

FPGA-Enhanced Real-Time Analysis of Network Traffic

Yu-Kuen Lai, Drixter Velayo Hernandez, Hsiang-Lun Hua, Bo-Shun Huang
Computer Networks & Systems Research Lab
Department of Electrical Engineering
Chung-Yuan Christian University, Chungli, Taiwan

THE 5th GLOBAL RESEARCH PLATFORM WORKSHOP

September 16-17, 2024 Osaka, Japan





Outline

- Motivation
 - Common Network Traffic Measurement Tasks
 - Challenges
- General Backgrounds
 - Sketch-based Data Stream Algorithms
 - Flow-size/Flow-length Distribution
- High-Speed Real-Time Processing
 - Estimating Weibull Model Parameters
 - CAIDA 2007 DDoS
 - Witty Worm
- Summary
- □ Q&A

CNSRL



Common Measurement Tasks

- Heavy Hitter Detection
- Shannon Entropy
 - One of the effective methods for Network Anomaly Detection
- Flow Size/Length Distribution Estimation plays a fundamental role in
 - Understanding network behavior
 - Optimizing performance
 - Managing resources, and
 - Detecting anomalies

CNSRL

Bohan Zhao, Xiang Li, Boyu Tian, Zhiyu Mei, and Wenfei Wu. 2021. DHS: Adaptive Memory Layout Organization of Sketch Slots for Fast and Accurate Data Stream Processing. In Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery. & Data Mining. ACM, Virtual Event Singapore, 2285 – 2293. https://doi.org/10.1145/3447548.3467353





Challenges

- Packet arrives at a rapid rate
- Per-Flow Traffic Measurement is challenging
 - Memory space
 - Processing time
- Key space
 - IPv4 address of 32-bit
 - High distinct number of flows
- How to process and compute statistics on packet streams in real-time?
- Typical Implementation
 - Software-based
 - Off-line approaches

CNSRL

A. K. Marnerides, A. Schaeffer-Filho, and A. Mauthe, "Traffic Anomaly Diagnosis in Internet Backbone Networks," B. Computer. Networks, vol. 73, no. C, pp. 224–243, Nov. 2014.





Sketch-based Data Stream Algorithms

- A probabilistic data structure
- "Processing a long stream of data items in one pass using small memory space" [Kumar et al., IMC 03]
- Widely used in networking and data stream processing
 - Capable of summarizing large amounts of data
 - Update/Query operations in a memory-efficient manner
- Especially useful in scenarios of
 - Infeasible to store all data due to resource constraints
 - High-speed Network Monitoring
 - Anomaly Detection, and Traffic Engineering

CNSRL





Flow Size/Length

- Flow
 - Flow is defined by selected attributes of layer 3-4 packet header.
 - For example, by 5-tuple of:

{IP_Src, IP_Dest, TCP_SrcP, TCP_DestP, Proto}

- Flow Size
 - The total amount of data transferred within a flow
 - In terms of the number of bytes
- Flow Length
 - The total amount of packets transferred within a flow

CNSRL



Security and Anomaly Detection

- Sudden Increase in Large/Small Flows
- Changes in Flow Length Distribution
 - DDoS attacks
 - Attackers generate large numbers of short flows to overload web servers or services
 - Network scanning
 - Numerous short flows could be a sign of network scanning
- Changes in Flow Size Distribution
 - Some malware or Worms may generate large flows
 - Attackers overwhelm network bandwidth with large amounts of data

CNSRL

電機工程學系



- Can we do high-level coarse-grained estimation to track Changes of Flow Length and Flow Size?
- At wire-speed of 100Gbps and beyond?



