

Analysis of Malware Dataset

Brad Miller

Outline

- Context & Data Overview
- Infrastructure
- Analysis

CONTEXT & DATA OVERVIEW

In an Ideal World...

- An evaluation dataset would include
 - Full analysis of every file that ever appears
 - Past, Present & Future!
 - The prevalence or importance of each file
 - Gold standard label given by magical oracle

Life in an Imperfect World

- Access to data in academia is limited
- The challenge is to produce solid results despite imperfections in the data
- **Our Goals: Examine the viability of a malware dataset for evaluating malware detection**

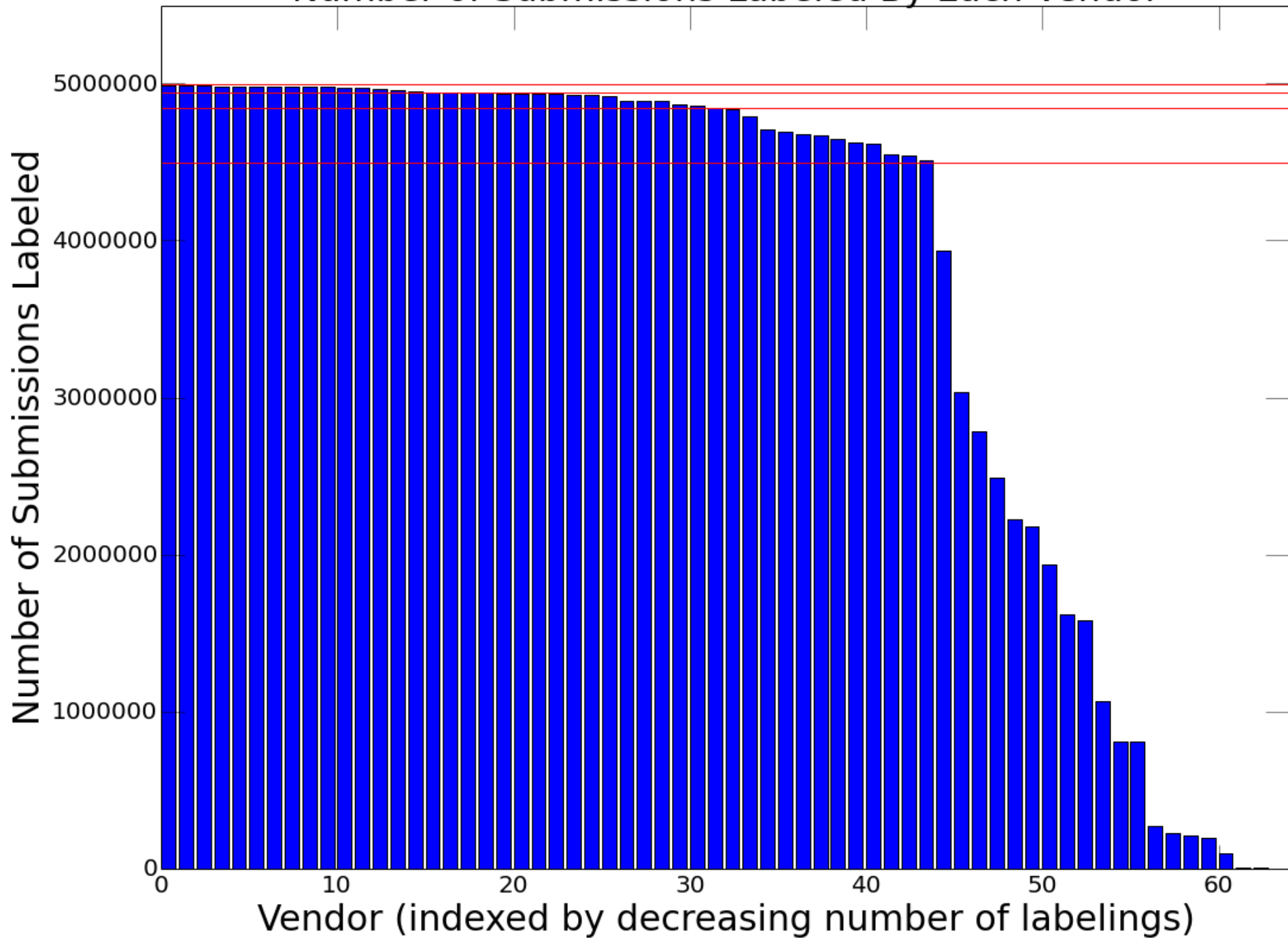
Data Origins

- Hashes provided by VirusTotal & AV Provider
- VirusTotal provides
 - Labeling from many AV vendors
 - Times when each hash is submitted to VirusTotal

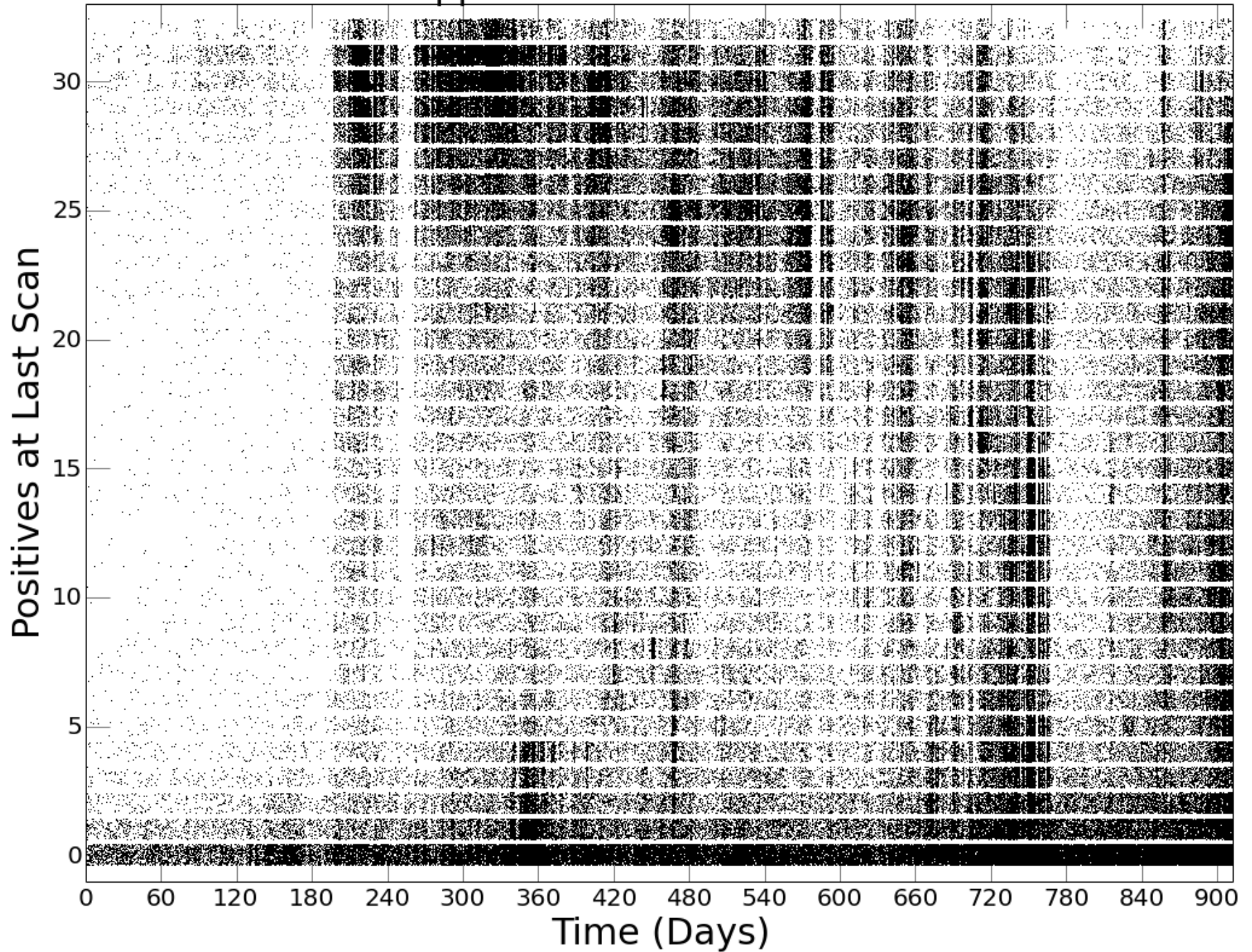
Data Overview

- 1.1M hashes with dynamic data available
- 400K hashes with multiple submissions
 - We will focus on these today
 - Multiple submissions allows us to observe AV vendor behavior over time

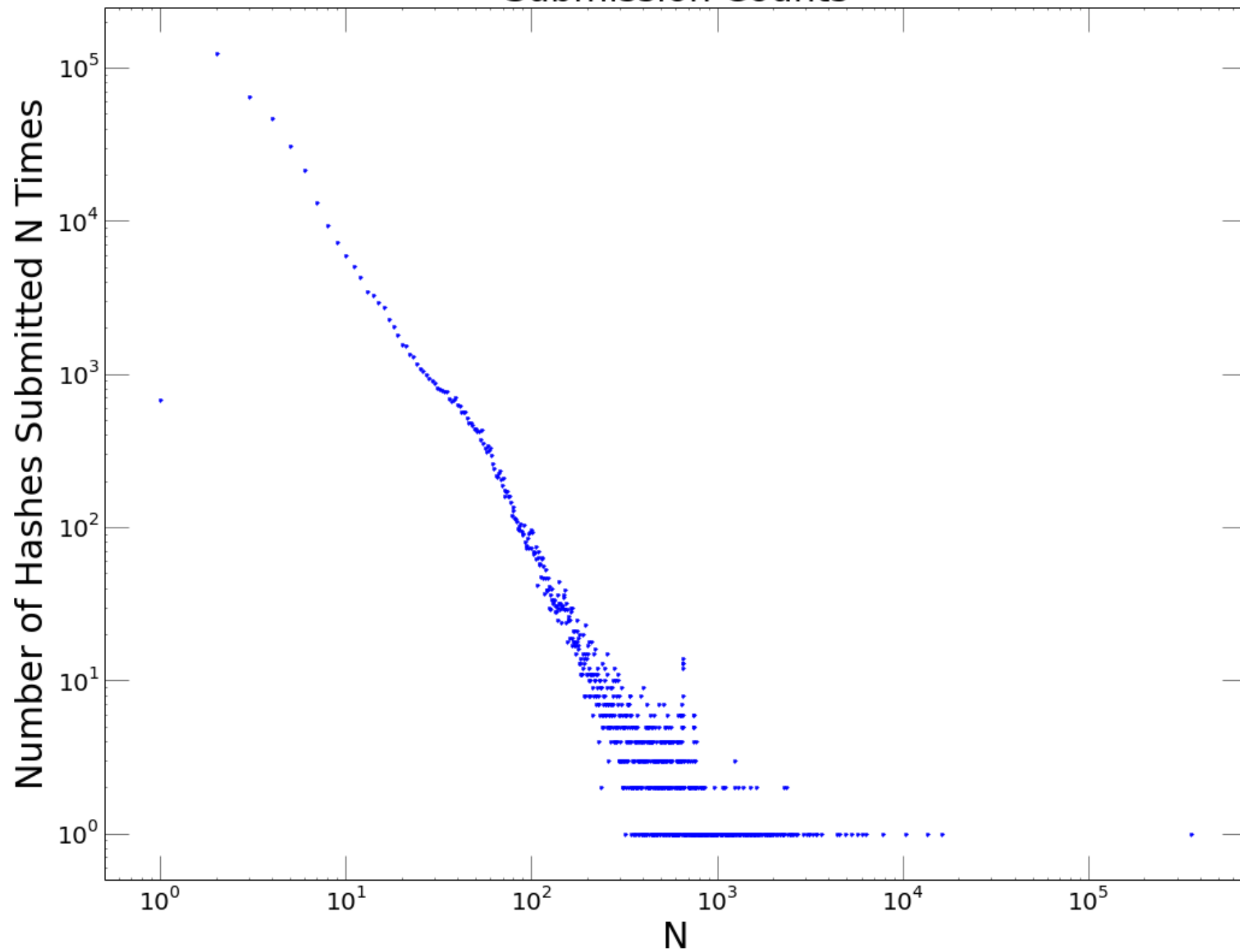
Number of Submissions Labeled By Each Vendor



