# Malware Analysis at Scale with ElasticSearch

## a.k.a playing around with a whole lot of malware

## Team Members

Taylor McCampbell
- Junior, COSC Major
- INL Intern

Rafer Cooley
- COSC Ph.D. Student
- INL Intern

## Background

- Amount of malware increasing at exponential pace
- Modern indicators are easily bypassed by changing minor behaviors in malware programs
- Big data problems have generated Big Data solutions
- Search engines have been developed to efficiently search through large streams of information
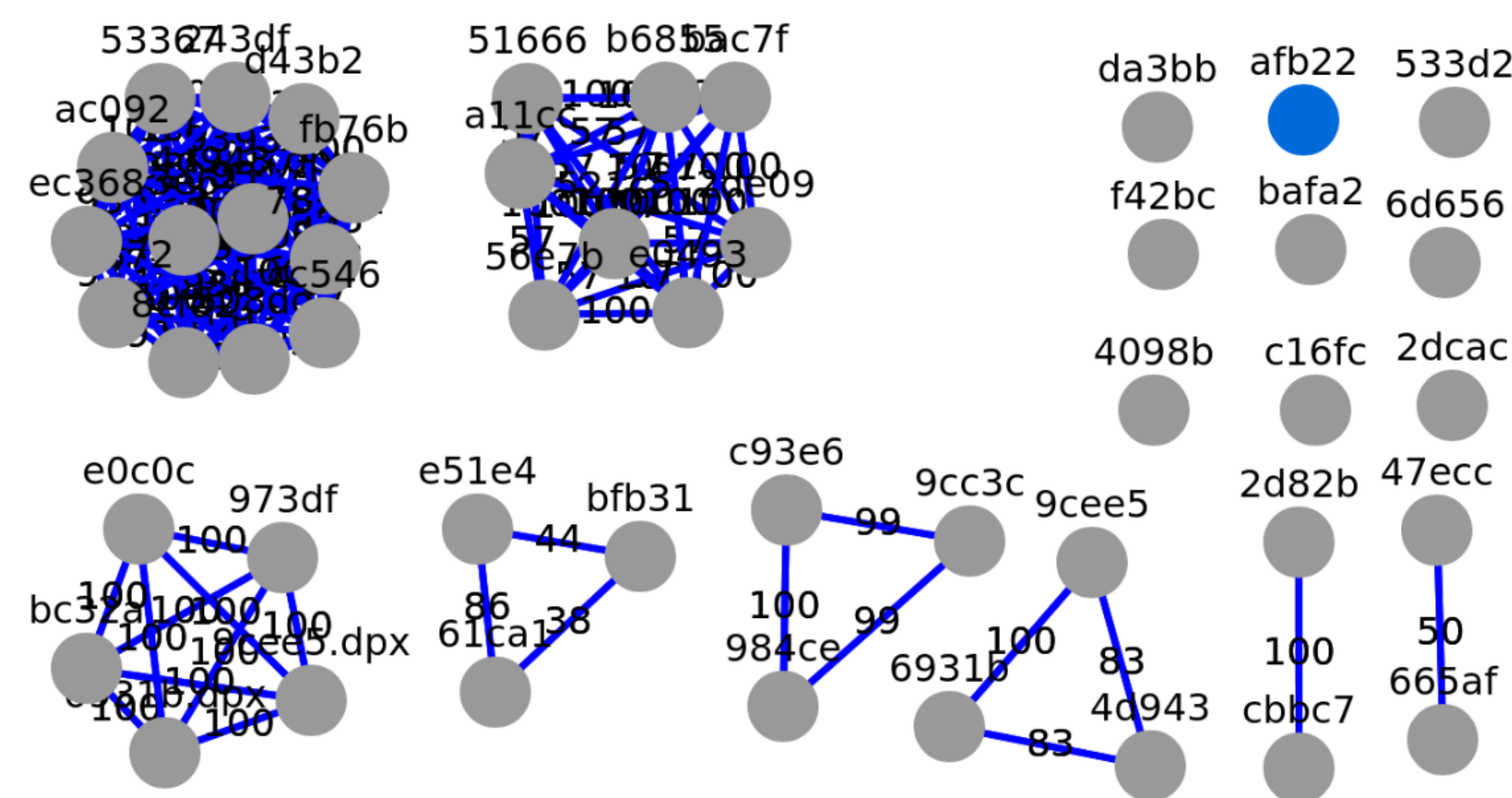
## Problem Statement

