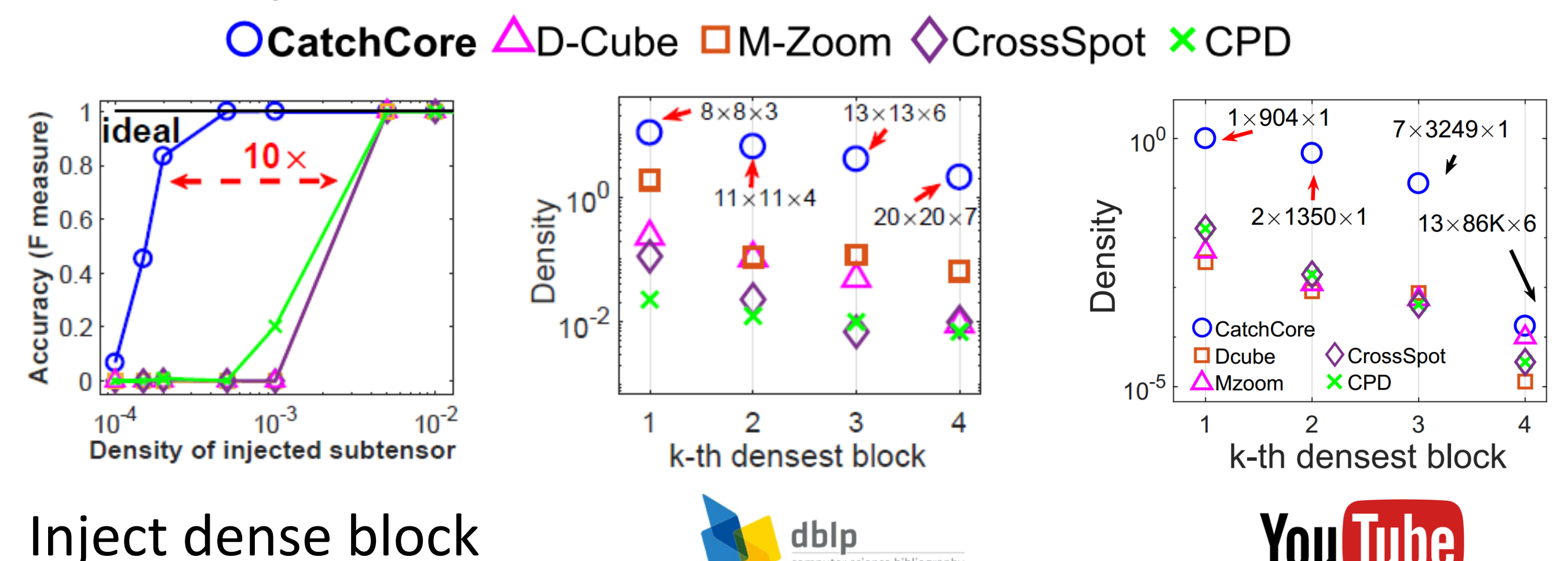
User: $\mathbf{x}_1 = [1, 1, 0]$
Item: $\mathbf{x}_2 = [0, 1, 1]$
Date: $\mathbf{x}_3 = [1, 0, 1]$
 $\mathbf{X}_{\mathcal{B}} = \{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3\}$

$M_{\mathcal{B}} = 36, V_{\mathcal{B}} = 2 \times 2 \times 2 = 8$

$$\mathcal{F}_{\phi}(\mathbf{X}) = \begin{cases} 0 & \mathbf{X} = \{\{0\}^{|\mathcal{R}_n|}; \forall n \in [N]\} \\ g(M_{\mathbf{X}}) - \phi \cdot h(S_{\mathbf{X}}) & \text{otherwise} \end{cases}$$

