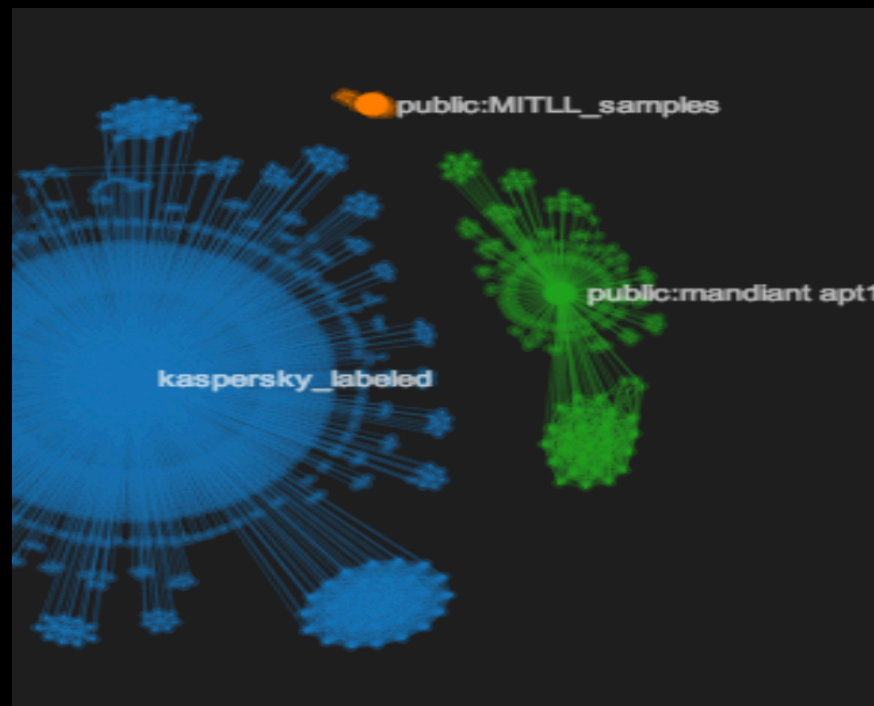


A Scalable, Ensemble Approach for Building and Visualizing Code-Sharing Networks over Millions of Malicious Binaries

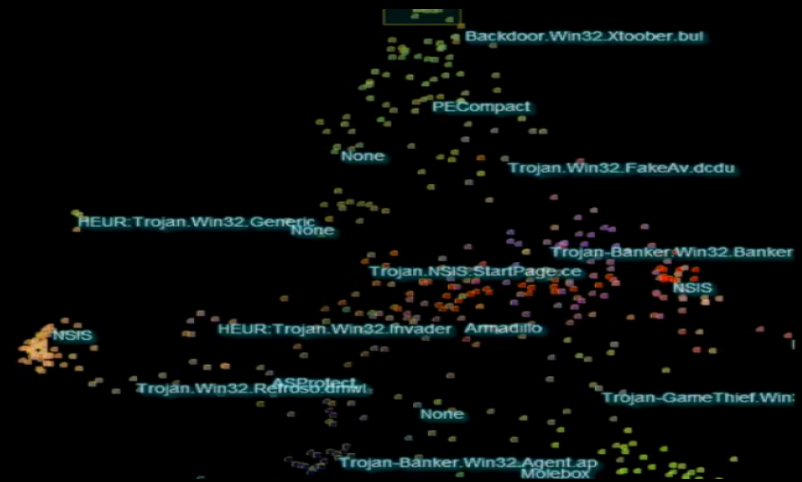
Joshua Saxe, Invincea Labs

Thanks to: Giacomo Bergamo, Robert Gove, Alex Long, Sigfried Gold



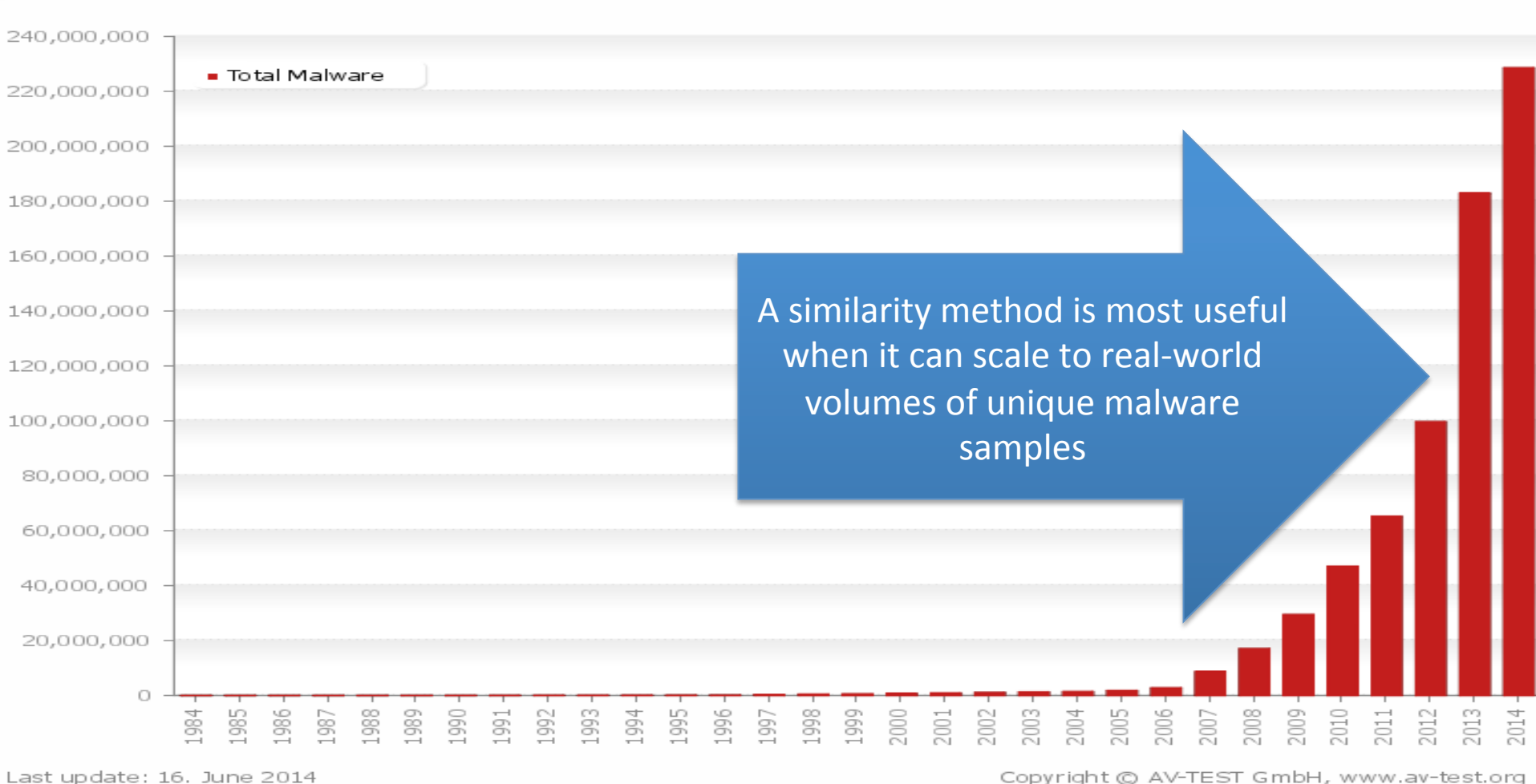
Identifying code-sharing networks: a key to advancing our ability to understand and detect malware

- Malware authors reuse work: graphical images, code, infection and persistence strategies
- Identifying reused work allows for analysts to reuse *analysis*
- Shared code networks are important pieces in the attribution puzzle
- Accurate shared code detection over huge volumes of malware is one piece of the next generation of malware *detection* approaches (we believe)



BUT: hard technical problems have stood in the way of accurately recovering malware code-sharing networks

Problem 1, scale: huge volumes of malware to evaluate for similarity



Problem 2: Accuracy, defeating multiple malware obfuscation methods

Blackhat malware author obfuscation strategy	Whitehat analyst feature recovery response	Limitations of whitehat approach
Packing	Dynamic analysis + memory dump	Unclear when we should dump memory
Code randomization	Abstract from code to control flow / data flow graph structures	No abstraction will defeat all code randomization techniques
Anti-debugging / anti-virtual machine	Static analysis	What if the sample is also very well packed?
Thin virtual machines and randomized bytecodes	Dynamic analysis + triggering mechanism	Triggering / code coverage is non-trivial and manual

Malware obfuscation makes it hard to always be right when we say whether or not two malware samples belong in the same lineage. Each individual similarity analysis technique can be defeated by a smart adversary.