CNSRL



Statistical Model of Flow Size/Length Distribution

- Highly-skew, Long-tail distribution
 - Lognormal, Pareto and Weibull are commonly used for the long-tail network traffic modeling [2] [24]
- □ Due to the complexity of Internet traffic, single distribution can not be fitted accurately [26]
- Weibull distribution is a flexible model
 - Can accommodate various shapes
 - Including the exponential and Rayleigh distributions as special cases
 - Tracking the change of the parameter
 - Real-time estimation of Weibull shape parameter is useful in traffic monitoring and anomaly detection [4]

CNSRL





Weibull and Log-Normal Model for Flow Length/Size Distribution

PARAMETER ESTIMATION

CNSRL



Weibull Model

- The probability density function of a Weibull random variable
 - \circ α : Shape
 - β: Scale
 - δ: Location

$$f(x) = \frac{\alpha}{\beta} \left(\frac{x - \delta}{\beta}\right)^{\alpha - 1} e^{-((x - \delta)/\beta)^{\alpha}}$$

CNSRL

雷俄工程學系



Method of Moments

- A technique for estimating the parameters of statistical distributions
- By equating the population moments to the sample moments
- Particularly useful for distribution models of
 - Weibull
 - Log-Normal

CNSRL



Estimating the First Moment m_1

- □ Take the *Flow Length* as an example:
- Assuming the flow id is the source IP address
- Packet counts of flows are updated in the counter of C[d][w]
- The total packets count of all IP flows is the total sum of each row

$$\sum_{id=0}^{n-1} x_{id} = \sum_{i=0}^{w-1} C[0][i]$$

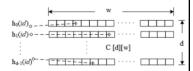
□ The first frequency moment F₁

$$F_1 = \sum_{i=0}^{w-1} C[0][i]$$

CNSRL

$$\frac{\overbrace{m_1}}{\sqrt{(cm_2)}} = \frac{\Gamma(1+\frac{1}{\alpha})}{\Gamma(1+\frac{2}{\alpha})-\Gamma^2[1+\frac{1}{\alpha}]}$$

$$m_1 = \mu = \frac{1}{n} \sum_{id=0}^{n-1} x_{id}$$



電鐵工程學系



Zeroth Frequency Moment F_o

$$m_1 = \mu = \sum_{i,d=0}^{n-1} x_{id}$$

- The distinct count of the observed flows: n
- lacksquare Can be estimated by using the sketch-based algorithms

 $n \approx \hat{F_0}$

- FM sketch [#flajolet_probabilistic_1983]
- HyperLogLog [#flajolet_hyperloglog_2007]

CNSRL



Estimation of the Sample Moments

 $lue{}$ Get a function with only the shape parameter α

$$\frac{m_1}{\sqrt{(cm_2)}} = \frac{\Gamma(1+\frac{1}{\alpha})}{\Gamma(1+\frac{2}{\alpha})-\Gamma^2[1+\frac{1}{\alpha}]}$$