Experiment Results

- Q1 Accuracy:** detect the densest block / HDS-tensors



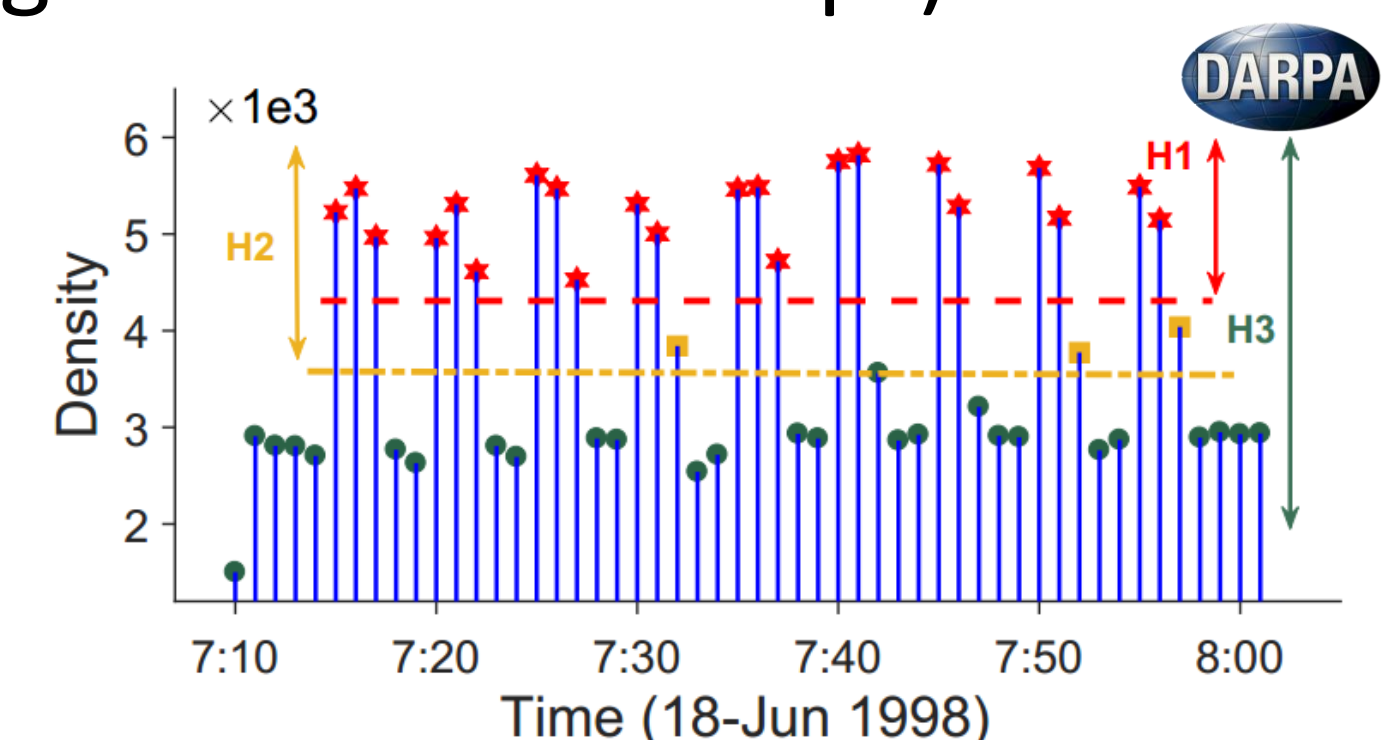
Inject dense block

- Q2 Effectiveness:** patterns detected in real-world tensors

- (1) TCP dump: (source IPs X target IPs X timestamps)

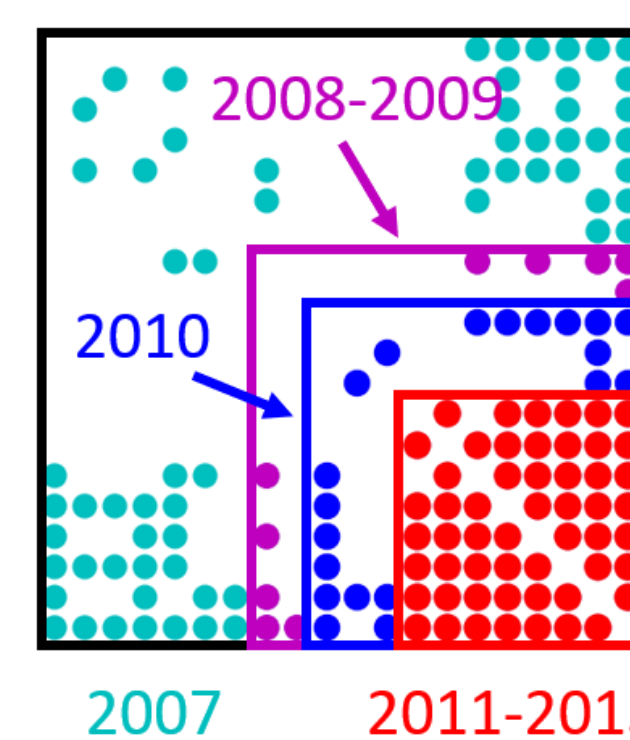
H	Shape	Ratio
1	$1 \times 1 \times 97$	100%
2	$1 \times 1 \times 100$	100%
3	$1 \times 1 \times 274$	100%
4	$15 \times 5 \times 24.7K$	87.0%
5	$171 \times 15 \times 29.2K$	85.4%