More detailed summary of state of the practice: surface feature similarity

- Ssdeep : file format independent byte-level homology detection
- Pehash, portable executable similarity with a Boolean range
- Imphash, portable executable similarity with a Boolean range
- These are the *best* approaches because they scale, but they are brittle and low-cost for adversaries to defeat

peHash: A Novel Approach to Fast Malware Clustering

Georg Wicherski
RWTH Aachen University
gw@mwcollect.org

Abstract

Data collection is not a big issue anymore with available honeypot software and setups. However malware collections gathered from these honeypot systems suffer from multiple problems, such as:

ware analysis [2]. However, current honeypot setups suffer from gathering multiple binaries with distinct message digest sums that belong to the exact same specimen and therefore pollute malware databases as well as automated analysing

SimHash: Hash-based Similarity Detection

Caitlin Sadowski
University of California, Santa Cruz
supertri@cs.ucsc.edu

Greg Levin
University of California, Santa Cruz
glevin@cs.ucsc.edu

December 13, 2007

1 Abstract

Most hash functions are used to separate and obscure data, so that similar data hashes to very different keys. We propose to use hash functions for the op-

an aid to search. A standard technique in similarity detection is to map features of a file into some high-dimensional space, and then use distance within this space as a measure of similarity. Unfortunately, this typically involves computing the distance between all



Tracking Malware with Import Hashing

By [Mandiant](#) on January 23, 2014



Invincea Labs, 2014

ssdeep - Latest version 2.10

Quick Links

- [Download ssdeep](#)
- [The ssdeep man page](#)
- [Changelog](#)
- [Quickstart Guide](#)
- [API documentation](#)
- [Sourceforge project page](#) - Home to ssdeep development

Current research: usually deeper but often non-scalable feature similarity

- Dynamic execution approaches have the “coverage problem”, and take a long time to execute each sample
- Deep static approaches are time consuming and can also have a coverage problem
- The call graph is not directly observable in the program text and can only be fully recovered using often intractable symbolic execution techniques
- Ultimately, deep methods are necessary but are very hard to scale and have inherent limitations

Behavioral Clustering of HTTP-Based Malware and Signature Generation Using Malicious Network Traces

Roberto Perdisci^{a,b}, Wenke Lee^a, and Nick Feamster^a
^aCollege of Computing, Georgia Institute of Technology, Atlanta, GA 30332, USA
^bDamballa, Inc. Atlanta, GA 30308, USA
perdisci@gtisc.gatech.edu, {wenke, feamster}@cc.gatech.edu

Scalable, Behavior-Based Malware Clustering

Ulrich Bayer*, Paolo Milani Comparetti*, Clemens Hlauschek*, Christopher Kruegel[§], and Engin Kirda[¶]

* Secure Systems Lab, Technical University Vienna
{pmilani, ulli, haku}@seclab.tuwien.ac.at

[§] University of California, Santa Barbara
chris@cs.ucsb.edu

[¶] Institute Eurecom, Sophia Antipolis
kirda@eurecom.fr

J Comput Virol (2011) 7:233–245
DOI 10.1007/s11416-011-0151-y

ORIGINAL PAPER

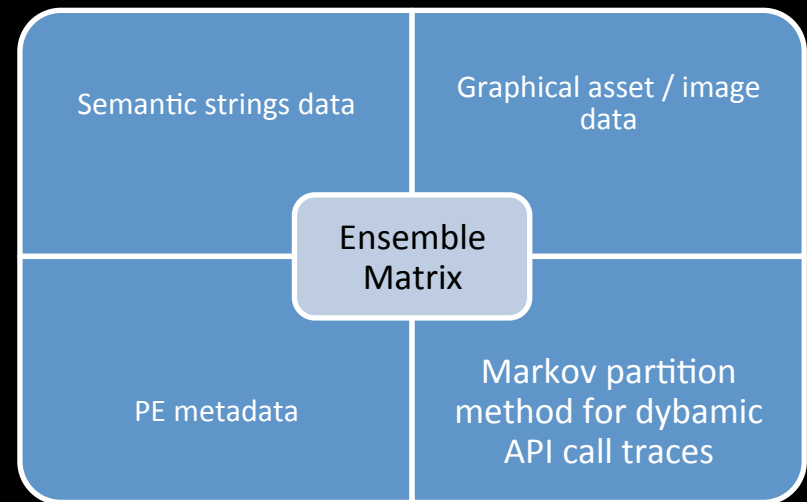
Malware classification based on call graph clustering

Joris Kinable · Orestis Kostakis

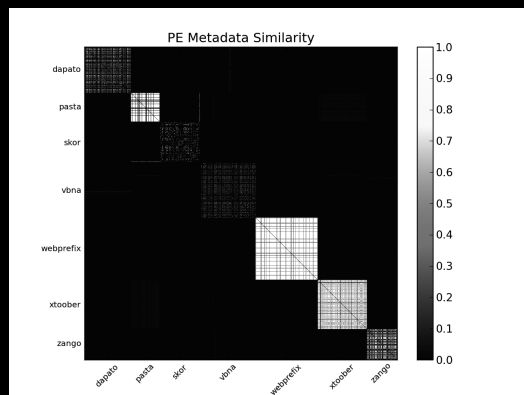
Our idea:

Scalable *and* deep shared-code network identification using ensemble code-sharing, locality sensitive hashing, and approximate nearest neighbor search

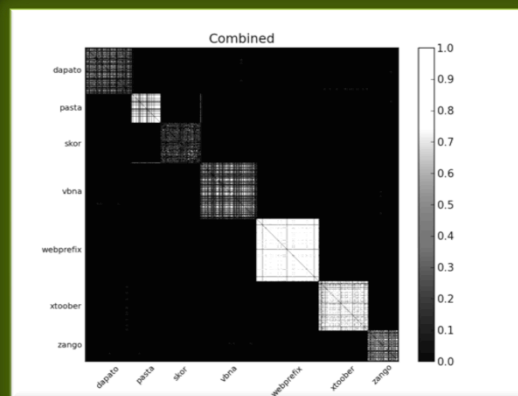
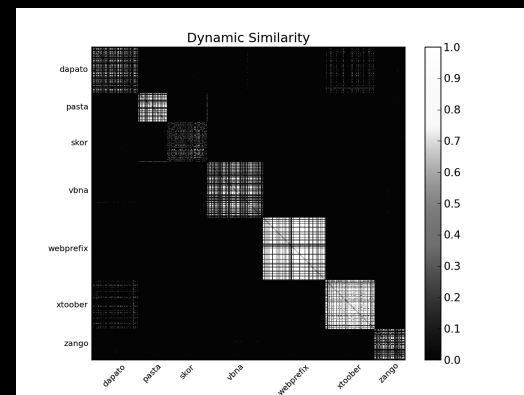
- While a malware author may succeed at obfuscating some subset of the malware feature domains, it's difficult to succeed in obfuscating across all software feature domains
- By integrating similarity sensors over multiple domains and combining sensor readings on a *pairwise* basis we can more accurately detect sample similarity
- ***Probabilistic data-structures (feature hash indexing and stochastic feature counting) allow us to scale to millions of samples***



