# Novel, Large Scale, High Speed Text and Data Analysis Research

Presented by

## Thomas E. Potok, Ph.D

**Group Leader**
**Applied Software Engineering Research**

# Who we are

**ASER** | Applied Software Engineering Research Group

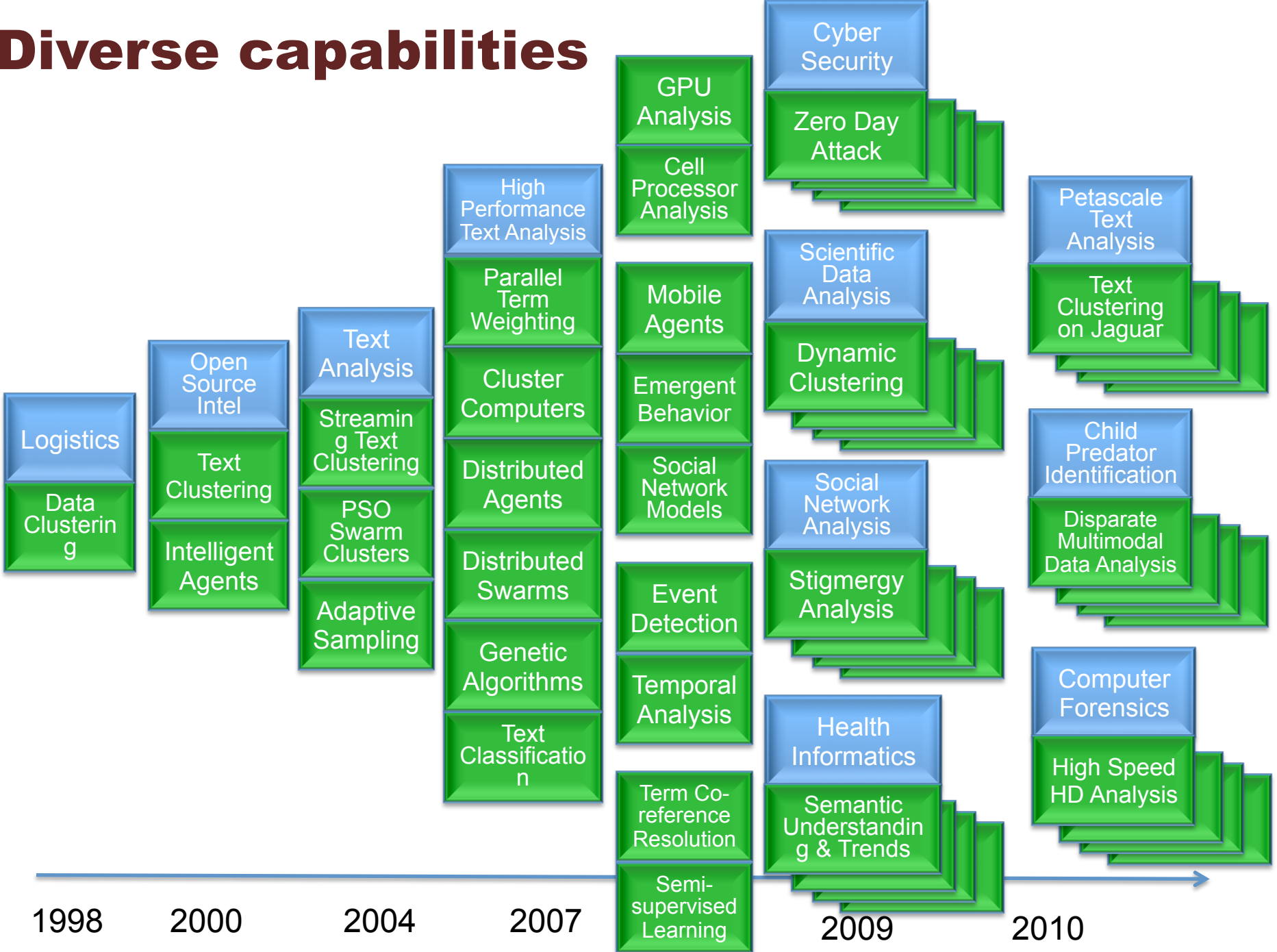- ## Group focus
  - 9 years of research in intelligent agents, emergent behavior, pervasive computing, machine learning, information retrieval, and knowledge discovery
  - In 2009 we organized 7 research workshops and published 26 papers and book chapters
  - Hands-on experience with DHS, Military, IC, and Industry

