# SPADA: SPark Anomaly Detection Ace

Antonia Affinito*, Alessio Botta*+, Luigi Gallo*, Mauro Garofalo*, and Giorgio Ventre*+

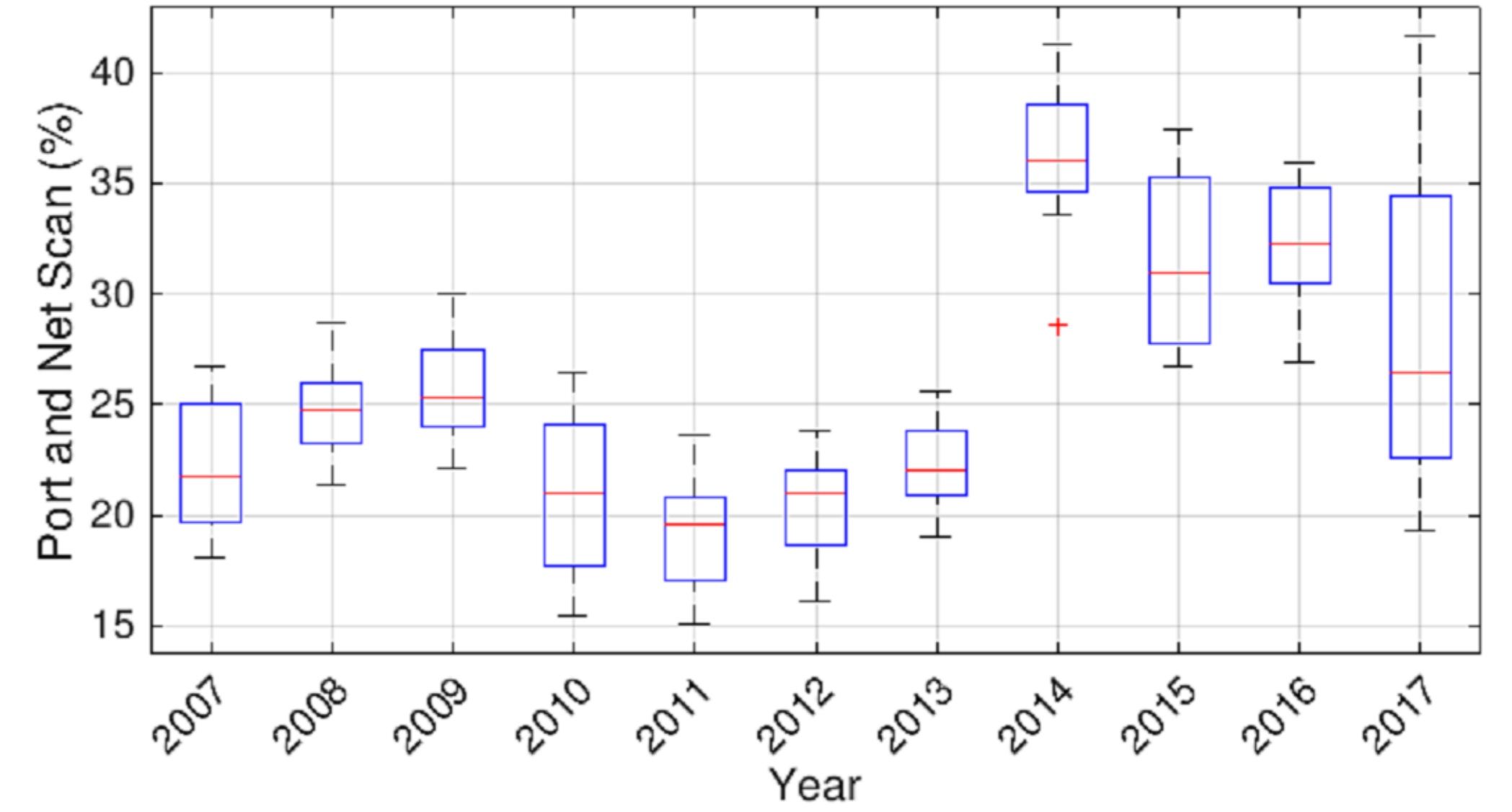*Department of Electrical Engineering and Information Technologies, University of Napoli "Federico II" – Italy

+ NM2 SRL, Italy

{antonia.affinito, luigi.gallo10}@studenti.unina.it {a.botta, mauro.garofalo, giorgio}@unina.it

## INTRODUCTION AND MOTIVATION

- An important category of network anomaly detection are port and net scan
  - The percentage of attacks and anomalous event in network traffic is constantly growing (see Figure 1)
- We present an approach to detect anomalies in high-speed networks working at flow-level
  - We use Apache Spark to cope with the problem of the large amount of data to be analyzed
  - We implement a simple threshold-based detection algorithm in Spark and test it by using several real traces


Figure 1: Malicious activities during the last 11 years

## ALGORITHM

- We consider the ratio between the number of flows generated and received by the same IP address together with other important features

$$\frac{FlowS}{FlowD} + \frac{\alpha}{PKTF} - \frac{BPP}{\gamma} - \beta\frac{FlowS}{IPC} > TH$$

- FlowS: number of generated flows;
- FlowD: number of received flows;
- PKFT: average number of packets per flow;
- BPP: average number of bytes per packet;
- TH: threshold.

