

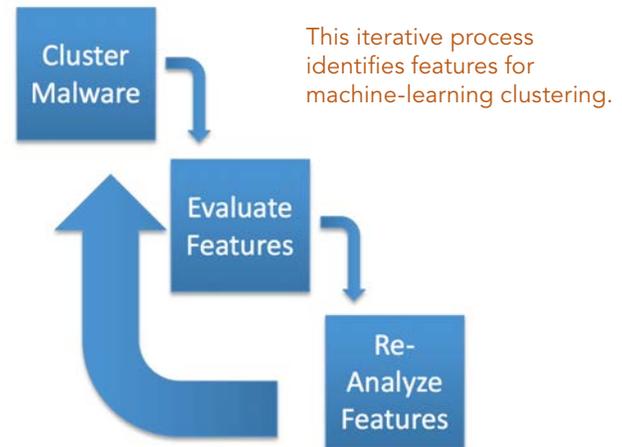


attribution, depending on the quality of the signatures. We have developed two approaches to solve these problems that take advantage of high-performance computing, machine learning, and bio-inspired techniques. These two approaches can be combined to generate better classifications and identifications of malware. The approaches will support the concepts of generating family trees of malware, which can then inform attribution.

## TECHNICAL APPROACH

The team at Pacific Northwest National Laboratory (PNNL) developed Machine Learning String Tools for Operational and Network Security (MLSTONES), a computational high-performance computing capability for applying biosequence analysis to cybersecurity applications. MLSTONES can be used to discover similarity in software, particularly for malware. The MLSTONES process creates “cyber protein” representations of malware and then uses protein alignment techniques to quickly generate families of proteins. With this method, we can create a single representation of an entire family of entities, significantly reducing the amount of data to analyze. MLSTONES has been applied to the analysis of a malicious software corpus to produce “family trees” of malware artifacts and a signature library. These family trees can be shared, to deploy the analytical technique into operational environments. Our high-performance computing implementation of this biosequence-based capability is ideal for malicious code analysis because we can infer the function of a “cyber protein” by its relationship to other similar proteins, the same process used in biology.

The team at Mississippi State University (MSU) is identifying features of malicious software that can be used to help define capabilities and relationships of malware. MSU’s focus is on dynamic features, those generated at runtime. The dynamic features are generated while executing the malware in a sandbox environment. During execution, any interactions with the operating system



are recorded and stored as features. Once execution is complete, memory is captured and analyzed to extract additional features that the sandbox might have missed at runtime.

These features are stored in a database and can be used to identify the malware capabilities and relate groups of similar samples. High-performance computing resources are necessary to effectively process the volume of new malware being discovered on a daily basis.

The combination of PNNL’s MLSTONES with MSU’s dynamic features enables the rapid characterization, classification, and detection of malware in the wild.

## IMPACT

Work on this project increases the likelihood of finding the source of an attack or widespread malware infection, as well as identifying new zero-day types of attacks. With the ability to leverage high-performance computing resources, analysts will be able to take current sources of potential new malware samples and quickly classify them by potential authorship or common source. We will also be able to identify the type of those new malware samples—even if they are not exact matches to standard signatures. New malware samples may be grouped with existing ones that have a high degree of similarity, giving analysts more information on specific threat actors and new types of malware. This benefits those who wish to attribute cyber attacks by providing more data that can be examined for operational mistakes by the attackers.

## Contacts

### Elena Peterson

Principal Investigator  
(509) 372-4573  
Elena.Peterson@pnnl.gov

### Dae Glendowne

Principal Investigator  
(662)-325-5964  
jg905@msstate.edu

### David A. Dampier

Director, Distributed Analytics  
and Security Institute  
(662) 325-0779  
dampier@dasi.msstate.edu

### John R. Johnson

Program Director  
(509) 375-2651  
John.Johnson@pnnl.gov