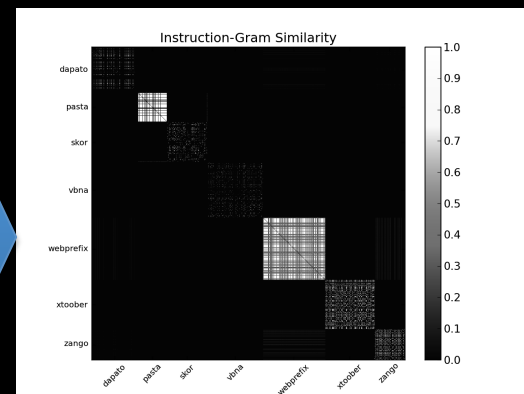
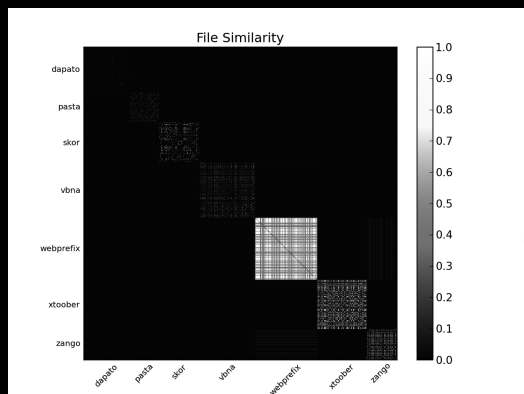
Evidence for the merits of the ensemble approach: four similarity matrices on the same ~2K samples (white boxes show families we have detected)



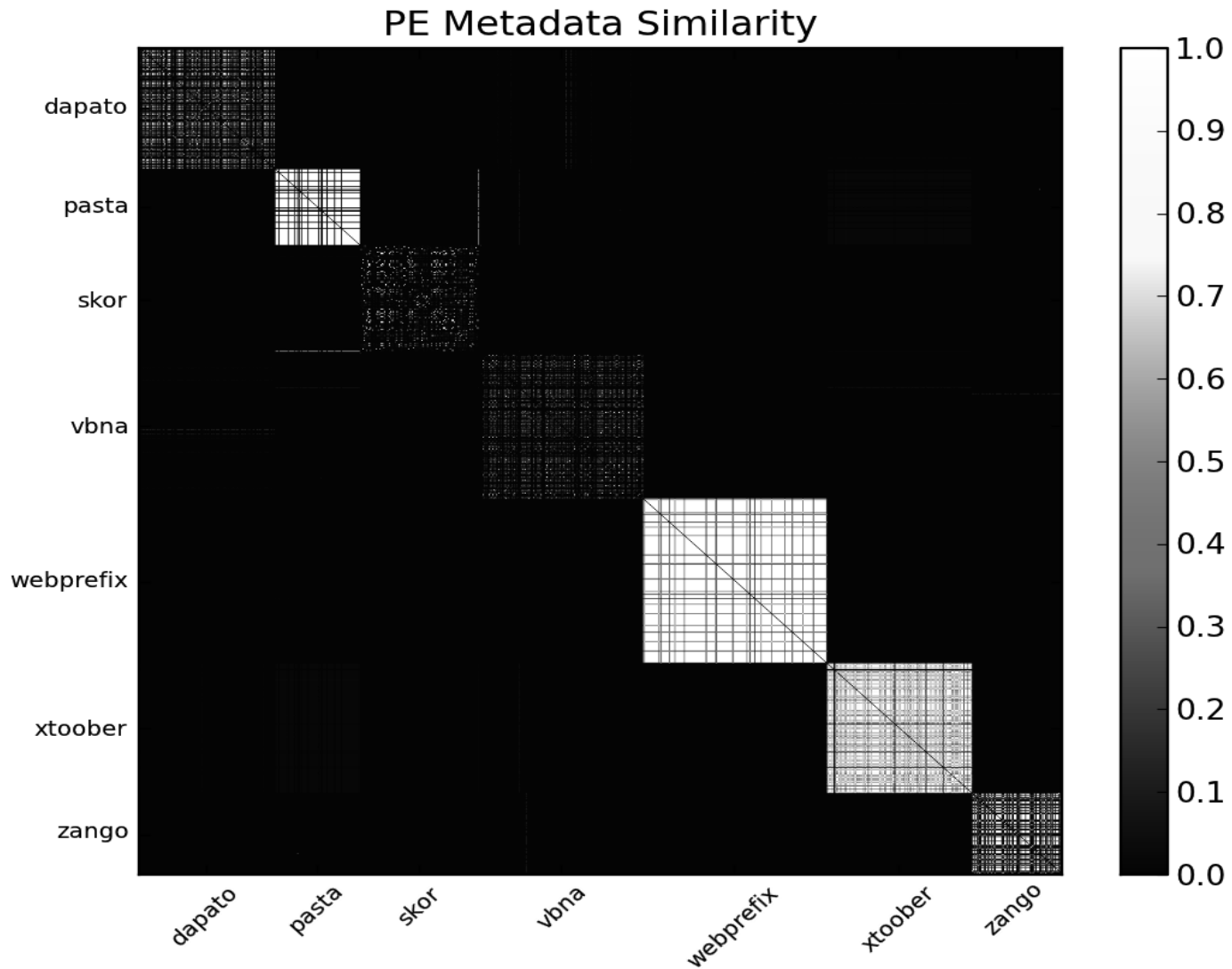
Combining matrices can produce better results