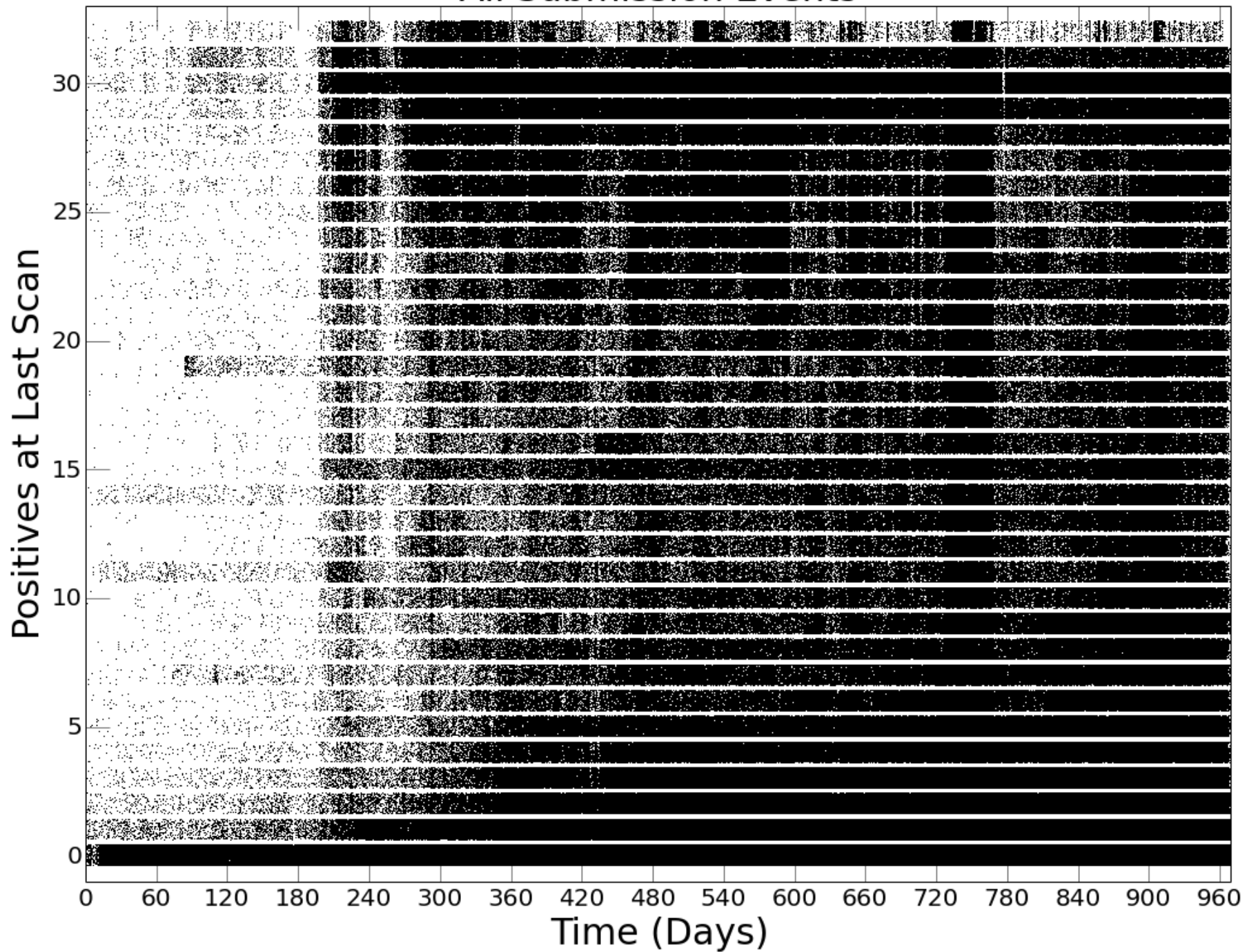
Appearance of New Hashes



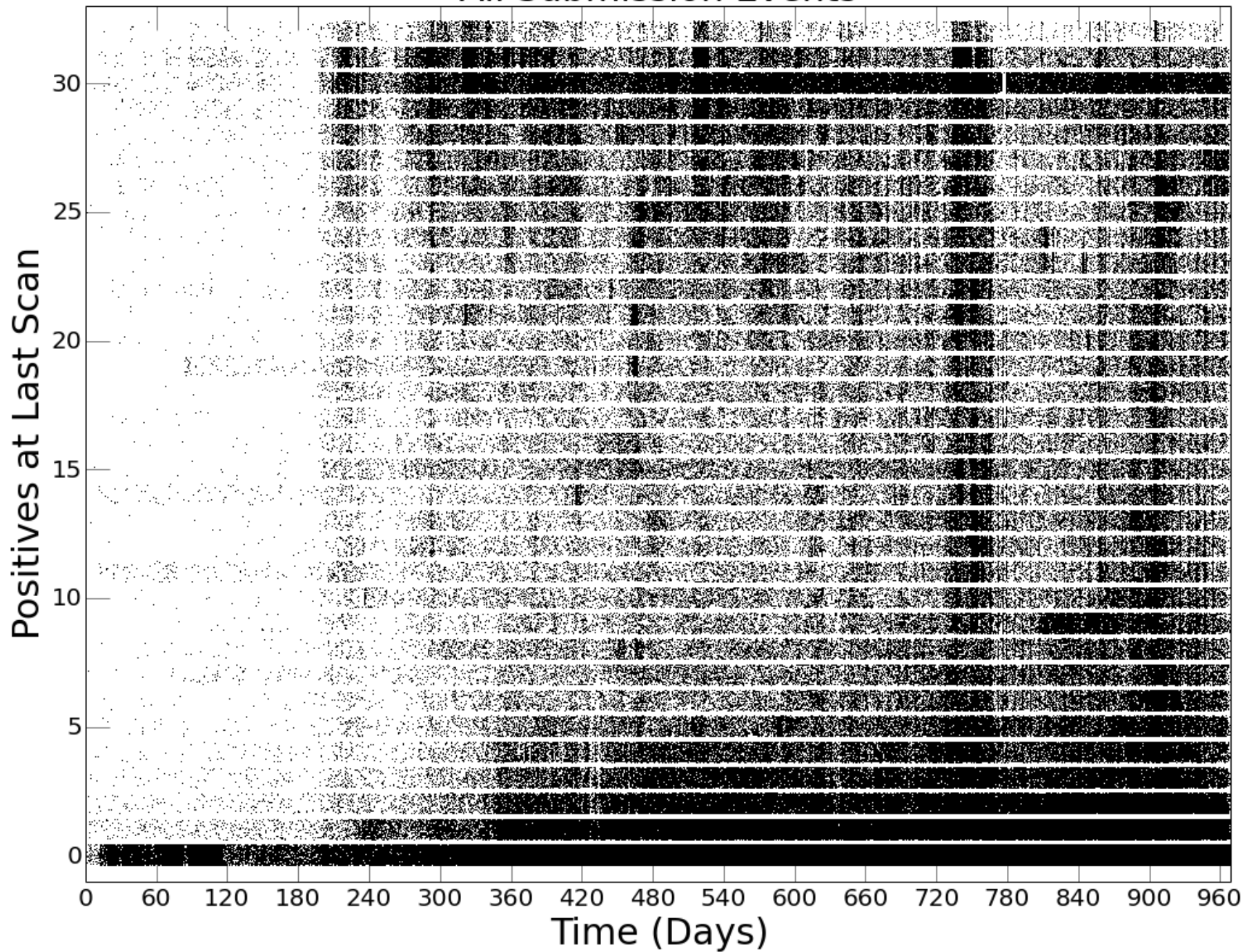
Submission Counts



All Submission Events



All Submission Events



INFRASTRUCTURE

Analysis Stack

matplotlib: Data Visualization

iPython Notebook:
Interactive Computation

Spark: Distributed (fast) Computation

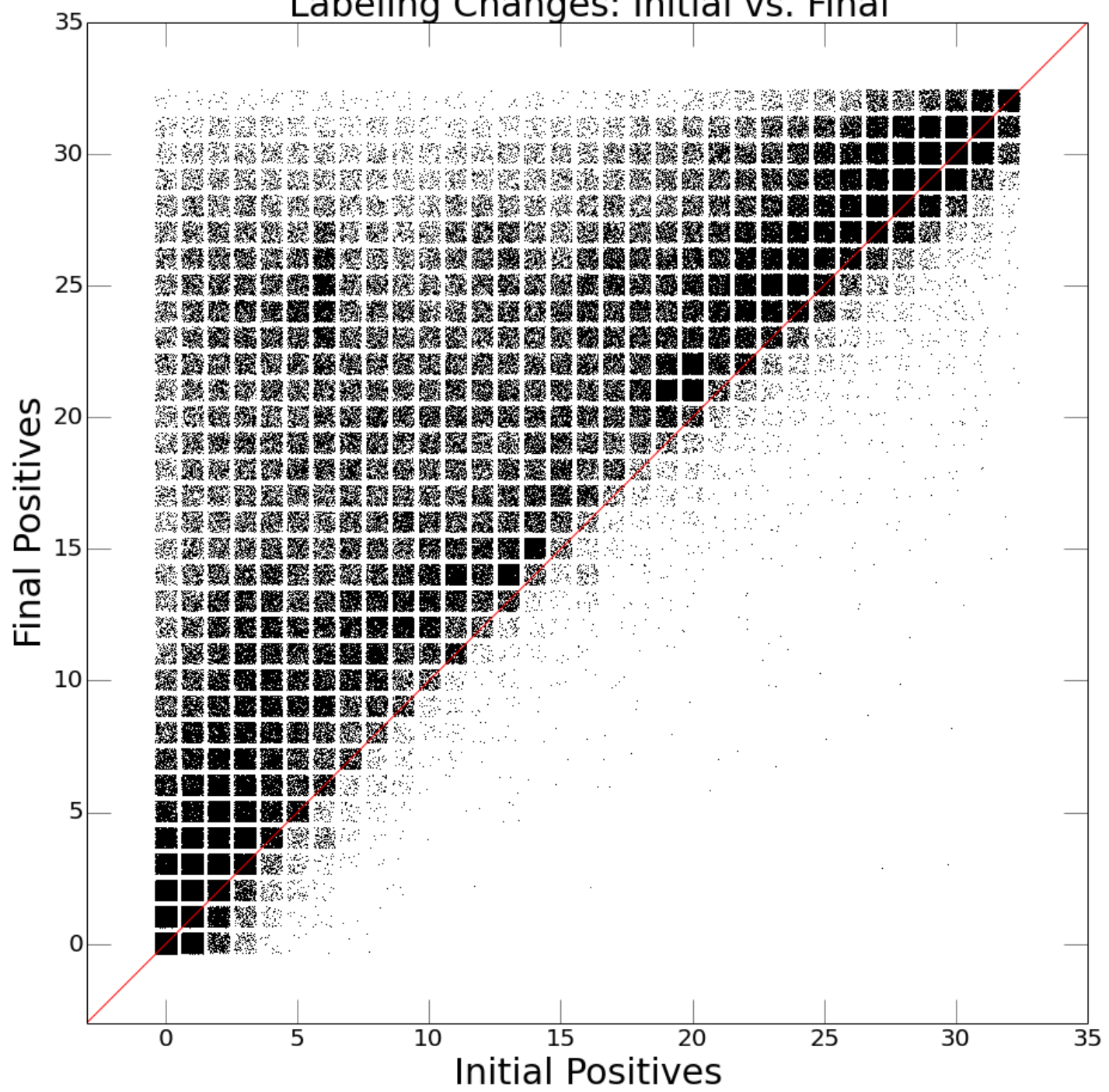
HDFS: Distributed Storage

ANALYSIS

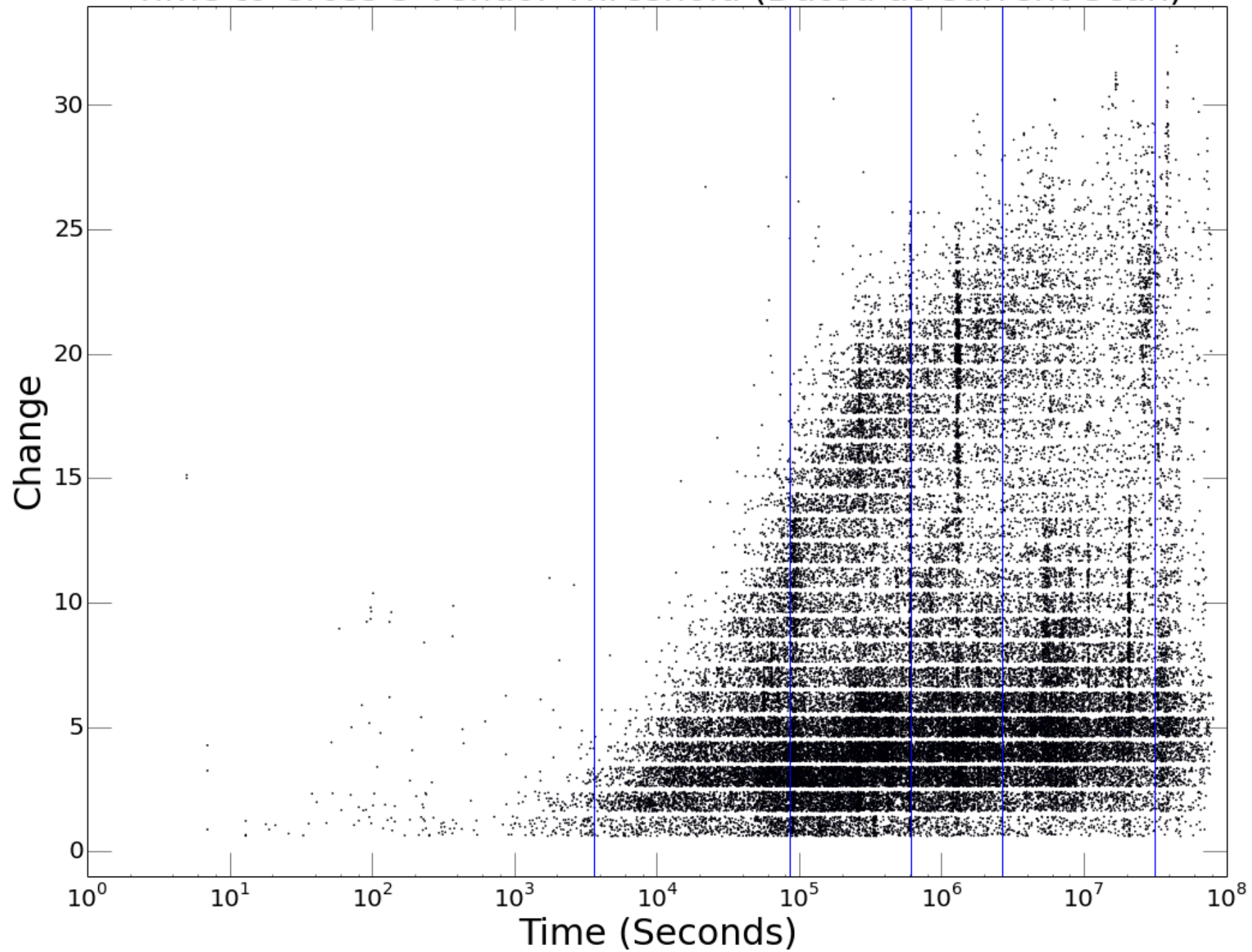