- ## Group success
  - $15M in research investment
  - 14 staff members, 9 PhDs in computer science
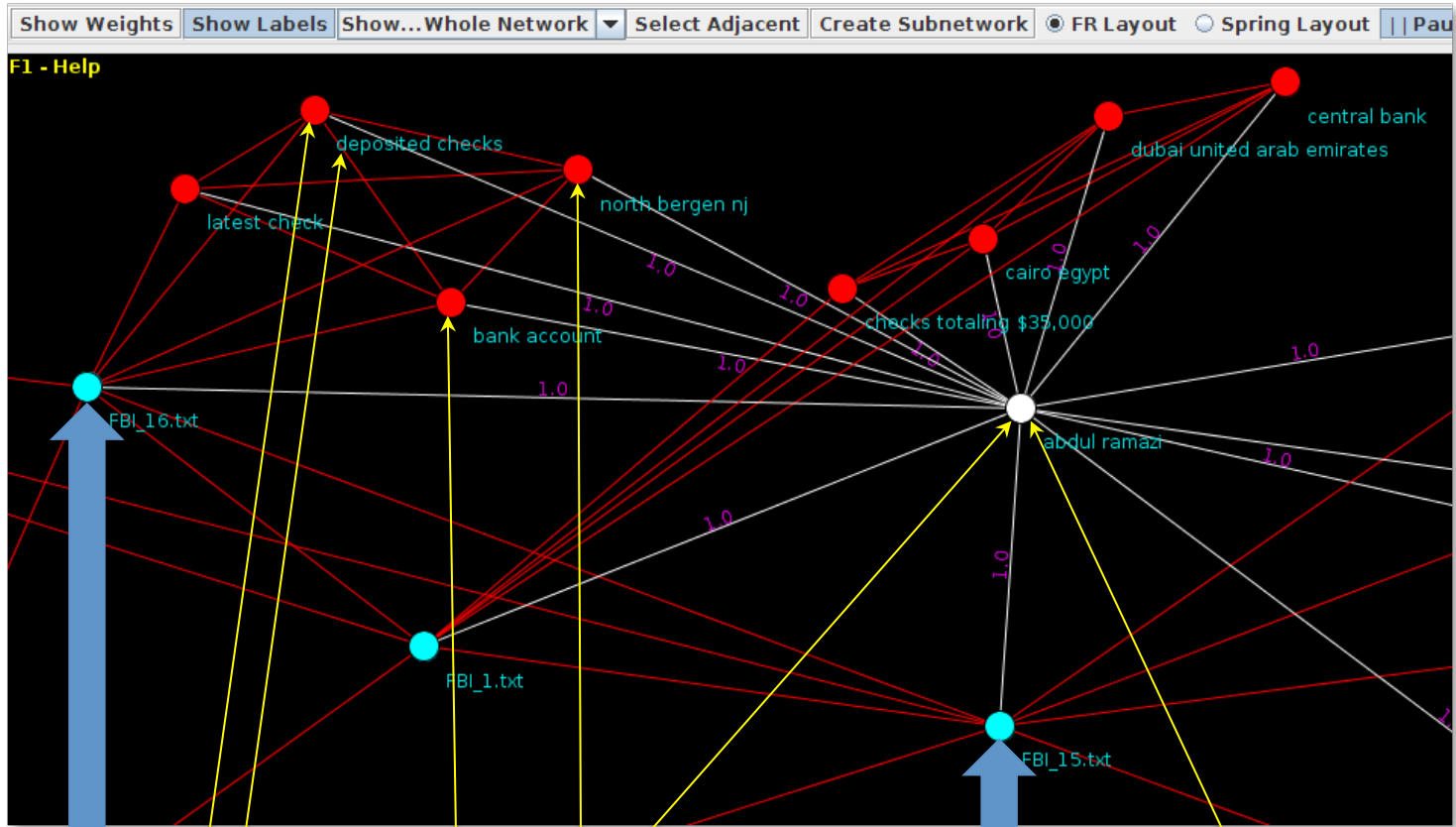  - R&D 100 Award (Oscars of invention) in 2007