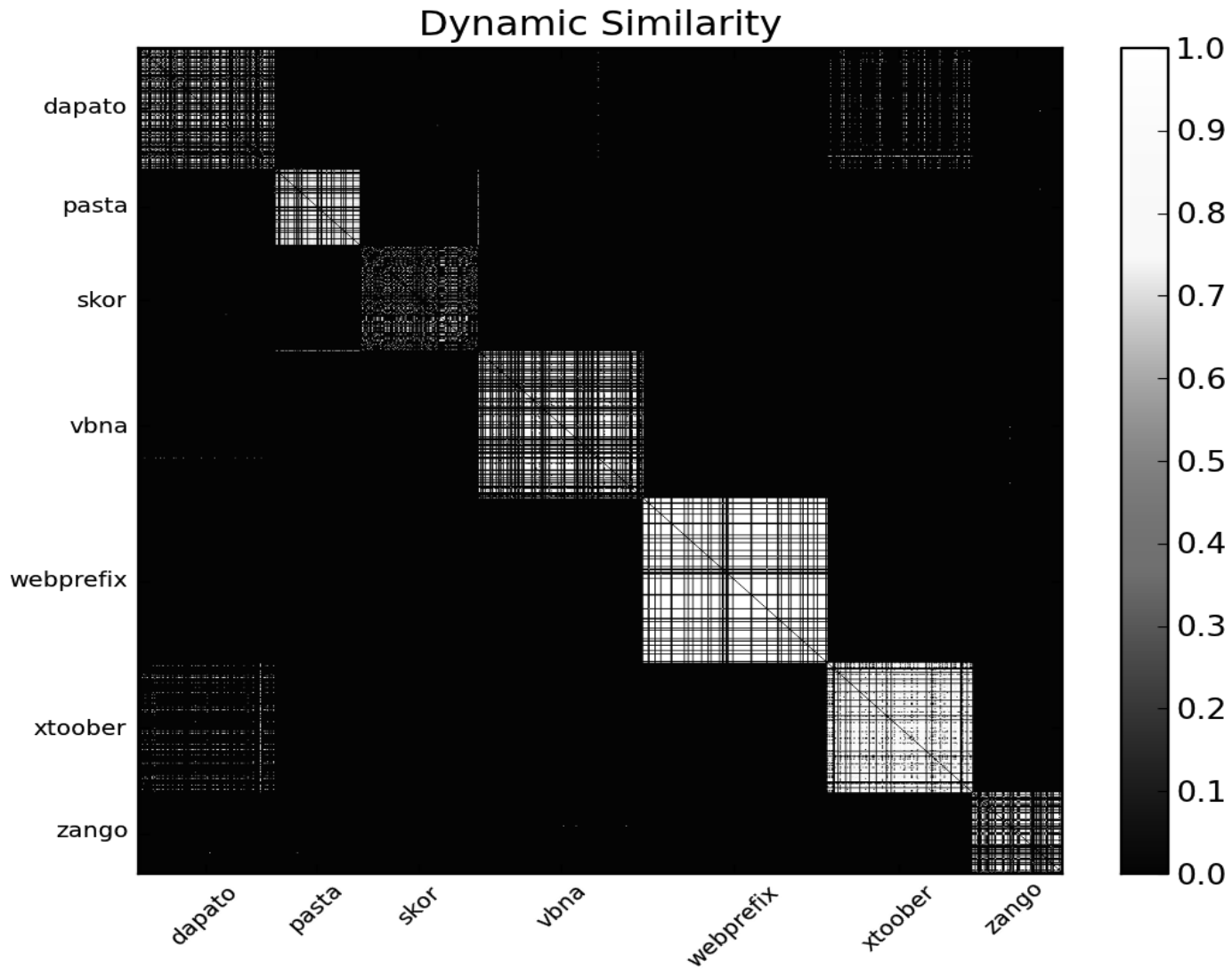
Looking at individual matrices can reveal what obfuscation techniques individual families use



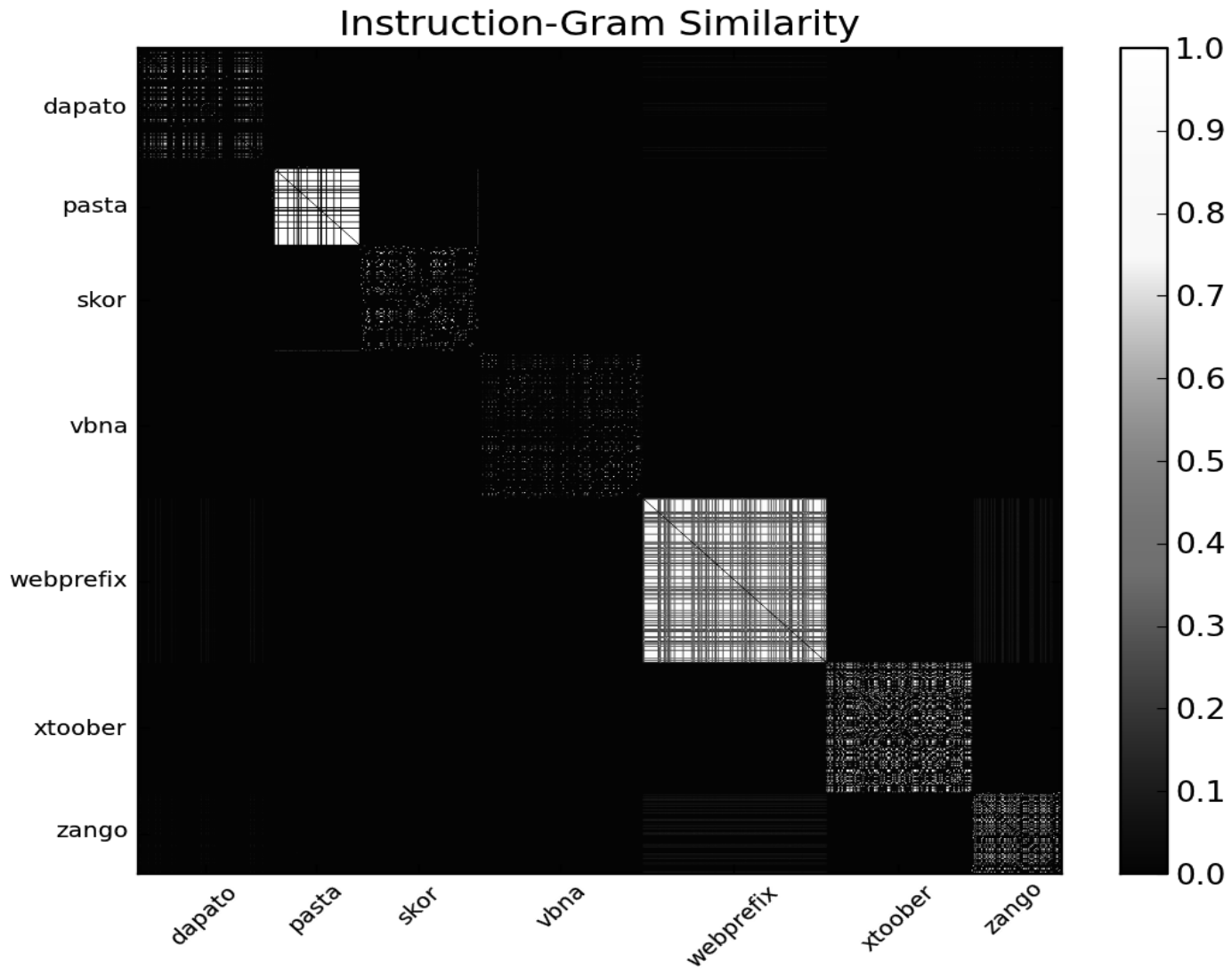
Evidence for the merits of the ensemble approach: four similarity matrices on the same ~2K samples
(white boxes show families we have detected)



