

ACM Web Conference 2025

## Pontus: A Memory-Efficient and High-Accuracy Approach for Persistence-Based Item Lookup in High-Velocity Data Streams

Weihe Li<sup>1</sup>, Zukai Li<sup>1</sup>, Beyza Bütün<sup>2</sup>, Alec F. Diallo<sup>1</sup>, Marco Fiore<sup>2</sup>, Paul Patras<sup>1</sup> <sup>1</sup>University of Edinburgh, Edinburgh, United Kingdom <sup>2</sup>IMDEA Networks Institute, Madrid, Spain



THE UNIVERSITY of EDINBURGH



# **Data Stream Processing**



- User behavior analysis ----> important to improve user experience
- Persistently accessed websites (*persistent items*) ----- user preferences

### Persistence



ltem	Persistence
eı	2
@2	1
e3	1
<b>e</b> 5	1



#### Fast processing speed

e.g. 10 Gb/s data stream: each item every 67 ns

#### Limited fast memory

# 11 Cache: around 64KB<sup>[1]</sup> Infeasible to store information for all items

[1] Li, W. and Patras, P. Tight-sketch: A high-performance sketch for heavy item-oriented data stream mining with limited memory size. ACM CIKM 2023. 4

- Sketches: compact data structure by hashing
  - Idea: hash data into limited space



[1] Zhang, Yinda, et al. "On-off sketch: A fast and accurate sketch on persistence." *Proceedings of the VLDB Endowment* 14.2 (2020): 128-140.









# Limitations of Existing Sketch Methods

#### Low detection accuracy under limited memory budgets

Persistent flows being evicted from the bucket by non-persistent ones due to the highly skewed traffic distribution.



# Limitations of Existing Sketch Methods

#### **Low Memory Efficiency**

Track multiple features per item to better protect persistent items.



[2] Li, W. and Patras, P. Stable-sketch: A versatile sketch for accurate, fast, web-scale data stream processing. ACM WWW 2024. 11

# **Our Contributions**

#### • Pontus

- A novel method for persistent item lookup
- High accuracy, high memory-efficiency and fast processing speed
- Deployable on the practical hardware, Tofino programmable switch

# Data Structure



# Update

Case 1:



|--|











Update

Case 2:





The incoming item can replace the tracked item only if its counter has decayed to zero.





Case 3:





# Only a single scan of all buckets is required to determine which bucket contains a value higher than the predefined threshold.

Query

# Evaluation

- Software -- CPU Platform (Intel(R) Core(TM) i5-1135G7 @ 2.40GHz processor, C++)
- Hardware -- Tofino Switch (P4)
- Traces -- CAIDA 2018 and 2019<sup>[3]</sup>

# Evaluation – Accuracy (CPU)



# Evaluation – Speed (CPU)



**Higher Speed!** 117.3% higher than On-Off



# Evaluation – Resource Usage (Tofino)

Resource	Usage	Resource	Usage	
Hash Bit	5.7%	Match Crossbars	4.6%	
Gateways	16.7%	Logical Table ID	21.9%	
VLIW Instruction	7.3%	SRAM	4.3%	
Total Average		8.2%		
		Limit	ed Overhea	

#### More results

- Formal analysis
- Evaluation on more tasks, like persistence estimation
- More results on the CPU platform and Tofino switch

• Code: <a href="https://github.com/Mobile-Intelligence-Lab/Pontus">https://github.com/Mobile-Intelligence-Lab/Pontus</a>

This research was supported by the SNS JU and the European Union's Horizon Europe research and innovation program under Grant Agreement No. 101139270 (ORIGAMI). Beyza Bütün is a Comunidad de Madrid predoctoral fellow (PIPF-2022/COM-24867). Weihe Li was partially supported by Cisco through the Cisco University Research Program Fund (Grant no. 2019-197006).