Potok_ASER_SC11

OAK RIDGE National Laboratory

# Diverse capabilities

**Logistics**
Data Clustering

**Open Source Intel**
Text Clustering
Intelligent Agents

**Text Analysis**
Streaming Text Clustering
PSO Swarm Clusters
Adaptive Sampling

**High Performance Text Analysis**
Parallel Term Weighting
Cluster Computers
Distributed Agents
Distributed Swarms
Genetic Algorithms
Text Classification

**GPU Analysis**
Cell Processor Analysis

Mobile Agents
Emergent Behavior
Social Network Models

Event Detection
Temporal Analysis

Term Co-reference Resolution
Semi-supervised Learning

**Cyber Security**
Zero Day Attack

**Scientific Data Analysis**
Dynamic Clustering

**Social Network Analysis**
Stigmergy Analysis

**Health Informatics**
Semantic Understanding & Trends

**Petascale Text Analysis**
Text Clustering on Jaguar

**Child Predator Identification**
Disparate Multimodal Data Analysis

**Computer Forensics**
High Speed HD Analysis

1998    2000    2004    2007        2009    2010

# Text analysis capability overview

| Capability | Capacity in documents | Search Engines | Natural Language Tools | Piranha |
|---|---|---|---|---|
| Search/Query | 100M+ | Yes | No | Yes |
| Unsupervised classification | 1M | Some | No | Yes |
| Supervised classification | 1M | No | No | Yes |
| Clustering (Full document analysis) | 100K | No | No | Yes |
| Term Frequency Analysis (Significant words) | 100K | Yes, but not available to user | Yes | Yes |
| Semantic Extraction (People, places) | 1000 | No | Yes | Yes |

Potok_ASER_SC11

OAK RIDGE
National Laboratory

# Piranha cluster view



Report Date: 1 April, 2003. CIA [From MI5]: On 30 March, 2003 the British Special Branch arrested Omar Bakri Qatada at his home at #11 St. Mary's Terrace, Paddington, London. Found in Qatada's bedroom was a small carton holding 10 ounces of Pentaerythritol [PETN] and Triacetone Triperoxide [TATP]. This is the same explosive that Richard Reid attempted to use on American Airlines #63 from Paris to Miami on 22 December, 2001. The BSB were alerted to follow and detain Qatada on the basis of information obtained from a respected moderate Moslem cleric in London, whose name was not provided in this report from MI5.

Report Date: 1 April, 2003. FBI: Abdul Ramazi is the owner of the Select Gourmet Foods shop in Springfield Mall, Springfield, VA. [Phone number 703-659-2317]. First Union National Bank lists Select Gourmet Foods as holding account number 1070173749003. Six checks totaling $35,000 have been deposited in this account in the past four months and are recorded as having been drawn on accounts at the Pyramid Bank of Cairo, Egypt and the Central Bank of Dubai, United Arab Emirates. Both of these banks have just been listed as possible conduits in money laundering schemes

# Term network with document links



16) Report Date 22 April 2003. FBI: Hani al Hallak, of **North Bergen, NJ**, has **deposited checks** in his **bank account** that were drawn on First Union Bank account number 1070173749003 in Springfield, VA, in the name **Abdul Ramazi**. The latest check is dated 16 April 2003 and was in the amount of $8500.