Investigative Targets

- Is this a good dataset to use for evaluating a malware detection system?
 - Keep in mind the properties of the ideal dataset
- How should we label samples?
 - Vendors may change their minds over time
 - Some vendors may be better than others

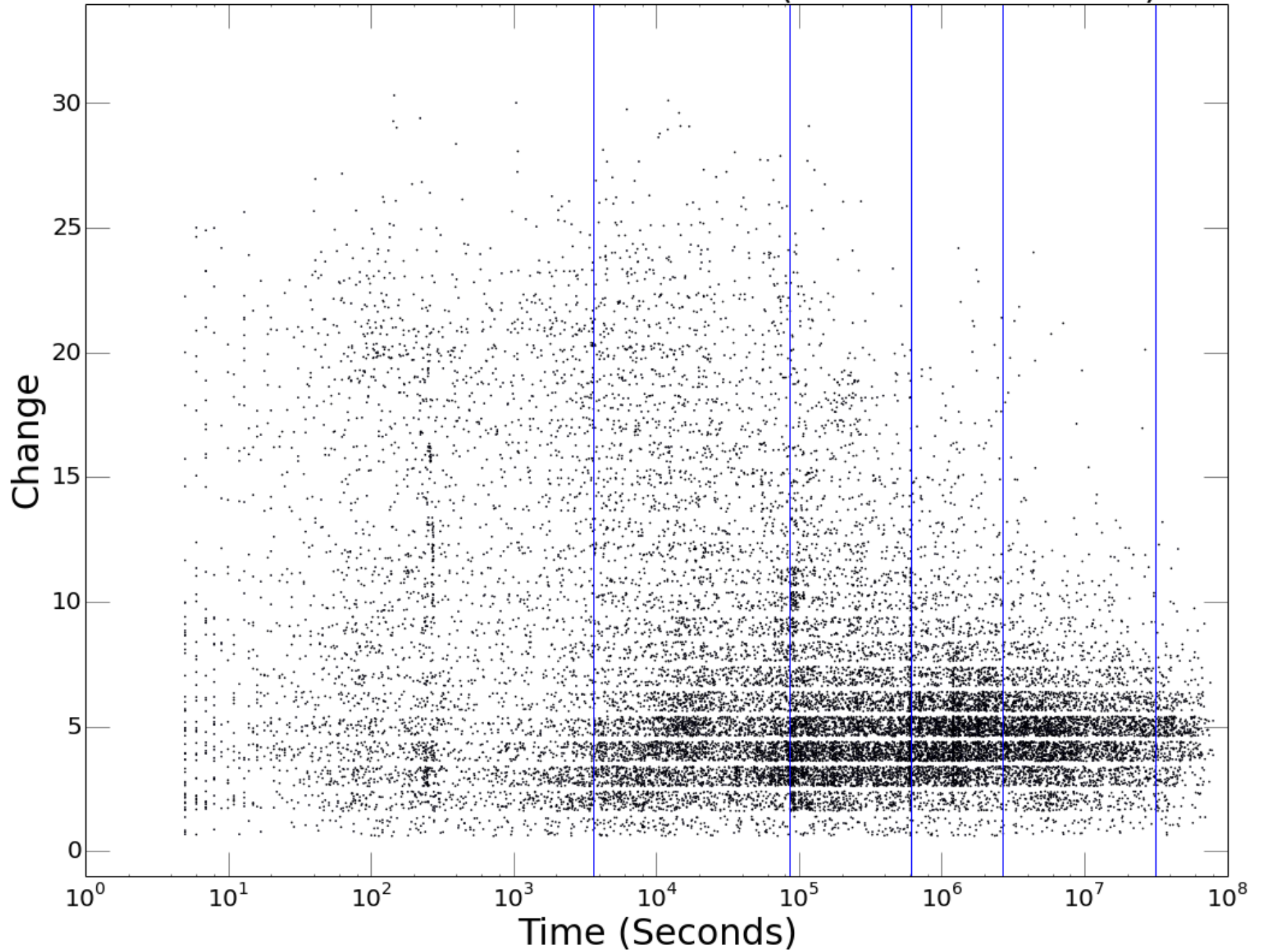
Labeling Changes: Initial vs. Final



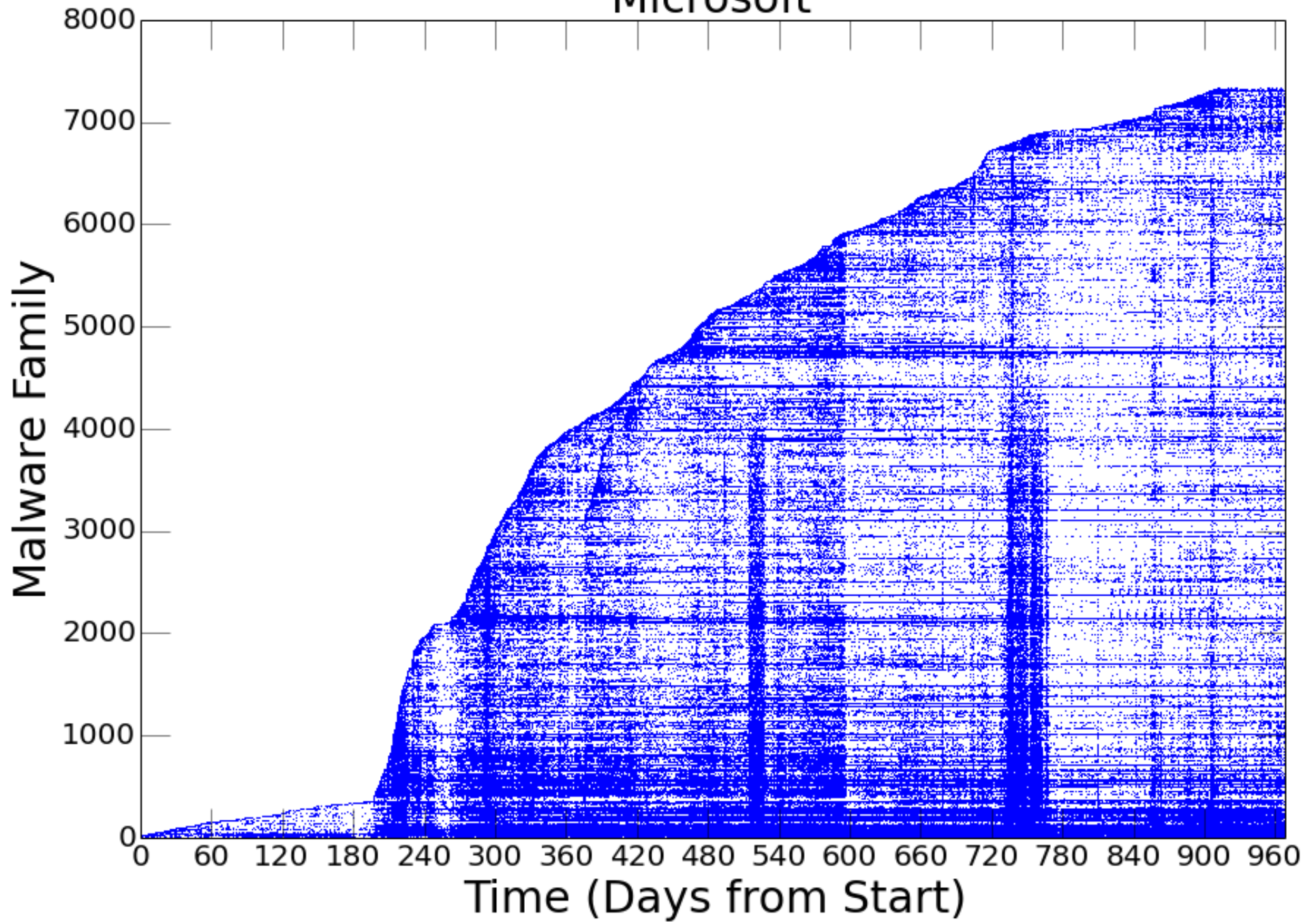
Time to Cross 5 Vendor Threshold (Dated at Current Scan)