With the exponential increase of malware being released everyday, an efficient and resilient scheme for malware indicators is needed. In order to store, analyze, and develop future indicators an efficient data management system is needed that can handle the amount of information necessary to accurately test said indicator schemes.
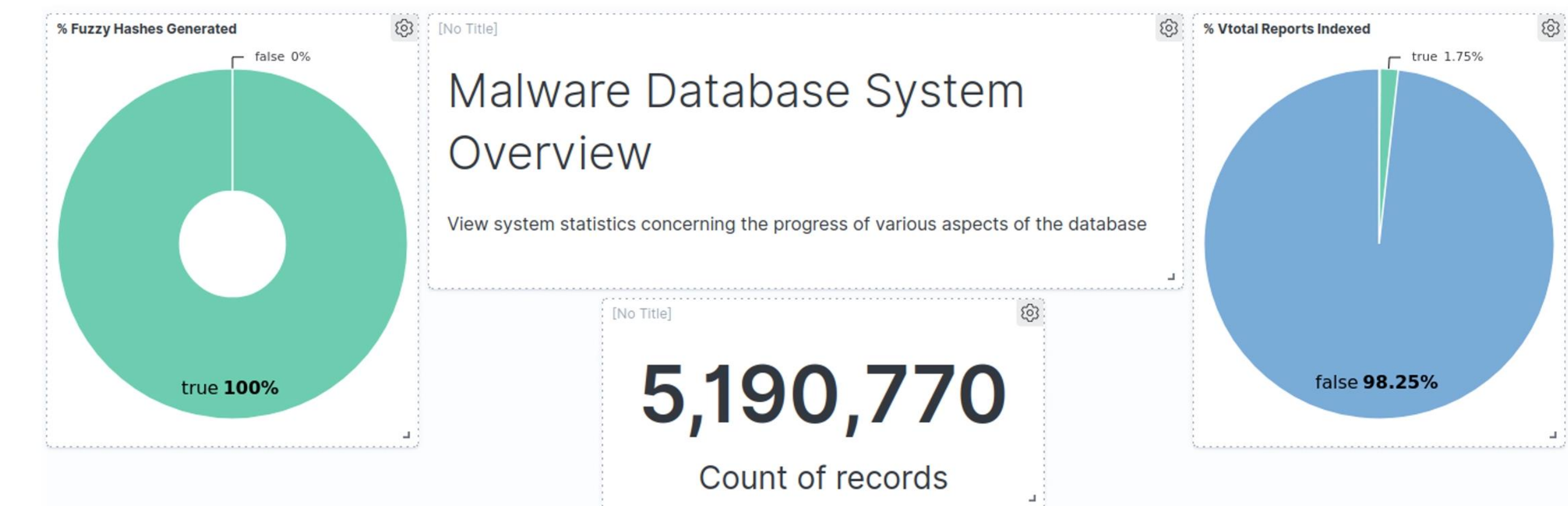
## Methods

- 27 TB of malware gathered from public sources
- VirusTotal reports gathered for each corresponding sample and uploaded into ElasticSearch
- SSDEEP and TLSH hashes and comparisons generated from sample set and stored in Neo4J
- Python based CLI program developed to interact with every feature



Darkside, ssdeep, pcode

## Results

- 5,190,770 Individual samples of Malware
- LSH used to detect similar malware with Neo4j
- Efficient and effective interaction with large scale data with ElasticSearch
- Automatic data collection and logging from open-source repositories
- Campaign visualization created to manage individual projects



## Challenges & Future Work

- Develop Neo4j to develop visualizations of custom indicator comparisons (custom LSH)
- Mathematically verify the integrity of LSH's
- Automatic VirusTotal report transforms to other threat indicator types (i.e. STIX)
- Develop search by library, function, etc capability
- Create easily deployable program (i.e. Dockerize)

UNIVERSITY of WYOMING

CEDAR