## TOOLS



- A platform for distributed processing of Big Data [2]
- Very fast both in storage and data processing because of *in-memory* processing

**MAWI** (Measurement and Analysis of the Wide Internet)
- An archive of traces of real traffic provided by the MAWI Working Group [4]
- Traffic captured every day from 14:00 to 14:15 on a transoceanic link

**MAWILab**
- An approach for the identification of network anomalies in MAWI [3]
- Uses four detectors: Principal Component Analysis (PCA), Gamma distribution, Kullback Leibler (KL) divergence, and Hough transformation

## SPADA



SPADA is a system able to run automatically all operations needed for the analysis of a traffic trace. Every day
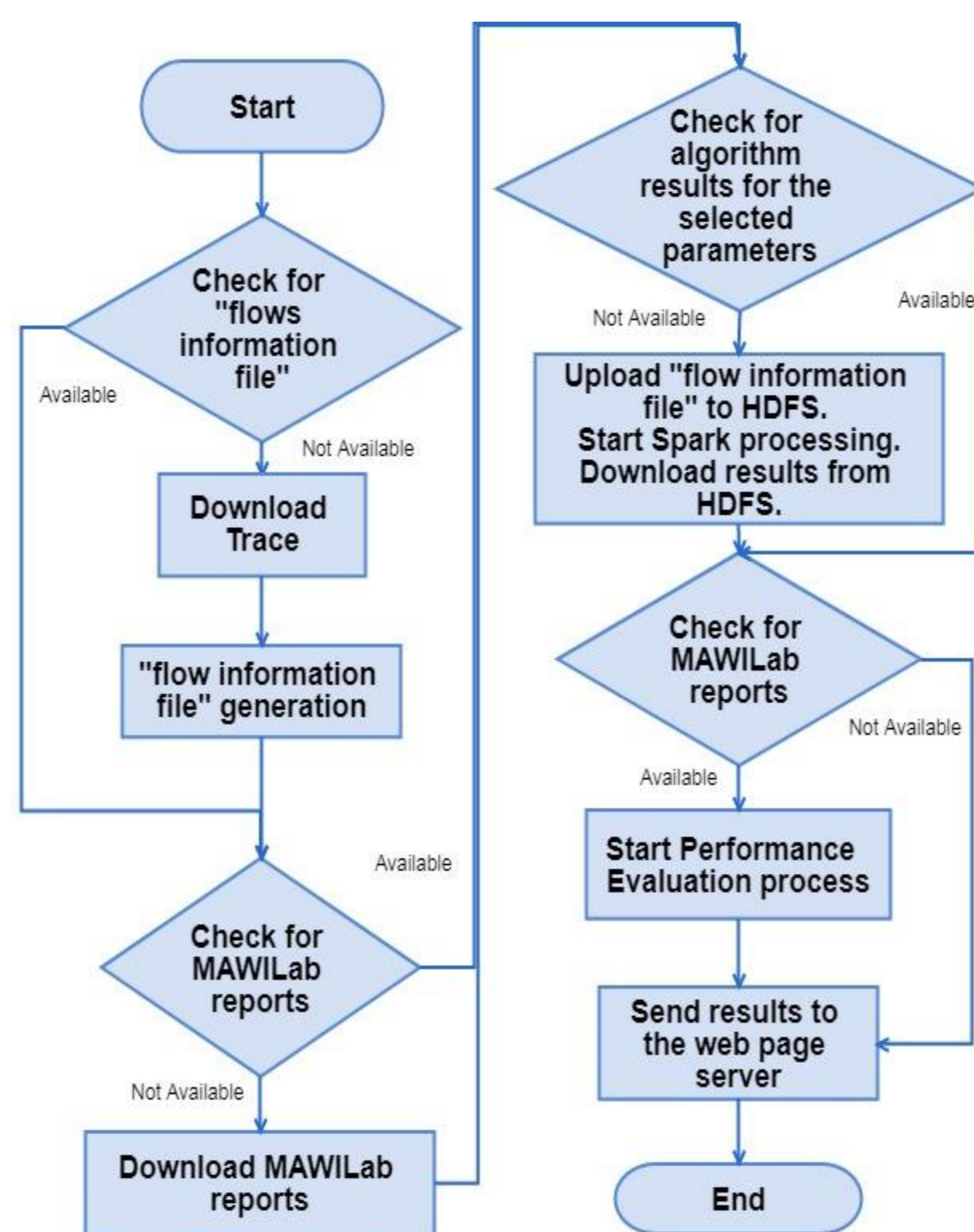
- It downloads the trace to be analyzed
- It generates a flow file using the TIE tool
- It run the algorithm on SPARK
- It publishes the results of the comparison on web portal