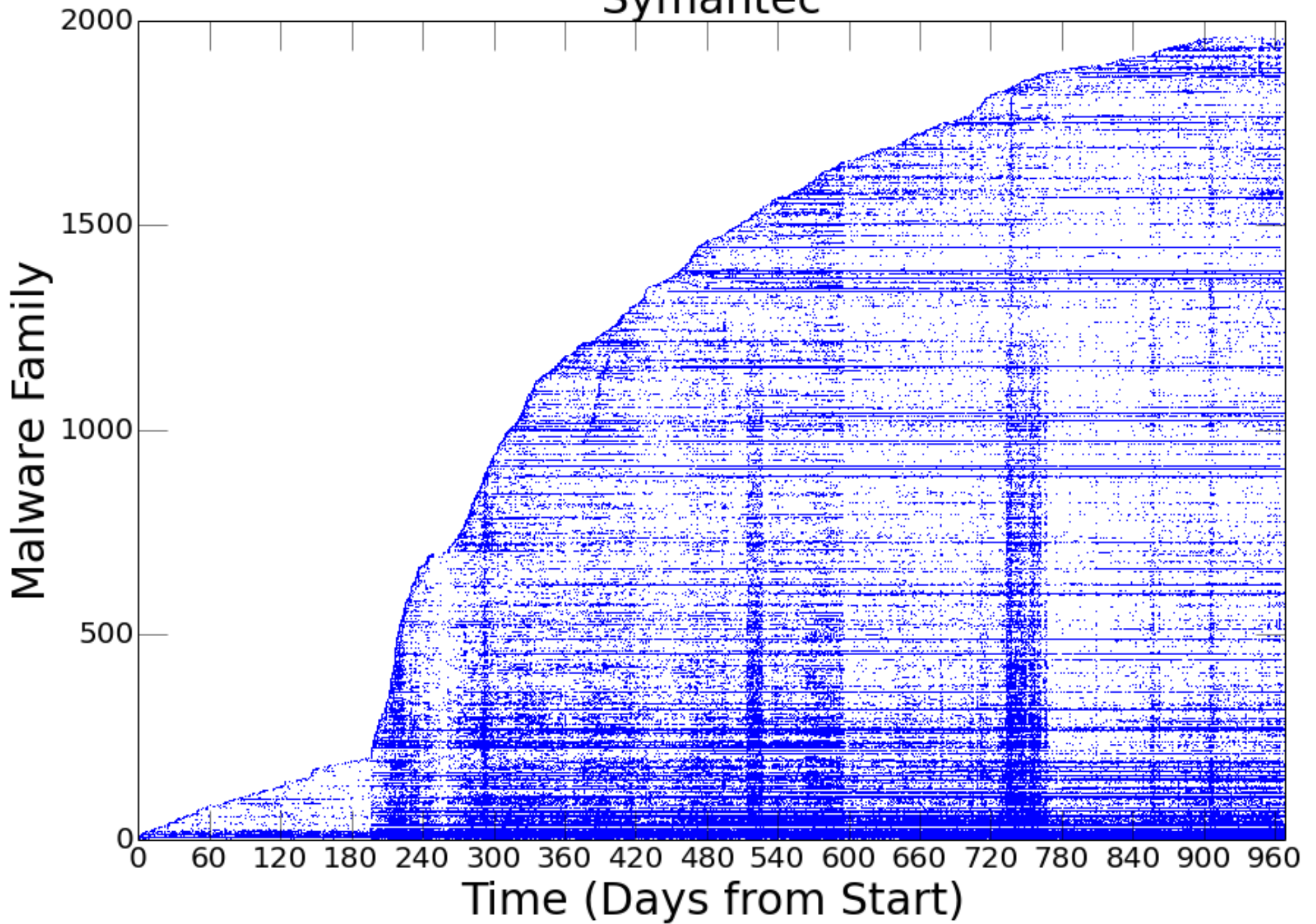
Time to Cross 5 Vendor Threshold (Dated at Prior Scan)