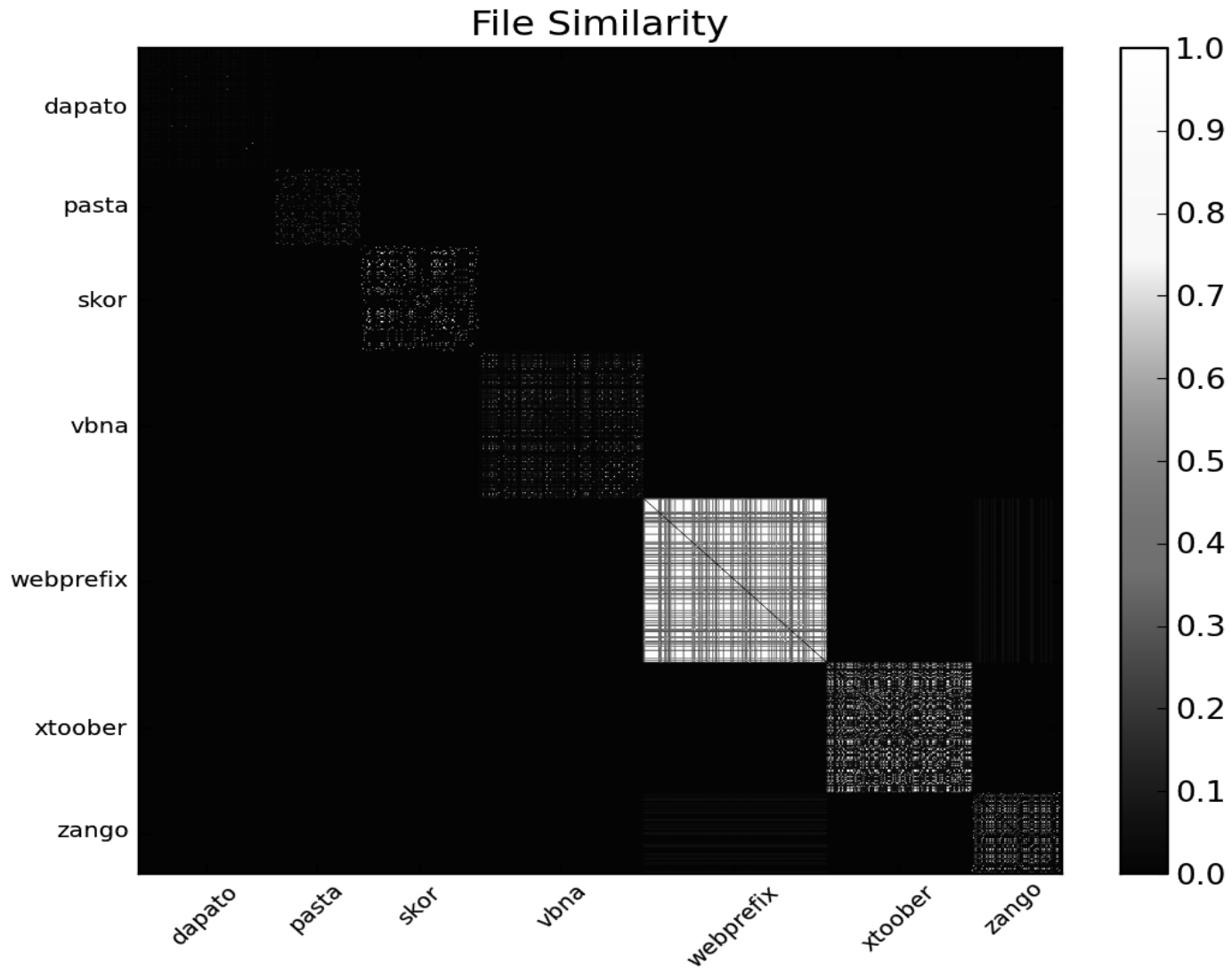
Evidence for the merits of the ensemble approach: four similarity matrices on the same ~2K samples
(white boxes show families we have detected)



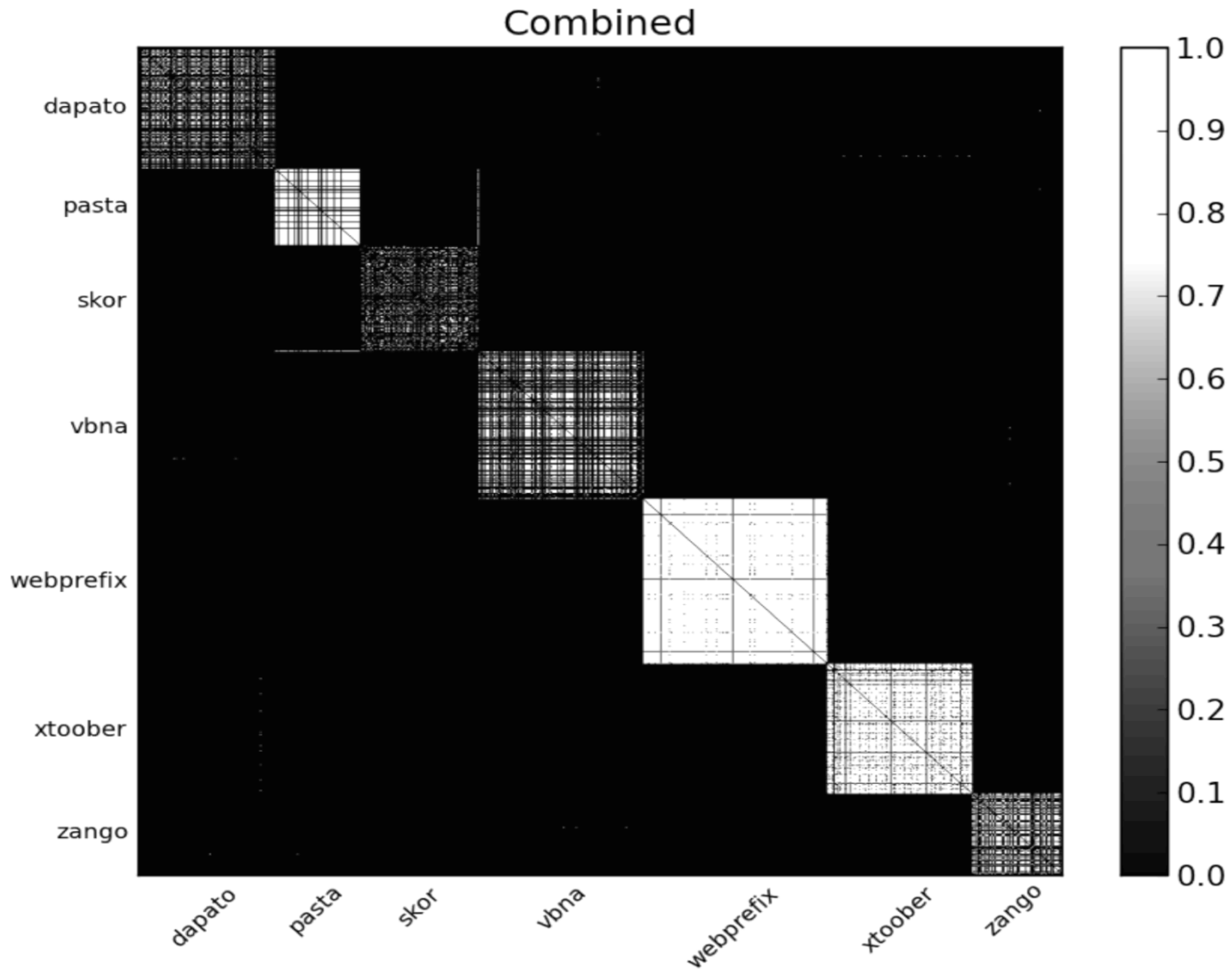
Evidence for the merits of the ensemble approach: four similarity matrices on the same ~2K samples
(white boxes show families we have detected)



Evidence for the merits of the ensemble approach: four similarity matrices on the same ~2K samples
(white boxes show families we have detected)

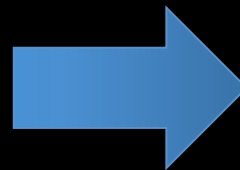
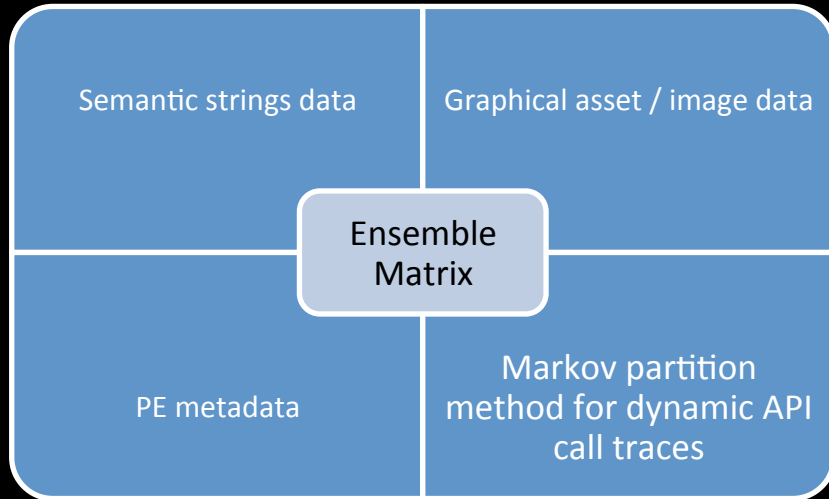