Top 5 HDS-tensors & attack ratio

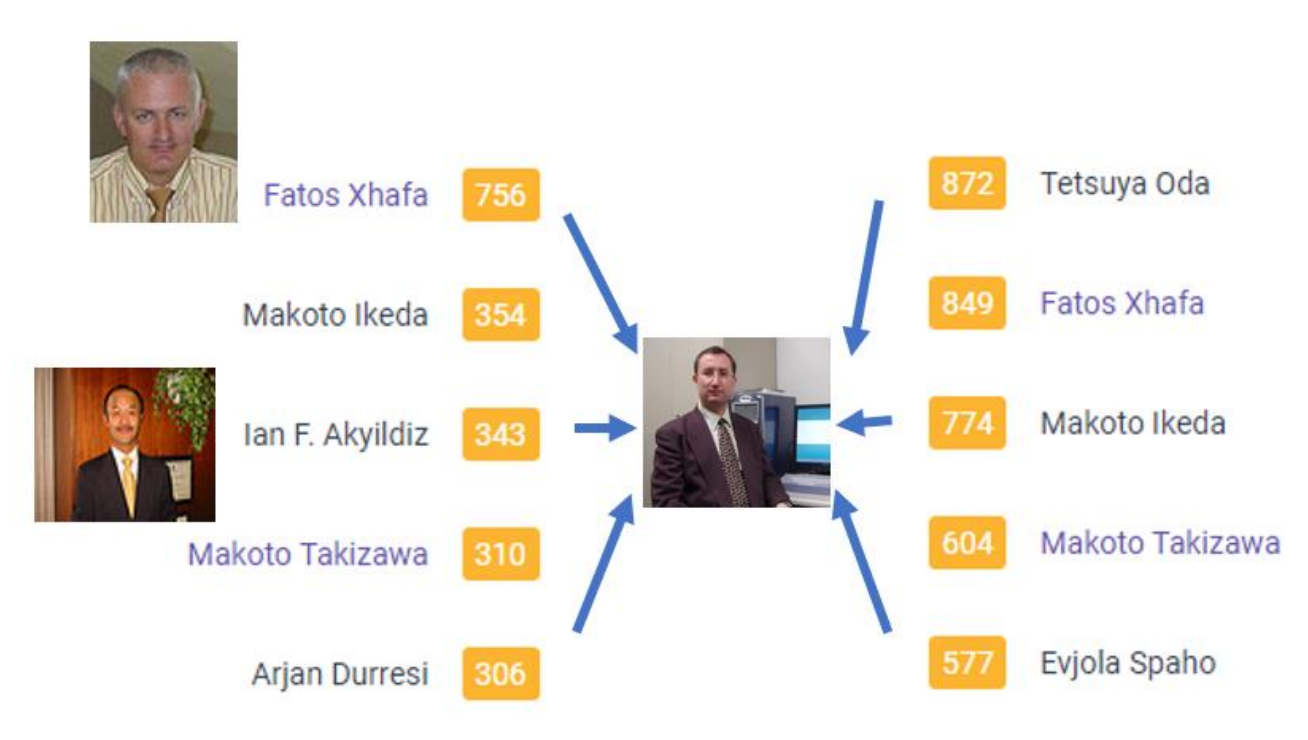


Periodical Neptune attack

- (2) DBLP co-authorship: (authors X authors X years)

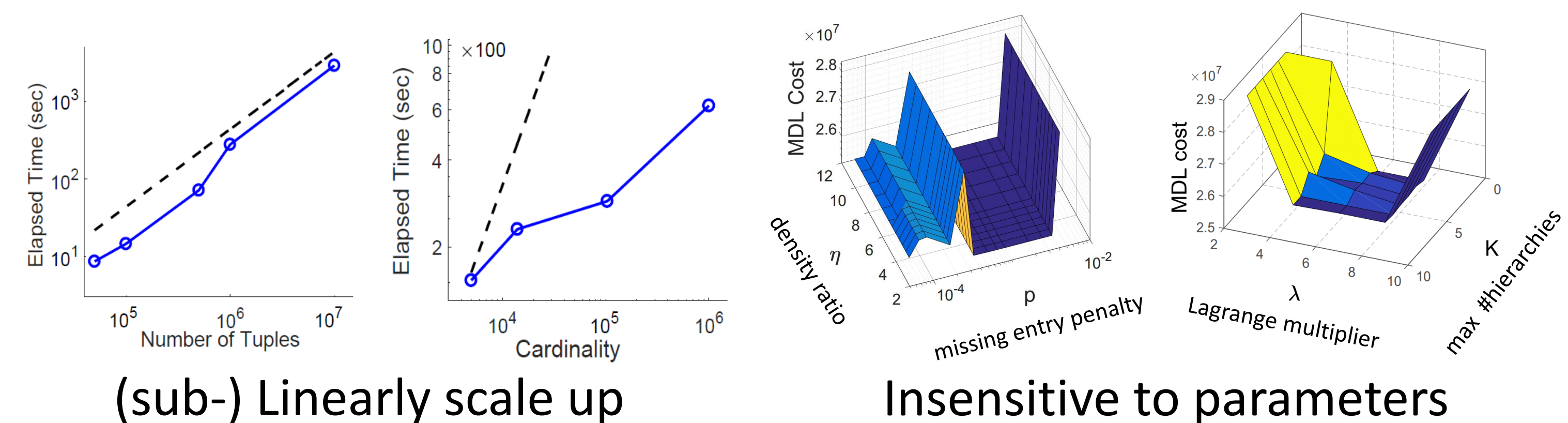


Top 4 HDS-tensors



Research group

- Q3 Scalability & Robustness:** CatchCore is linear in all aspects of input, achieves the optimal results for wide parameters range



(sub-) Linearly scale up

Insensitive to parameters