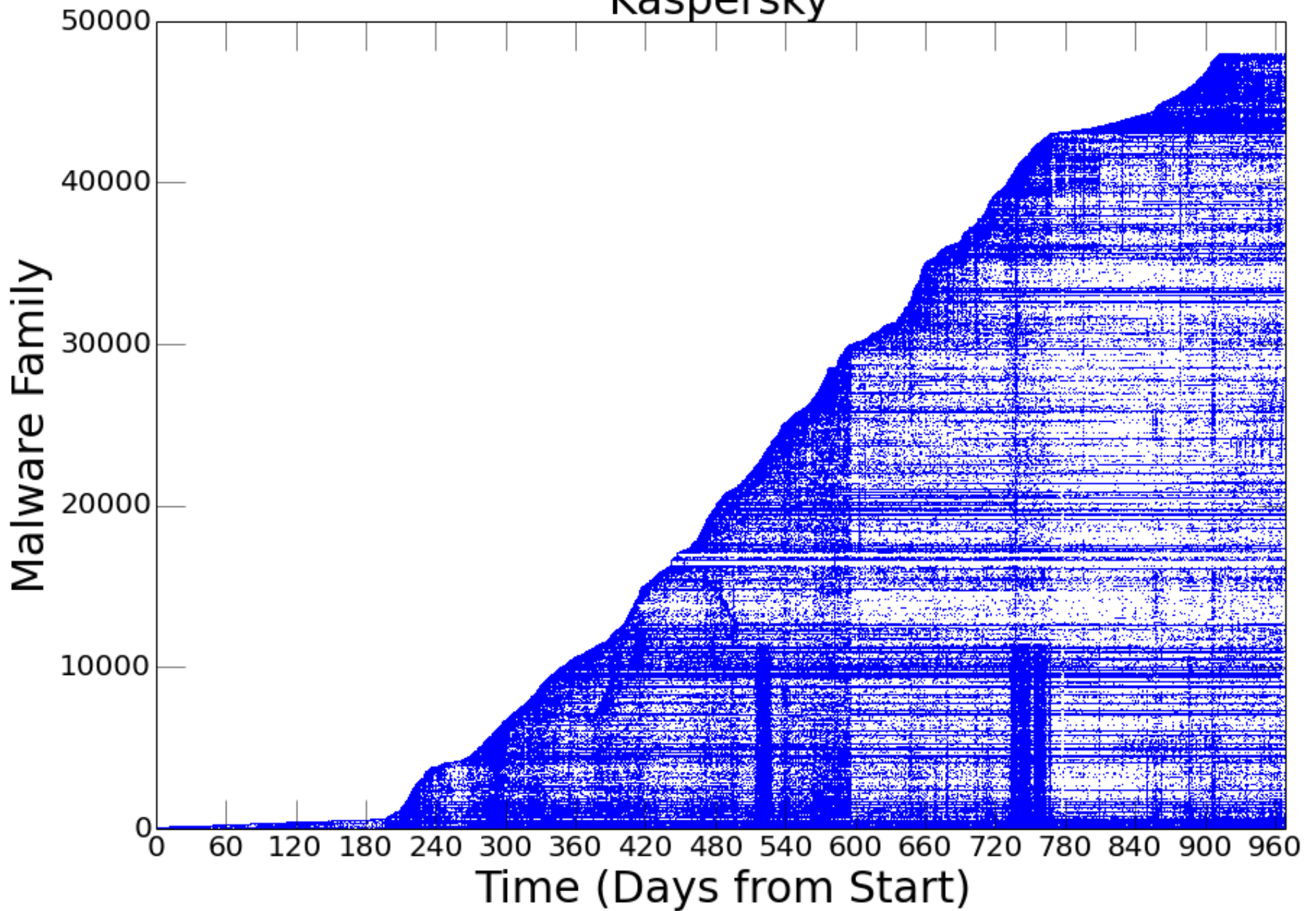
Microsoft

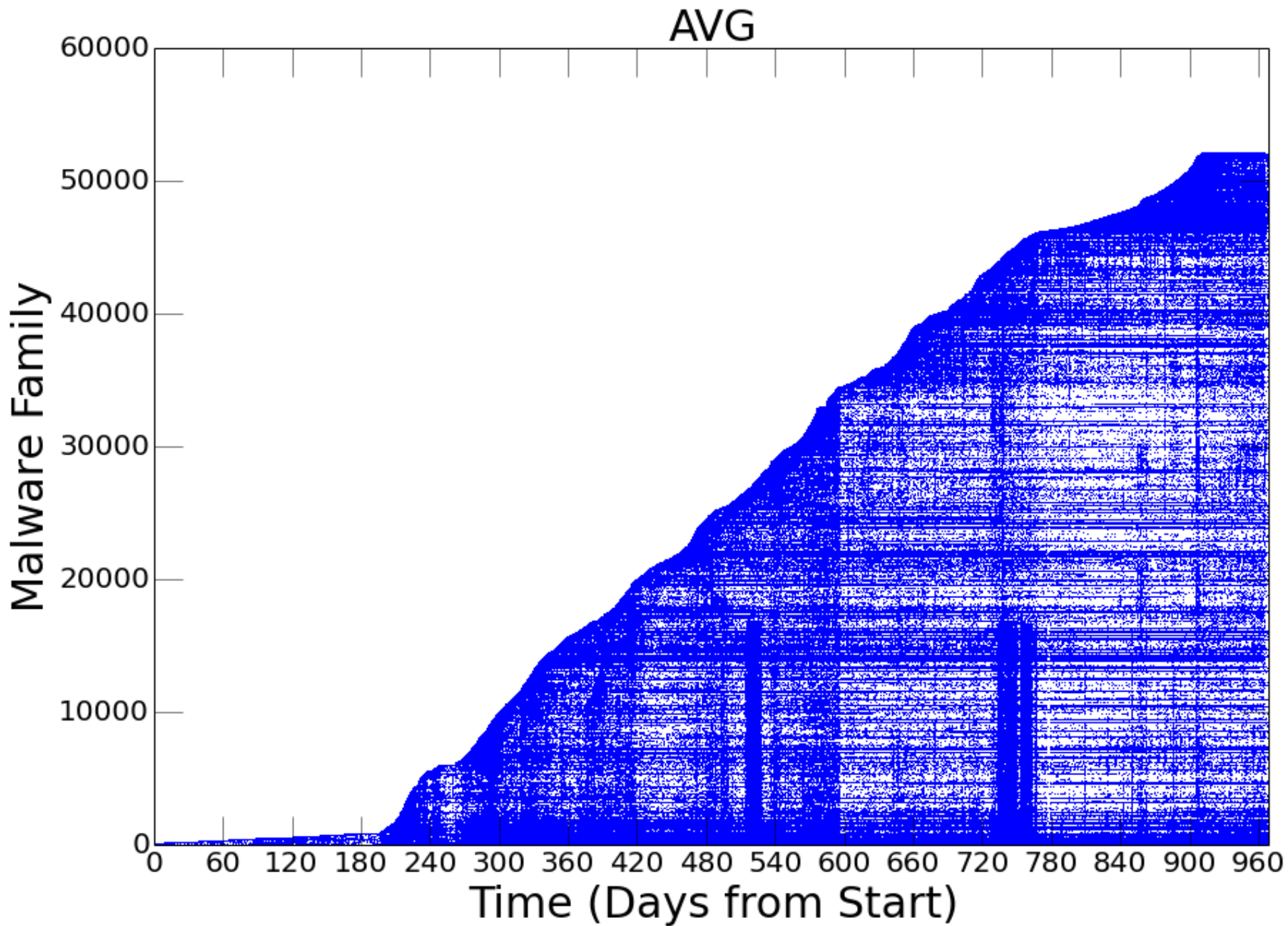


Symantec

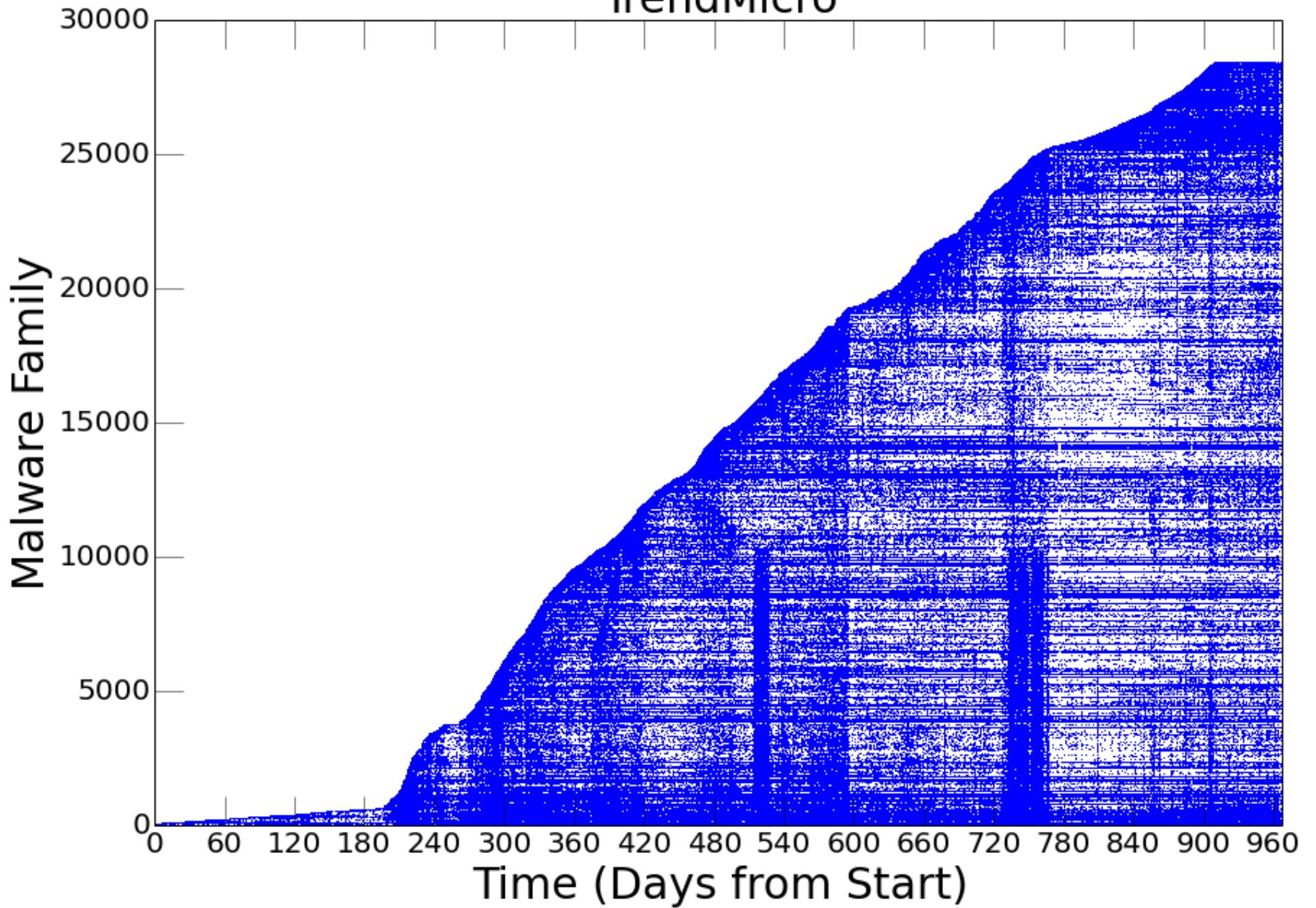


Kaspersky

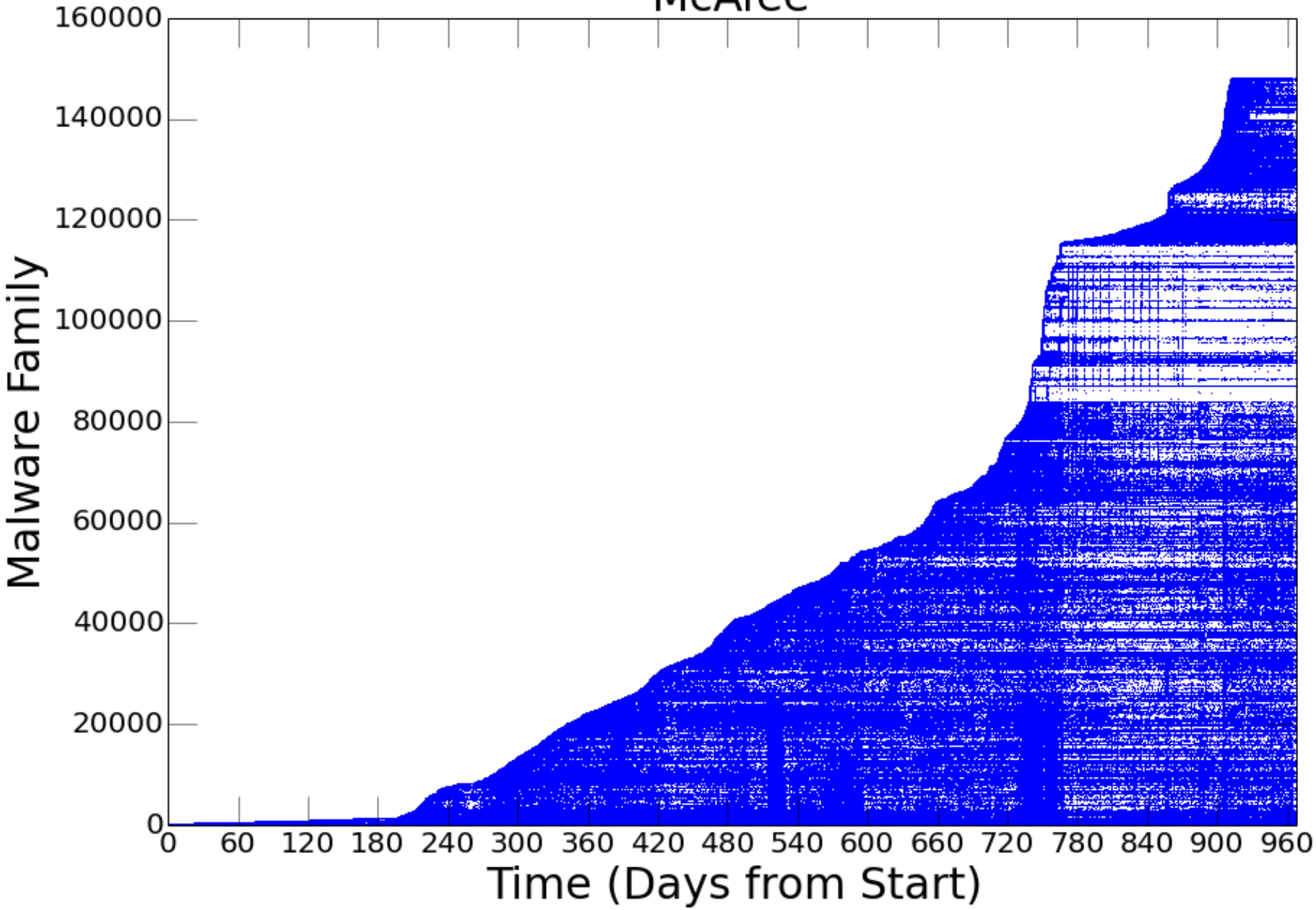