Evidence for the merits of the ensemble approach: four similarity matrices on the same ~2K samples
(white boxes show families we have detected)

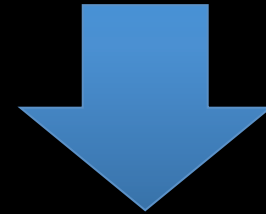


Scalable similarity network detection prototype design

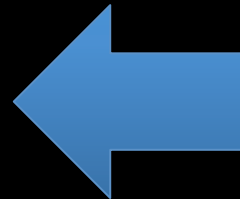
Static and dynamic feature extraction



Feature “sampling” and indexing approach, “OR” technique over matrices to get final similarity matrix



Scalable locality sensitive hashing based similarity index



Visualization



Deep dive into how we calculate sample similarity: Step 1, assigning malware features appropriate weights

Goal: Not all features extracted from malware should be considered equally when looking for shared code / shared provenance relationships between samples, the features that are more *rare* should be considered to be more important when they match

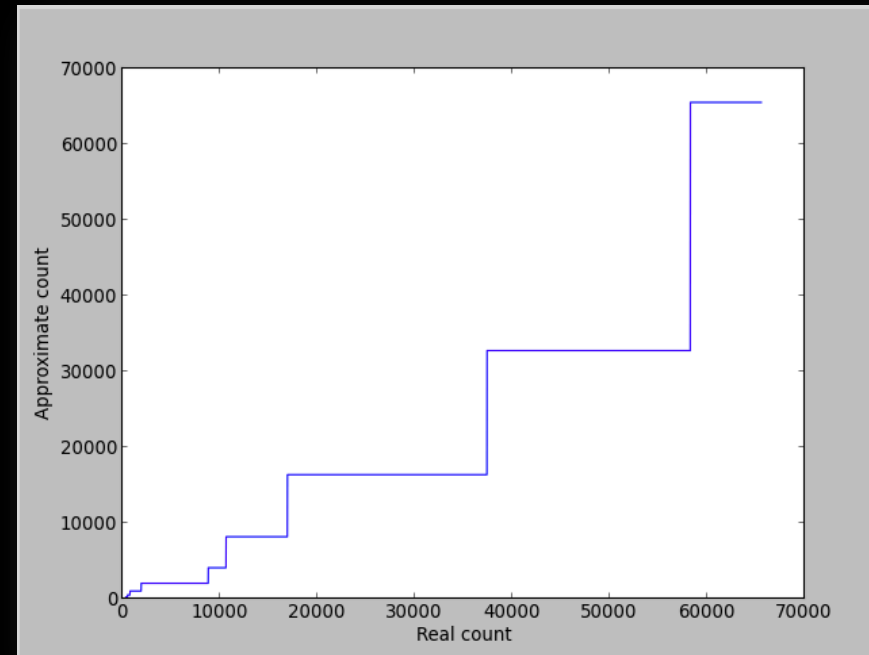
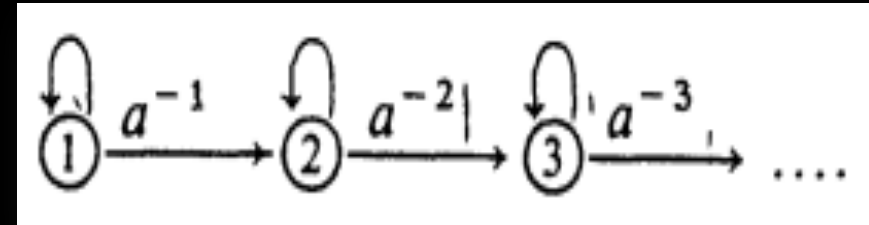
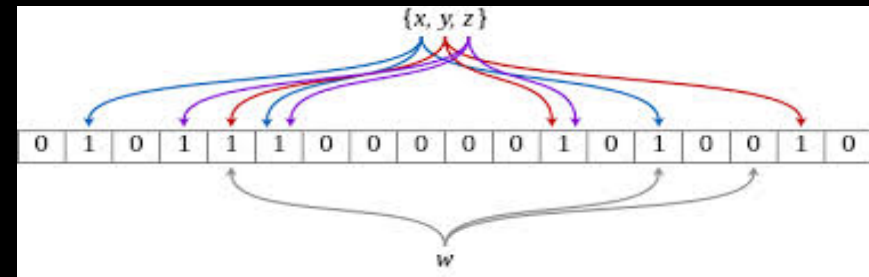
Problem: Counting occurrences of features at scale is hard; unbounded number of features, unbounded number of observations, bounded memory and CPU resources (*example: we estimate there would be 5.1 billion unique printable string features for a 30 million sample dataset*)

Solution: Morris counting for compact approximation of item counts

- Store only exponent values in counter
- Randomized counter increment value approximates true magnitudes

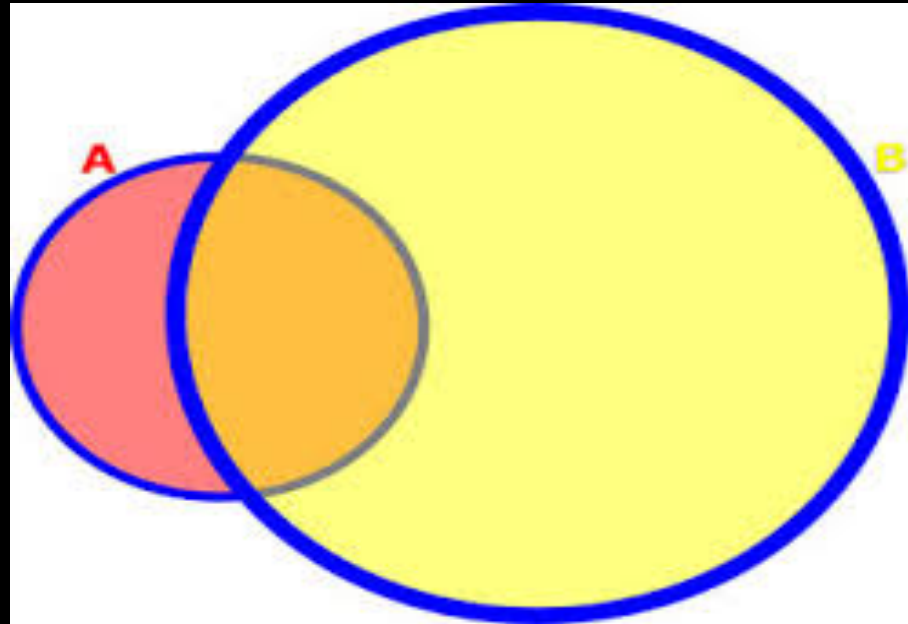
and ... adapted bloom filter for indexing of Morris counters

- Shared Bloom filter's properties; no false negatives, prone to false positives
- 8-bits per register, one hash function, 32MB data structure