Visit us at:

## spada.comics.unina.it

The system allows to
- Detect anomalies on high-speed networks;
- Carry out longitudinal analysis on the anomalies in one of the most used trace reports
- Provide to the scientific community an always-updated reference for comparison
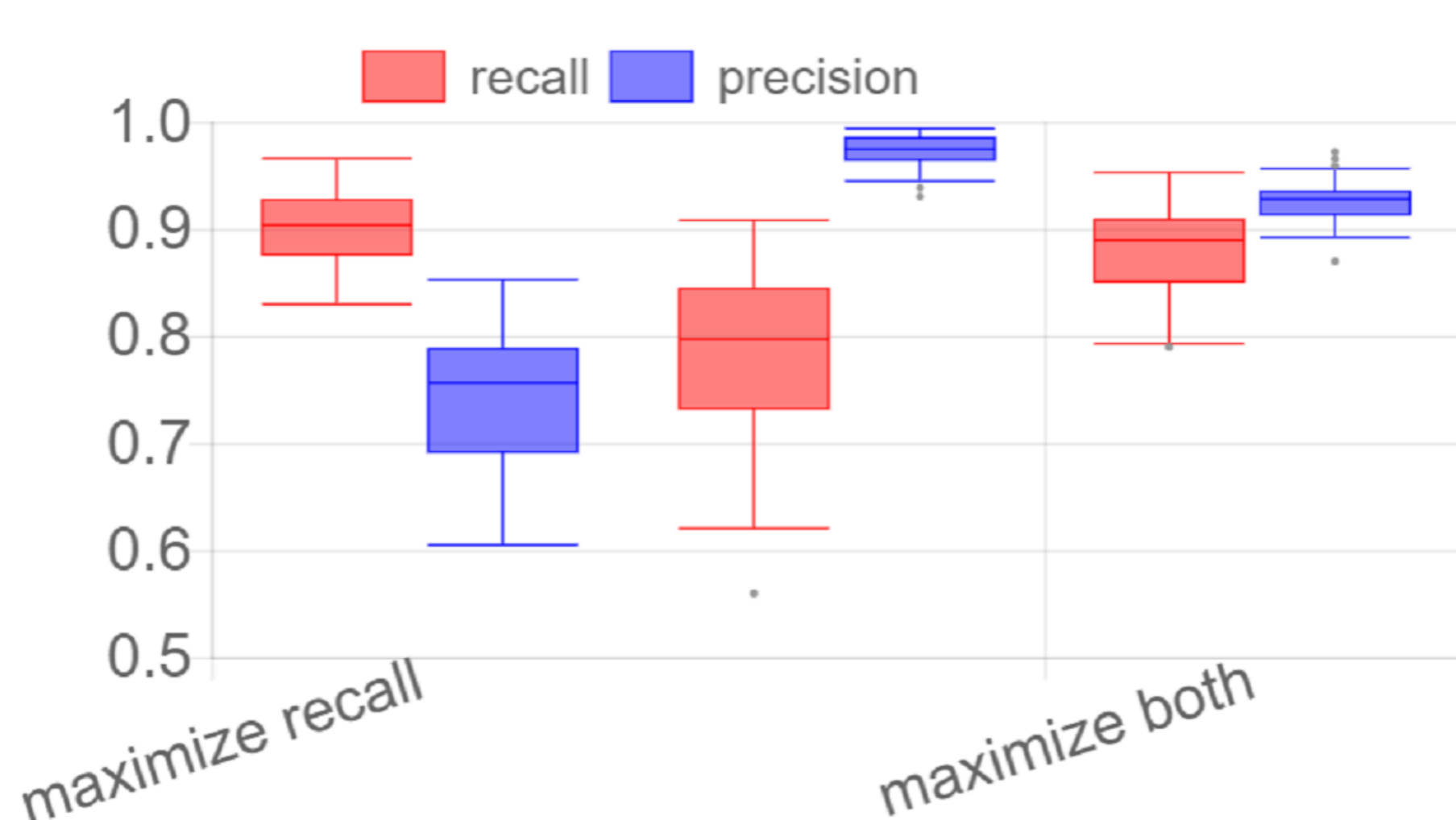