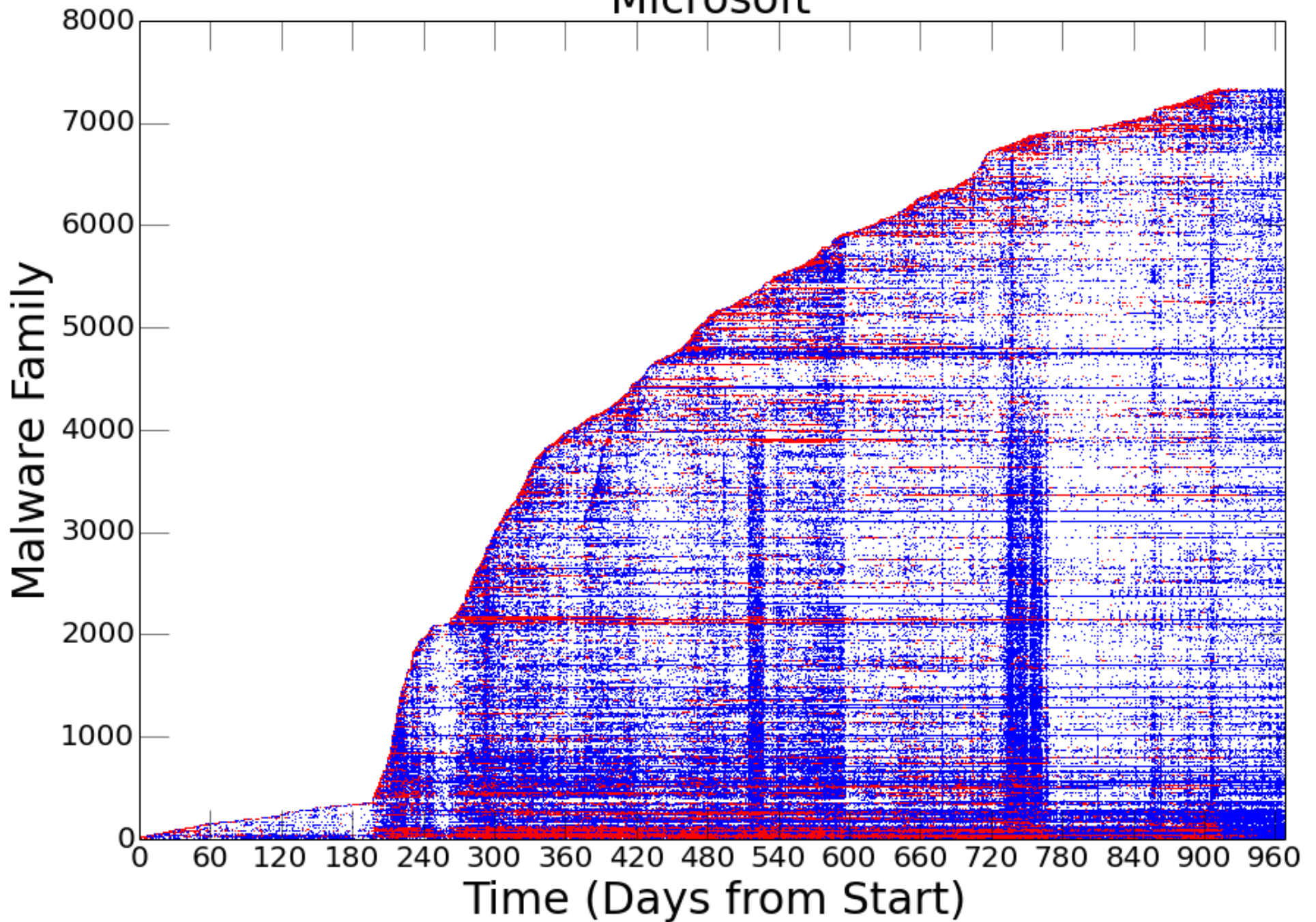
TrendMicro



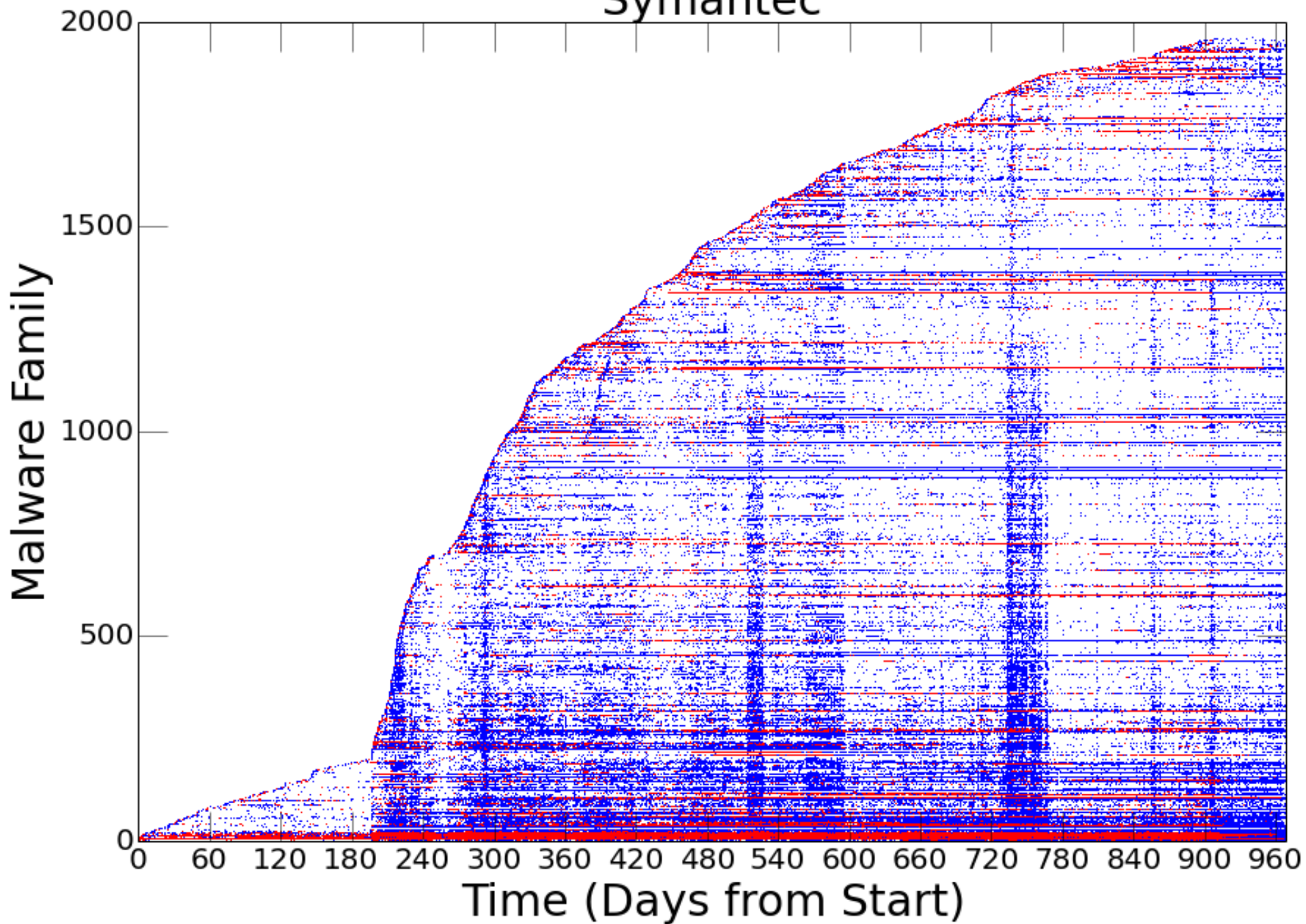
McAfee



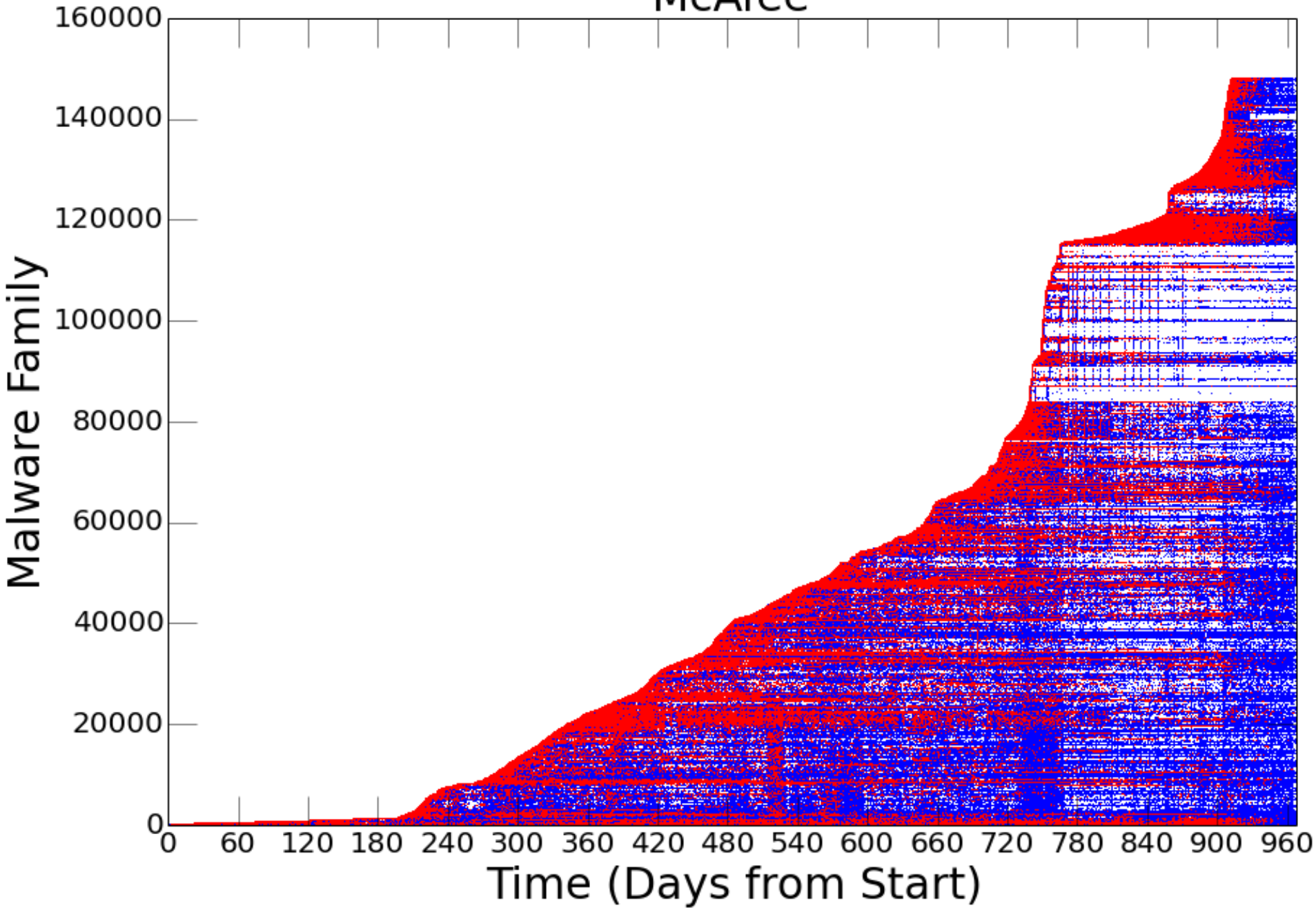
Microsoft



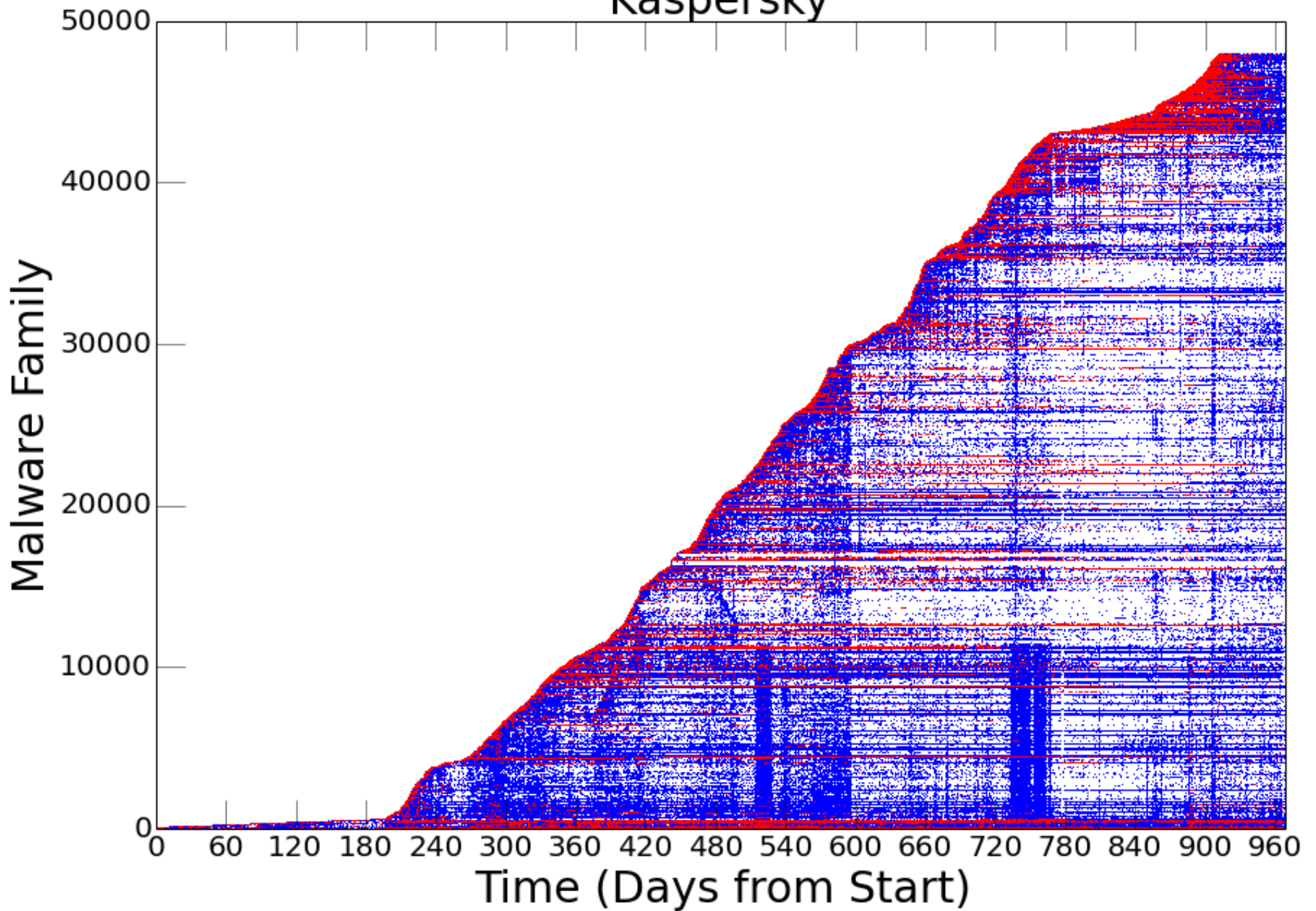