15) Report Date 20 April 2003. FBI: Mukhtar Galab has an account at the Virginia National Bank in Charlottesville, VA. Bank records say he has deposited several checks in the last three months, totaling $13,000, drawn on account number 1070173749003 held by **Abdul Ramazi** at the First Union Bank in Springfield, VA

Potok_ASER_SC11

# Petascale text analysis

- **ORNL's Jaguar is the fastest computer in the world**
  - 255,000 cores  -10PB (13,400 1TB drives) of Storage  -362TB of memory

- **Google has indexed 1 Trillion unique URLs, but has not analyzed the content of the information**

- **We are currently developing petascale text analysis techniques to cluster (deep analysis) of 1 trillion documents using Jaguar**

Potok_ASER_SC10

# Exfiltrate the data from real targets

## Hacker's Machines

173.232.185.102
Poland

231.187.219.203
Brazil

168.217.122.006
Australia

240.034.149.210
Arizona

090.198.241.046
France

060.96.230.017
Oak Ridge National Laboratory

Seattle, WA
Starbucks

**NEW YORK STATE DRIVER LICENSE**

KB00000
U.S. Department of Energy
Oak Ridge Operations Office

Visitor

Q

C

TEST
MAINTENANCE

000001

Oak Ridge National Laboratory

## Target Network

**Jane Doe**
**Visitor Control**
Names, SSN, DOB, dates of visits, clearances, authorized buildings

**Bob Smith**
**Computer Scientist**
Vulnerabilities of SCADA systems

- **Download all files from Smith and Doe's computers**
- **Sell Smith's data to highest bidder – Iran or North Korea**
- **Find open visitor requests for a visitor with a clearance**
- **Forge a driver license**
- **Pick up the badge at Visitor Control**

Potok_ASER_SC11

OAK RIDGE
National Laboratory

# New approach

| Semi-Supervised Learning | Cluster Analysis | Particle Swarm Optimization | Temporal Analysis | Categorization |
|---|---|---|---|---|

**Analysis Interface**

**Cyber Data Fusion and Interface**

| Feature Extractor | Splunk |
|---|---|

Suspect Data

Network Metrics

Network Alerts

Computer Logs

Network Packets

Snort Intrusion Detection System Alerts

ORNL Developed

Third-Party Software

# Find attacks on vulnerable computers

# Find exploits



Hacker's Machines

- 173.232.185.102 Poland
- 231.187.219.203 Brazil
- 168.217.122.006 Australia
- 240.034.149.210 Arizona
- 090.198.241.046 France

Seattle, WA Starbucks

Target Network

**Bob Smith**
**Computer Scientist**
Vulnerabilities of SCADA systems

**Jane Doe**
**Visitor Control**
Names, SSN, DOB, Dates of visits, clearances, authorized buildings

## OAK RIDGE CYBER ANALYTICS

Correlation Engine Control
Cyber Analysis Configuration
Cyber Event Analysis

20 mins

Alerts

Time Range

From: Hour: 0  Min: 0
To: Hour: 0  Min: 0

Key Words

IP Address:
Keyword:

Search
<<  >>

Event Listing | Cluster Analysis | Swarm Analysis | Categorization

Model: IP Address Historical
Iterations
Similarity: 35
Iteration: 130

Run    Continue

Five external machines behaving suspiciously

Current node: 152
x: 134
y: 416
Src IP: 86.164.180.1
Src Port: Source_port N/A
Dest IP: Dest_ip N/A
Dest Port: Dest_port N/A
Time: time N/A
Behavior Code: 1
Behavior: 86.164.180.1  [Behavior unknown and alarm action is: (spp_frag3) Fragmentation overlap

Third Party Proprietary

# Current capabilities

- **TRL 7+ in 18 month**
  - Live Sprint and AT&T feeds
  - 3+ terabytes of network traffic daily

- **Extraction of network features – Alerts plus network traffic metrics**

- **Machine learning – Incorporate new insights into analysis**

- **Cluster analysis – Group similar alerts**

- **Particle swarm optimization – Group similar computer behavior**

- **Temporal analysis – View time vs. IP Addresses vs. alert type**

- **Categorization – Group alerts by defined category**

OAK RIDGE National Laboratory

# Child predators



Somer
Thompson
2002 - 2009



- This is Somer. She was born April 5th 2002 in Columbia, South Carolina.

- She enjoyed swimming and dancing. She often thought of herself as becoming a ballerina.

- On October 19th 2009 at 2:45pm Somer left Grove Park Elementary school to meet up with her brother and sister for their daily mile walk home.

- She ran ahead of her siblings and got separated from the group

- This was the very last time Somer Thompson was seen alive.

- Her body was found in a Georgia landfill 2 days later. At that time, there were no suspects.