Step 2, comparing malware samples for similarity: weighted Jaccard calculation

- Weighted Jaccard measures the size of the intersection of A and B, weighted by the importance of the features in the intersection (we use rarity, with some caveats, as a proxy for importance)
- Good, intuitive measure of similarity between two sets of features extracted from malware samples
- Thus would be appropriate to compute 'weighted Jaccard' similarity matrix for each of our ensemble components
- ... but it's slow to compute for every pair of malware samples



```

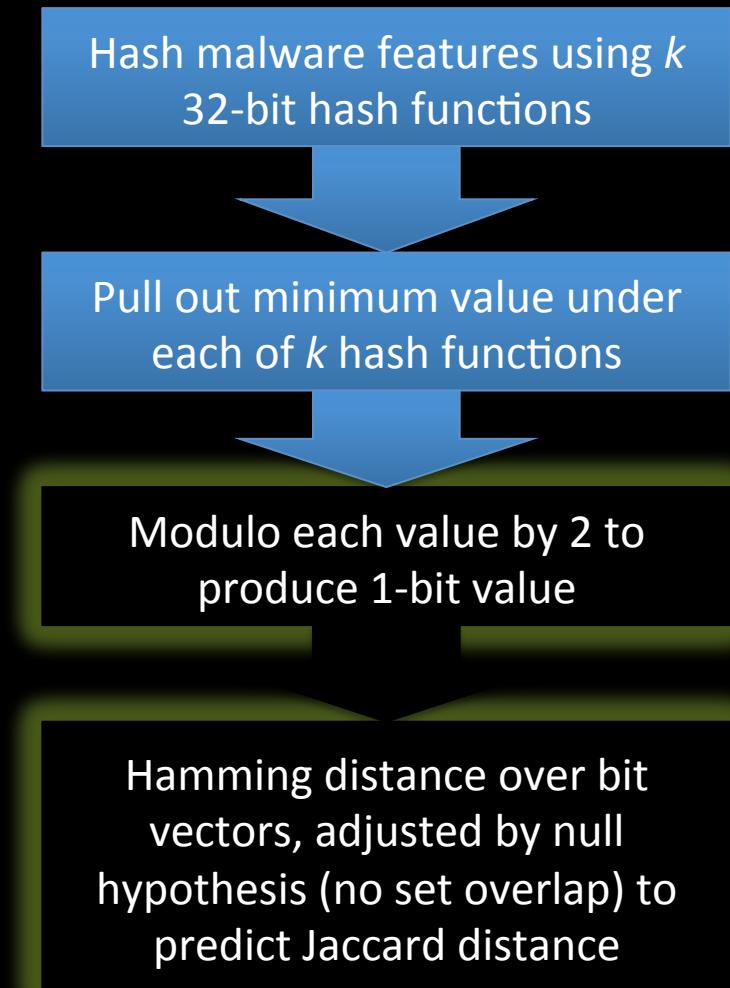
1 1.000
2 0.573 1.000
3 0.444 0.600 1.000
4 0.491 0.558 0.604 1.000
5 0.499 0.620 0.527 0.574 1.000
6 0.534 0.586 0.571 0.534 0.655 1.000
7 0.471 0.444 0.434 0.434 0.608 0.600 1.000
8 0.448 0.458 0.519 0.456 0.525 0.621 0.580 1.000
9 0.643 0.576 0.534 0.552 0.477 0.563 0.574 0.541 1.000
10 0.576 0.508 0.492 0.559 0.508 0.569 0.571 0.600 0.623 1.000
11 0.552 0.559 0.478 0.535 0.571 0.562 0.534 0.521 0.727 0.621 1.000
12 0.484 0.541 0.500 0.484 0.536 0.539 0.571 0.684 0.539 0.556 0.611 1.000
13 0.322 0.421 0.453 0.418 0.491 0.561 0.481 0.460 0.452 0.517 0.471 0.569 1.000
14 0.475 0.508 0.545 0.509 0.625 0.672 0.585 0.717 0.596 0.571 0.542 0.737 0.623 1.000
15 0.500 0.478 0.527 0.524 0.609 0.652 0.541 0.578 0.576 0.627 0.613 0.627 0.576 0.627 1.000
16 0.455 0.519 0.453 0.444 0.623 0.630 0.531 0.509 0.533 0.596 0.613 0.618 0.608 0.623 0.774 1.000
17 0.410 0.517 0.447 0.518 0.484 0.548 0.441 0.579 0.532 0.556 0.571 0.661 0.574 0.632 0.585 0.566 1.000
18 0.466 0.517 0.474 0.518 0.596 0.600 0.547 0.607 0.548 0.556 0.528 0.590 0.635 0.704 0.594 0.604 0.614 1.000
19 0.443 0.452 0.509 0.526 0.649 0.581 0.574 0.641 0.508 0.534 0.587 0.611 0.696 0.641 0.630 0.649 0.755 1.000
20 0.492 0.525 0.558 0.607 0.590 0.578 0.544 0.532 0.613 0.611 0.630 0.609 0.452 0.607 0.565 0.500 0.565 0.641 0.597 1.000
21 0.524 0.523 0.508 0.574 0.581 0.574 0.581 0.574 0.603 0.581 0.603 0.581 0.603 0.581 0.581 0.581 0.581 0.581 0.581 0.581 1.000
22 0.482 0.536 0.473 0.566 0.596 0.593 0.588 0.579 0.548 0.610 0.641 0.590 0.619 0.615 0.534 0.552 0.596 0.649 0.618 1.000
23 0.483 0.508 0.517 0.508 0.574 0.574 0.492 0.567 0.554 0.609 0.632 0.619 0.508 0.617 0.597 0.600 0.606 0.574 0.581 0.578 0.551 0.586 1.000
24 0.541 0.500 0.568 0.559 0.653 0.716 0.589 0.556 0.705 0.620 0.639 0.606 0.500 0.620 0.637 0.596 0.615 0.571 0.541 0.604 0.708 0.623 0.552 1.000
25 0.473 0.525 0.459 0.525 0.617 0.578 0.526 0.583 0.563 0.585 0.620 0.635 0.579 0.633 0.662 0.636 0.516 0.574 0.500 0.619 0.636 0.702 0.578 0.750 1.000
26 0.462 0.463 0.446 0.532 0.569 0.559 0.517 0.493 0.582 0.612 0.714 0.500 0.532 0.493 0.614 0.632 0.569 0.576 0.574 0.614 0.643 0.606 0.594 0.621 1.000
27 0.492 0.508 0.525 0.593 0.596 0.552 0.452 0.547 0.615 0.646 0.653 0.529 0.536 0.594 0.672 0.567 0.554 0.530 0.585 0.631 0.623 0.655 0.641 0.580 0.606 0.657 1.000
    
```

Figure 2. Genetic similarity coefficient matrix for 27 individuals of *S. pallens*, based on 79 PCR fragment patterns obtained from 10 RAPD primers. Values were calculated by Jaccard coefficient. Individuals number 1 to 11 are from Central Brazil (DF) and 12 to 27 are from the Northeast (RN).