■ The first moment about the origin

$$\hat{m_1} = \frac{F_1}{\hat{F_0}}$$

■ The second moment about the mean

$$c\hat{m}_2 = \frac{\hat{F}_2}{\hat{F}_0} - (\frac{F_1}{\hat{F}_0})^2$$

CNSRL

雷俄工程學系

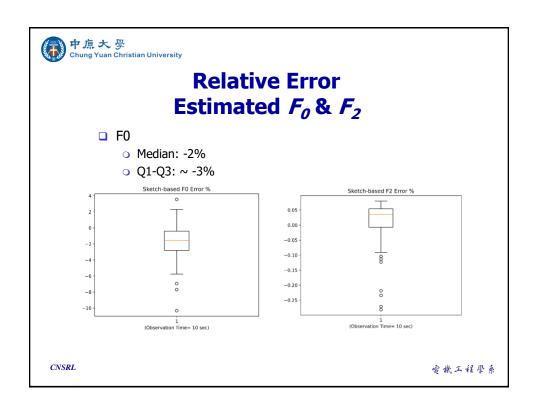


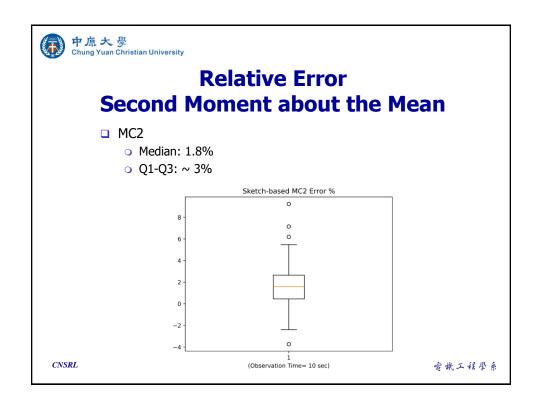
FPGA Functional Simulations

- Count-Min Sketch
 - \circ cm d = 7
 - \circ cm_w = 65536
- FM Sketch PCSA
 - o fm_n = 512
 - o fm_bits = 32
- Trace: MAWI, 200403201400.dump.gz

CNSRL

200403201400.dump.gz_obsT10_flow_length_sip_2024-08-18





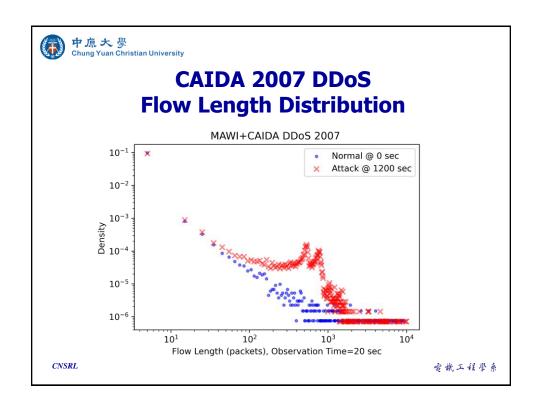


Synthetic DDoS Network Traffic

- CAIDA+MAWI
- CAIDA 2007 DDoS Trace (Attacking):
 - The attacker sent many ICMP Echo request packets to the Victim
 - Randomly created source IP addresses
 - Four DDoS attacking traffic are selected from the CAIDA 2007 DDoS trace (to-victim).
 - 20070804_140436.pcap
 - 20070804_140936.pcap
 - 20070804_141436.pcap
 - 20070804_141936.pcap
- MAWI DITL 2019 Trace (Background)
 - https://mawi.wide.ad.jp/mawi/ditl/ditl2019/201904091800.html

CNSRL

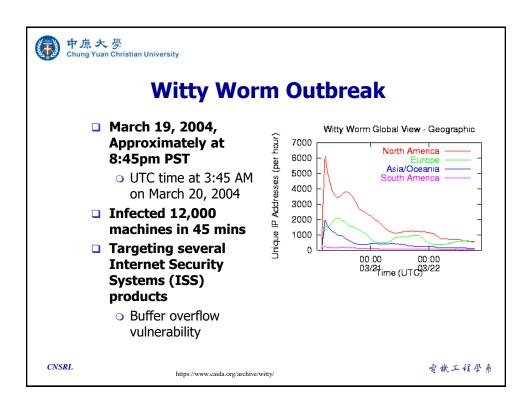
雷俄工程學系





WITTY WORM

CNSRL 電鐵工程學系





Witty Worm Characteristics

- The infecting host sends
 - o 20,000 **UDP** packets
 - Random destination IP address
 - Random size between 796 and 1307 bytes
 - o from source port 4000