OAK RIDGE
National Laboratory

# Child predators (continued)



## Jarred Harrell

**He abducted Somer that day.
He was charged with raping
and murdering her.**

**"…Police had possession of Harrell's computer
and camera two months before Thompson
disappeared, …'they found Harrell's camera
containing a video with horrifying images of
Harrell molesting a three year-old child',"
reports ABC. -** Saturday, 17 April 2010 10:54 ABC News' 20/20

OAK RIDGE
National Laboratory

# Fact



**If Jarred Harrell's computer
and camera had been
forensically analyzed
in a timely manner,
Somer would still be alive.**

Potok_ASER_SC11

OAK
RIDGE
National Laboratory

# Sampling



Ignores system files

Immediately flag known illegal files

Catalogs files of interest

- **Sample suspect's disk to find a predator immediately**

- **Assuming illegal files are grouped together becomes a search problem**

- **Full analysis can be completed later**

Potok_ASER_SC11

# Artemis – 30 child pornography arrests in Knoxville, TN

- **Working with the *Protect* and *Internet Crimes Against Children* (ICAC) Task Force in Knoxville, TN**

- **A prototype thumb drive system to rapidly analyze information on a hard drive and identify illegal material**

- **When booted loads Knoppix and Artemis to allows the officer to scan specific directories or the entire disk for images and documents.**

- **Image analysis looks for high percentage of skin tones, faces, or an illegal hash values**

- **Text analysis looks for documents that are similar to a provided set.**
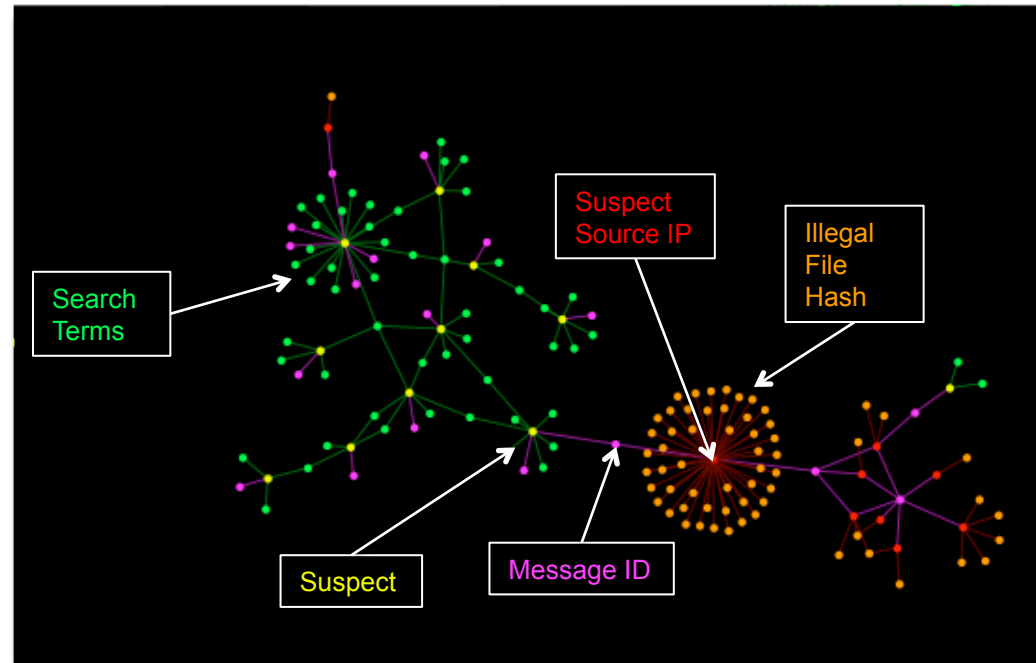
Potok_ASER_SC11

OAK RIDGE National Laboratory

# ASLAN – Graph Analysis of Large file sharing network

- You can make part of your computer public, so that others can download your files

- The software keeps track of the IP addresses of the computers that download the files