## PRELIMINARY RESULTS
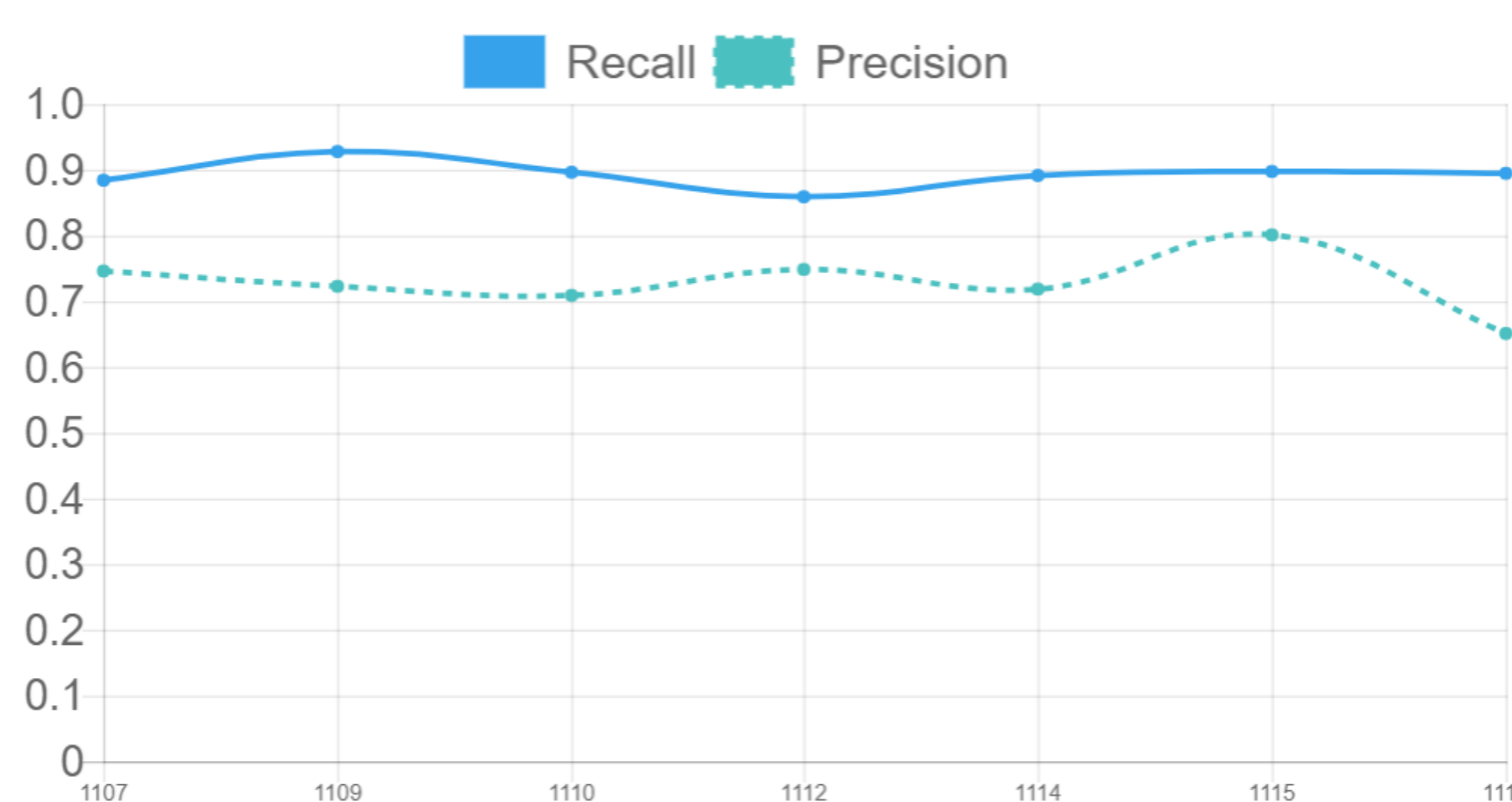

Figure 2: maximizing recall and/or precision


Figure 3: Last 7 days of SPADA

- Figure 2 and 3 show SPADA performance measured on 100 traffic traces with different values of the parameters
- The algorithm is more effective in detecting malicious scanning activity than MAWILab
- It allows to obtain a new ground truth starting from MAWILab

## Acknowledgement

## References

[1] SPADA documentation- SPark-based Anomaly Detection Ace. http://spad7a.comics.unina.it/
[2] Big data analytics on apache spark- https://link.springer.com/content/pdf/10.1007%2Fs41060-016-0027-9.pdf
[3] MAWILab – Home- http://www.fukuda-lab.org/mawilab/
[4] P. Borgnat, G. Dewaele, K. Fukuda, P. Abry, and K. Cho. Seven years and one day: Sketching the evolution of internet traffic. In IEEE INFOCOM April 2009
[5] Casas