Symantec



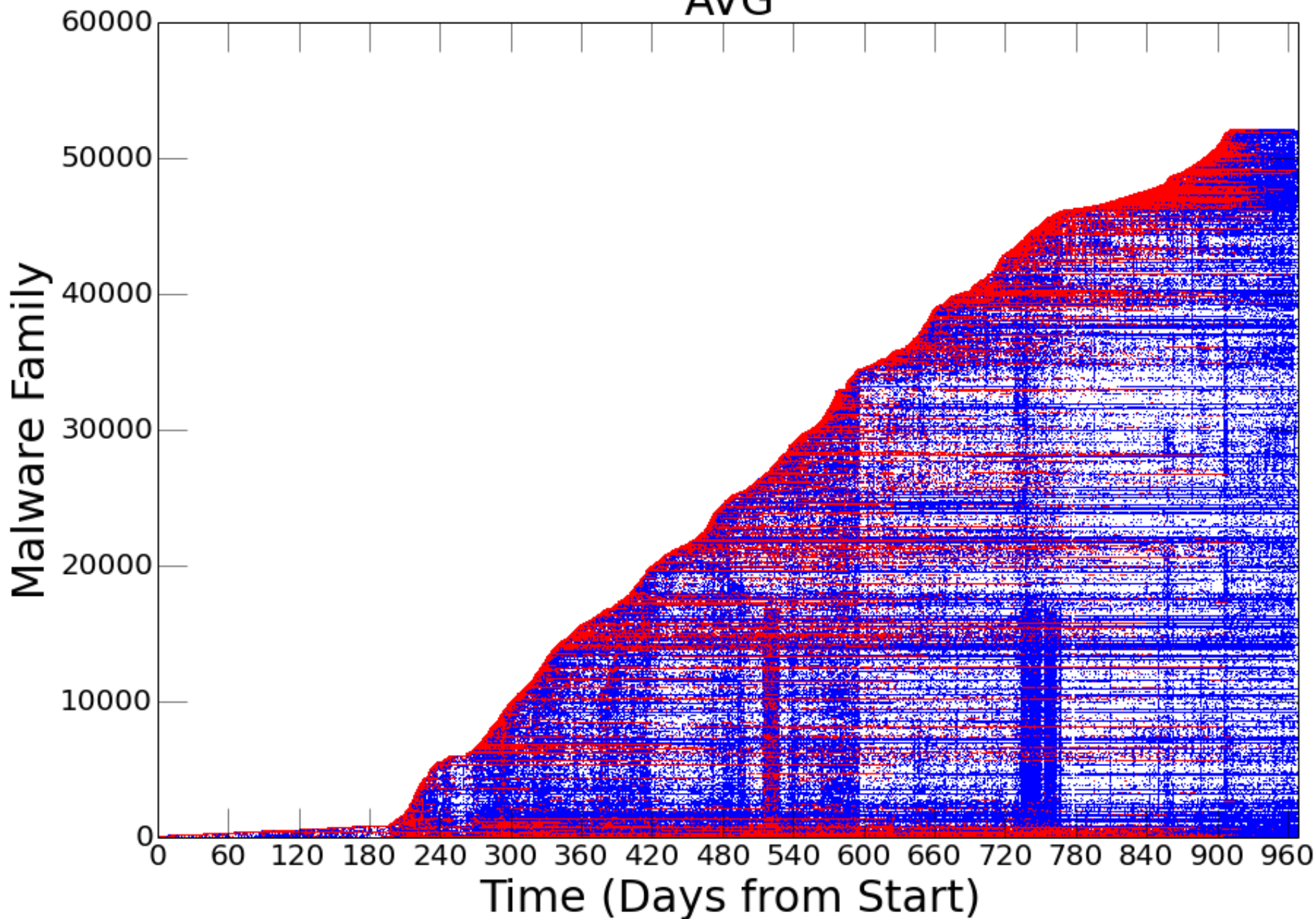
McAfee



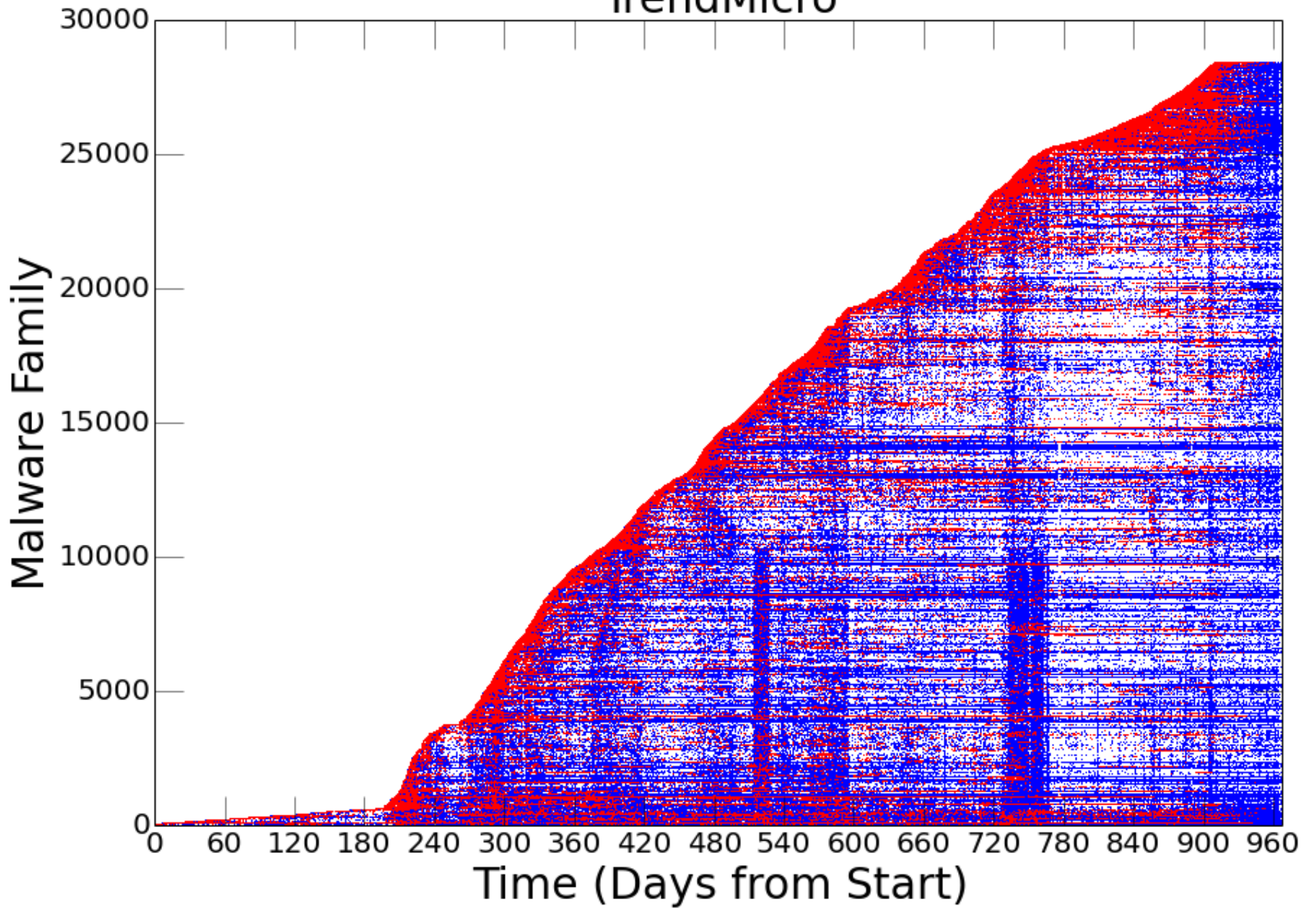
Kaspersky



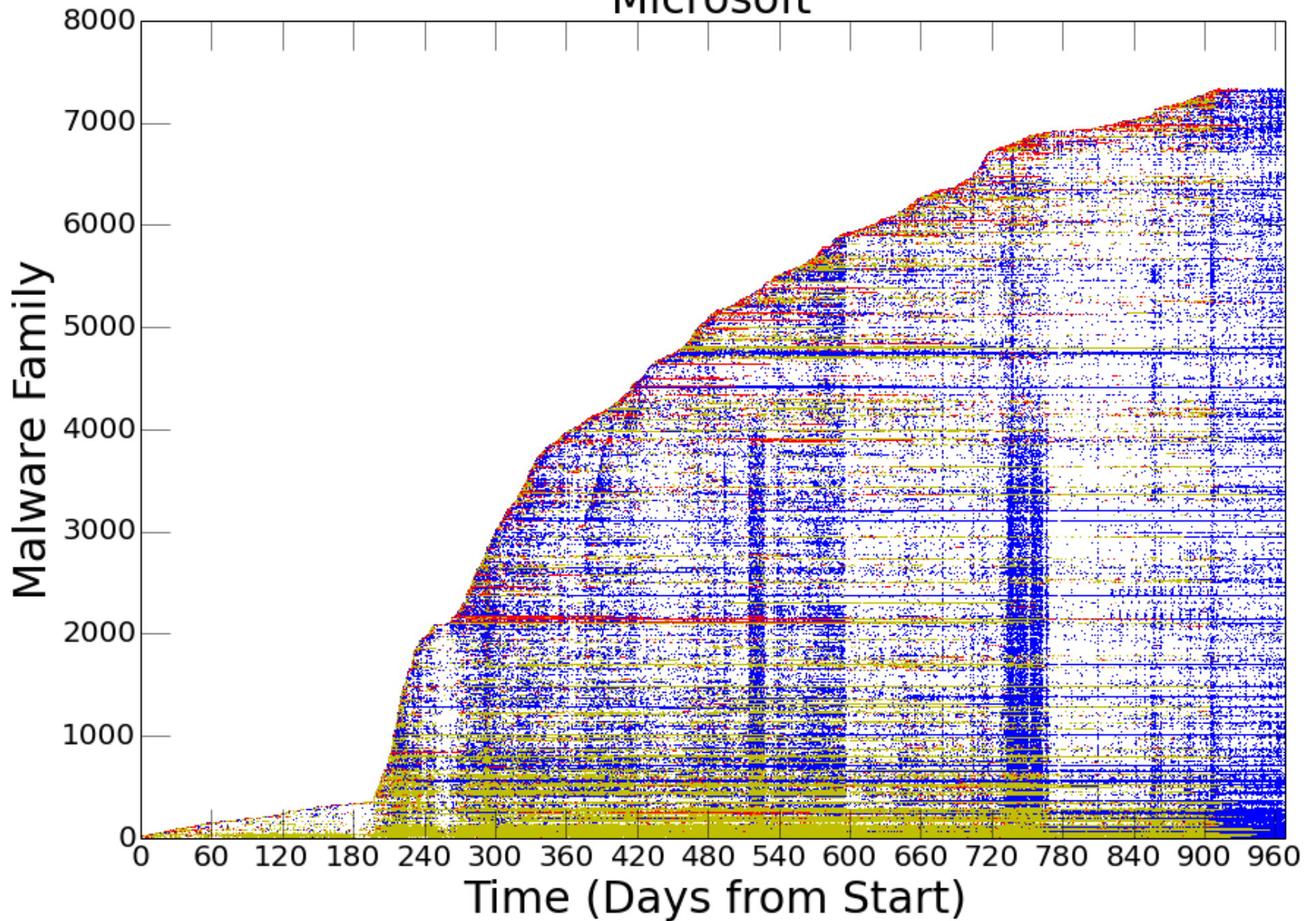
AVG



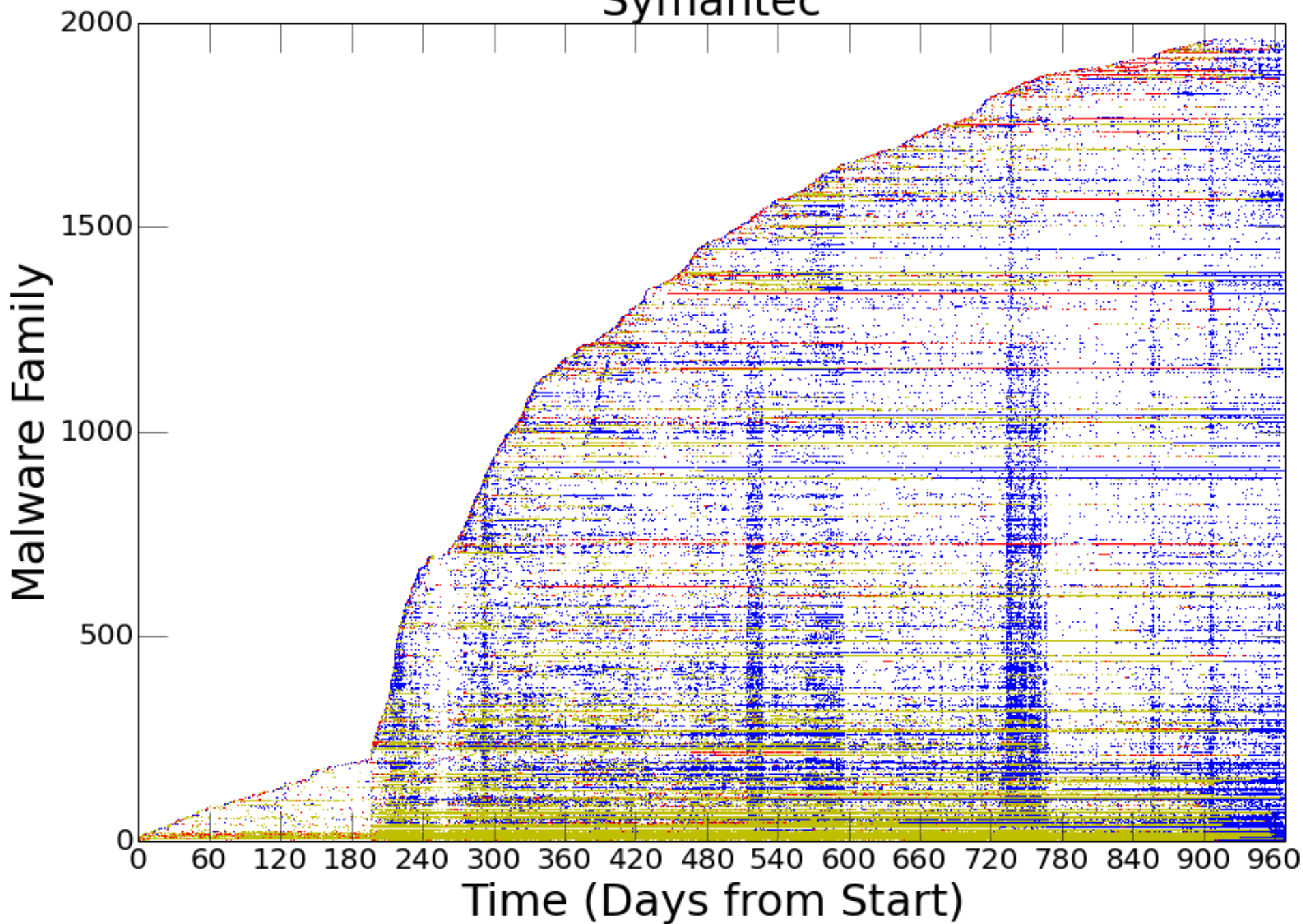