CNSRL

CNSRL

https://www.caida.org/archive/witty/

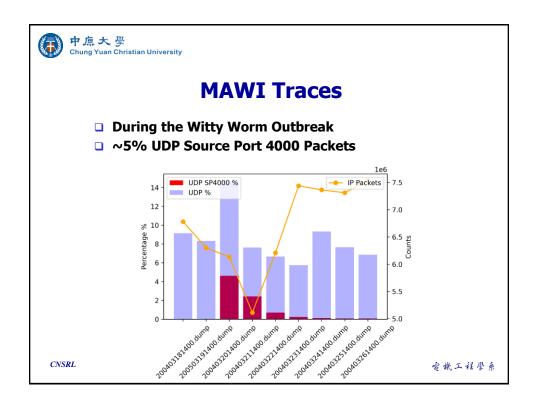
電機工程學系

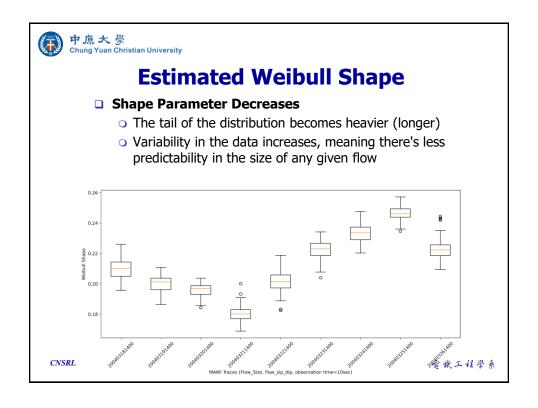


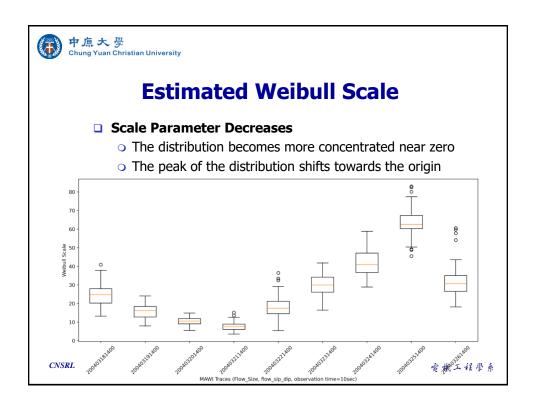
MAWI 2004 Data Set Witty Worm

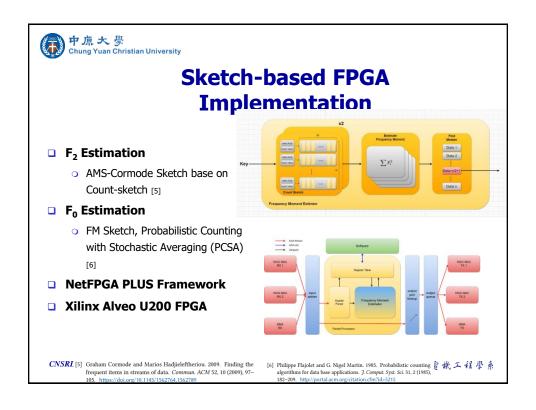
- Outbreak begins at UTC Time 3:45 AM on March 20, 2004
- □ Captured by the daily 15-min traces

MAWI Trace 2004	UTC Time 2004	
200403181400.dump.gz	03/18 · 05:00:01 AM	
200403191400.dump.gz	03/19 · 05:00:01 AM	
200403201400.dump.gz	03/20 · 05:00:01 AM	****
200403211400.dump.gz	03/21 · 05:00:01 AM	****
200403221400.dump.gz	03/22 · 05:00:01 AM	**
200403231400.dump.gz	03/23 · 05:00:01 AM	*
200403241400.dump.gz	03/24 · 05:00:01 AM	*
200403251400.dump.gz	03/25 · 05:00:01 AM	
200403261400.dump.gz	03/26 · 05:00:01 AM	





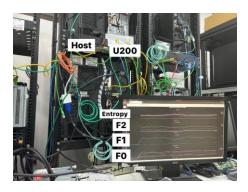






Xilinx FPGA Demo

- Replays the synthetic trace from a 2-port 100Gbps NIC
- Observation time of 30 seconds



CNSRL

電機工程學系



Summary

- Real-time online processing
 - o exceeding 100Gbps line speed
- □ For less than 132K bytes of memory space
 - Count Sketch: 128K Bytes
 - FM Sketch PCSA: 4K Bytes
- □ Flow-Length/ Flow-Size Distribution
 - Weibull Parameter Estimation
 - Log-Normal Parameter Estimation
- Sketch-based Implementation
 - Xilinx Alveo U200 FPGA /w NetFPGA PLUS Framework
- Suitable to implement on the high-speed network equipment





CNSRL