Universal problem in similarity network detection: computing pairwise similarities is **quadratic**

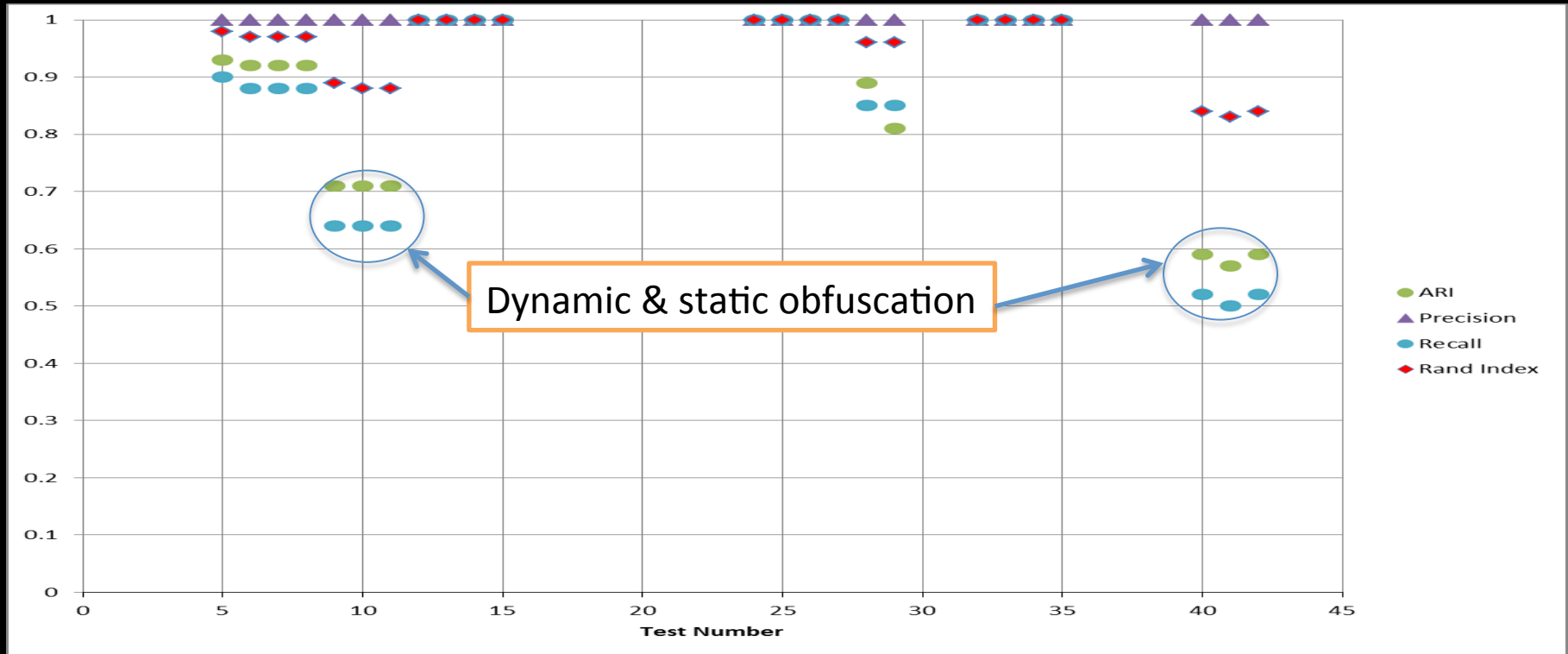
Step 3, scaling pairwise similarity calculations: Binary version of weighted 'Minhash' algorithm

- Minhash is a dimensionality reduction technique
- Allows for fast *approximate* computation of pairwise set similarities
- We've modified the minhash algorithm to further reduce feature vectors to binary vectors
- We can then use fast XOR based Hamming-distance calculation to approximate minhash, which in turn approximates weighted Jaccard

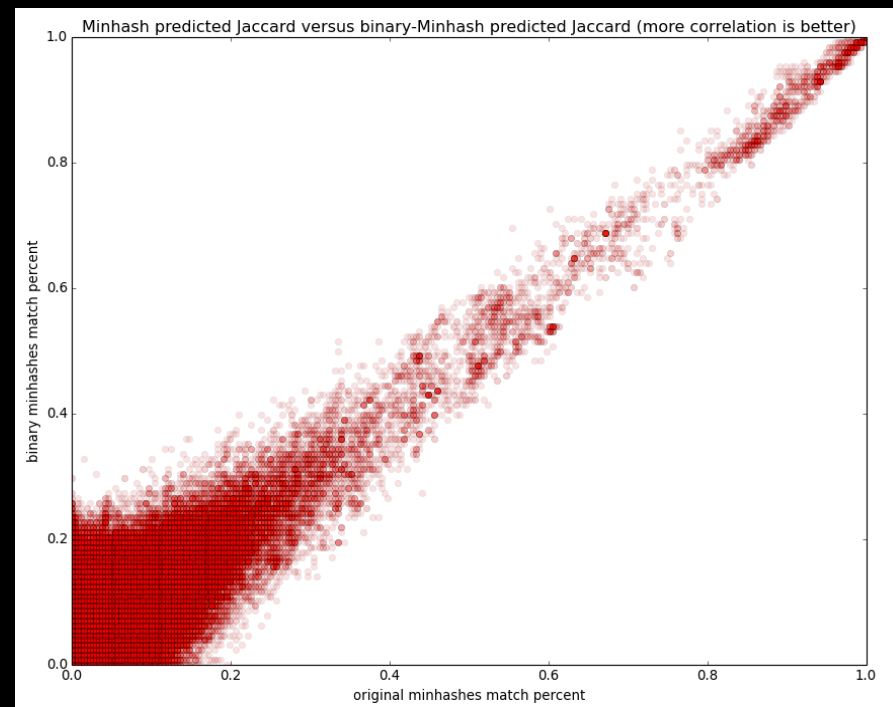
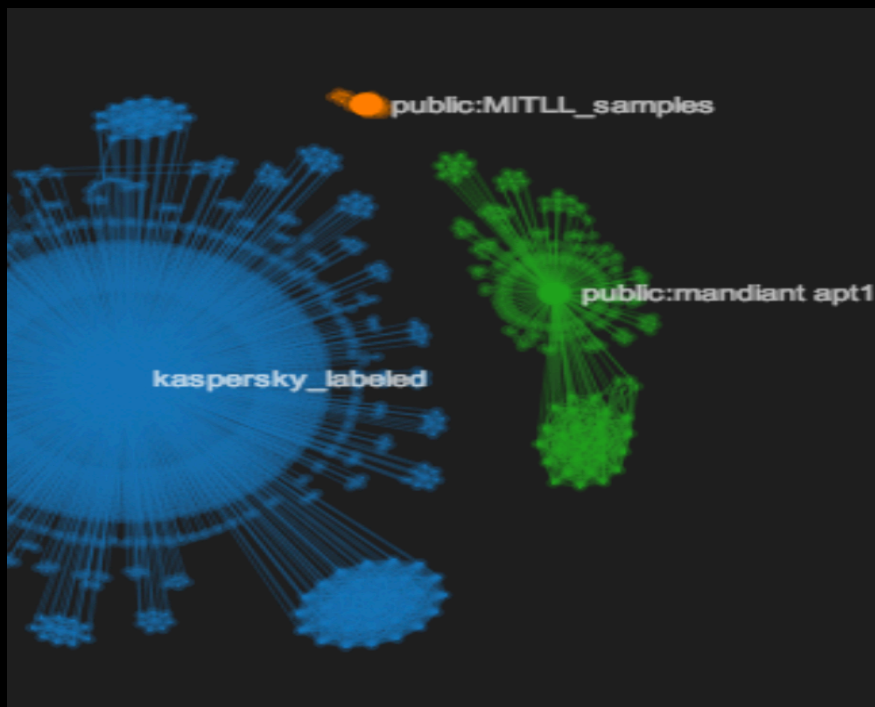


Overall Similarity-Network based Clustering Performance

- 10k extremely well labeled samples used in test
- Consistent results in tests where the evaluated samples the same, regardless of sample size
- Perfect performance on 10 out of 21 clustering tests
- Precision of 1.0 on all tests
- Dramatically outperformed 'ssdeep' baseline (about double average adjusted rand index)
- Highly scalable; we've clustered 2M samples and could cluster more with appropriate hardware



Demonstration of malware similarity network / clustering tool



Questions and Discussion