TrendMicro



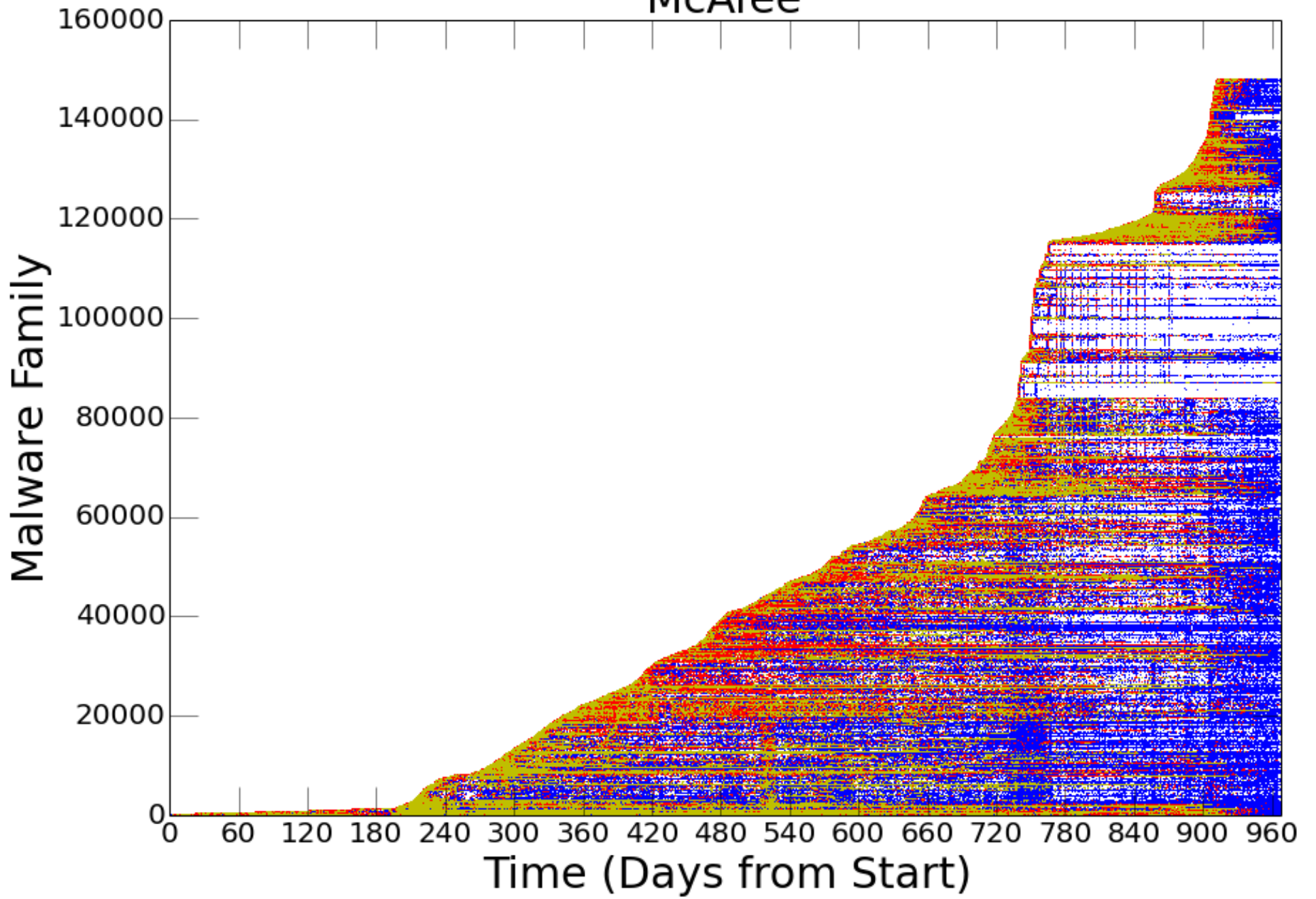
Microsoft



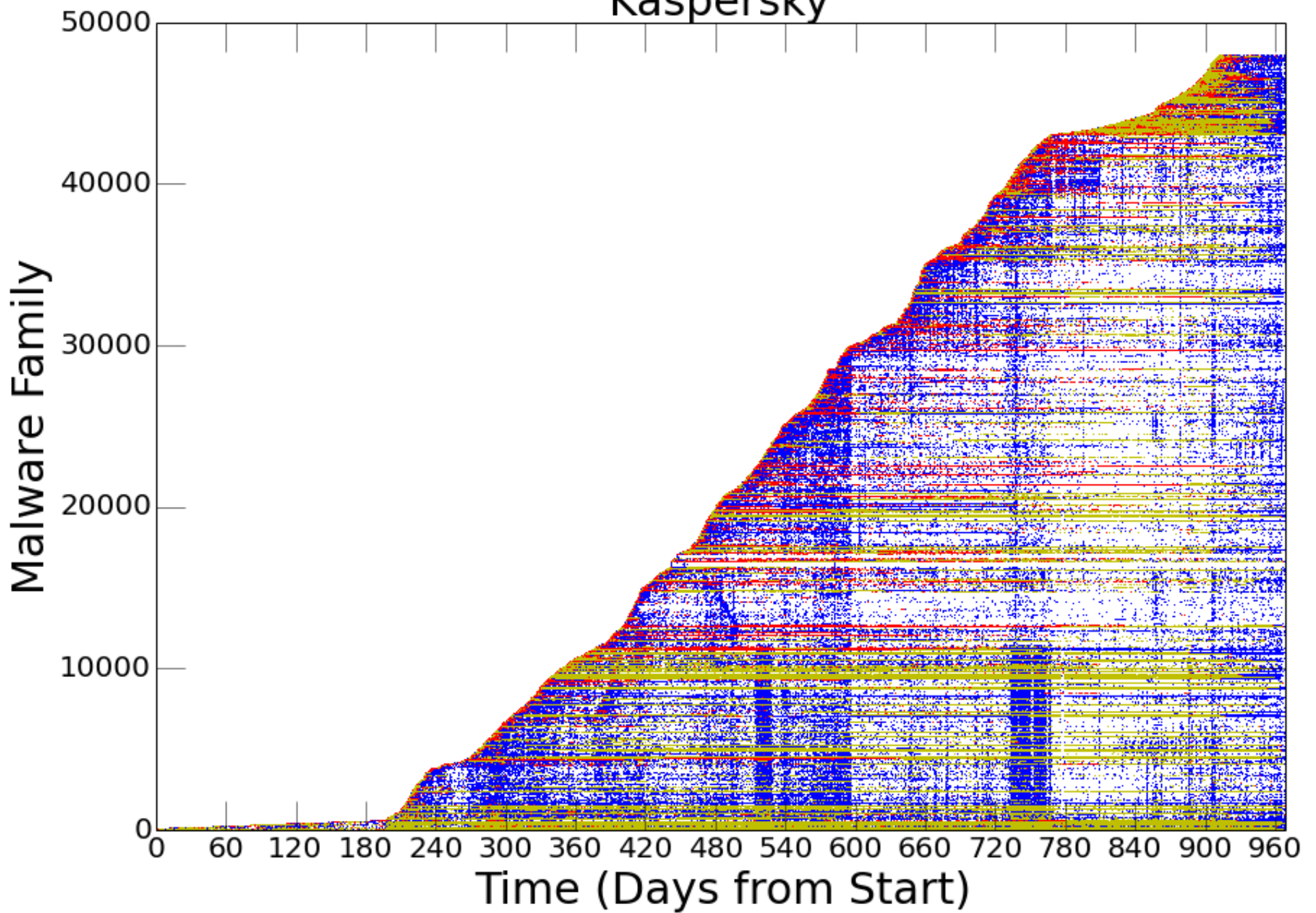
Symantec

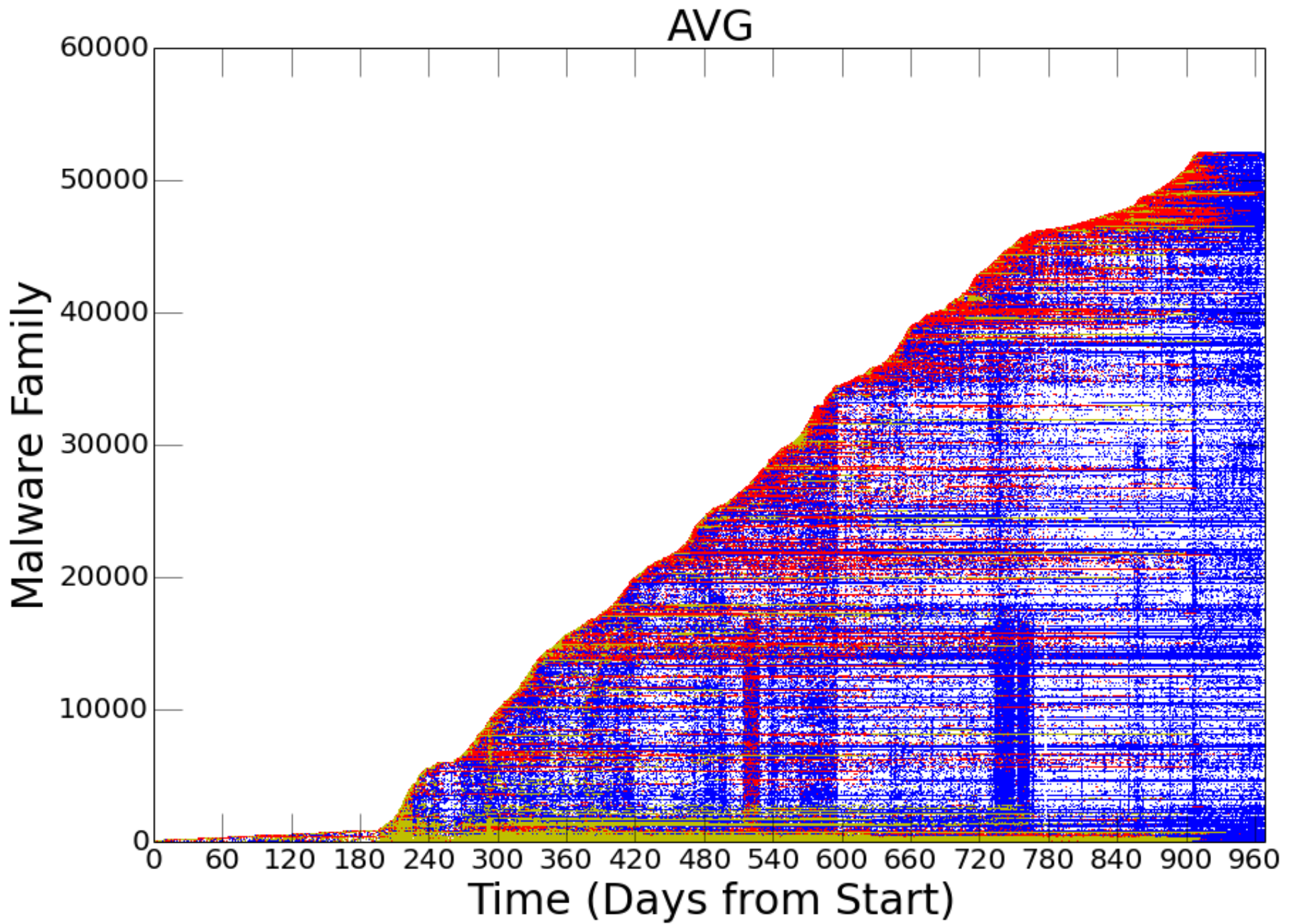


McAfee



Kaspersky





TrendMicro