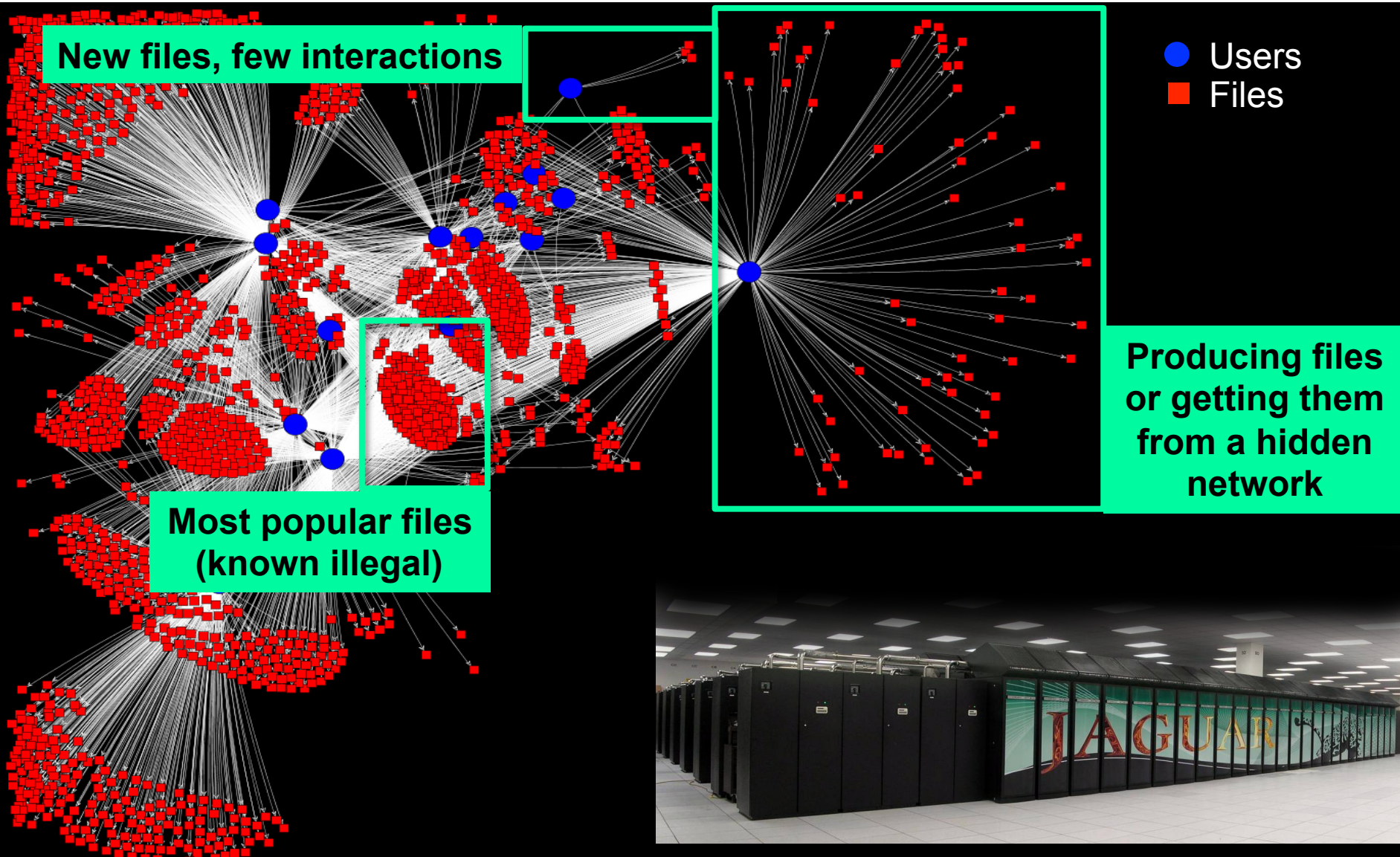- Social networks can be built to see who is sharing from whom

File Sharing Software

Social Network Graph

# Use supercomputers to find the worst of the worst



New files, few interactions

Most popular files (known illegal)

Producing files or getting them from a hidden network

Users
Files

# Summary

- **Difficult national issues**

- **The challenge is in analyzing large volumes of unstructured data**

- **We are proven experts in analyzing large volumes of unstructured data**

- **Piranha and ORCA provides new ways of analyzing and correlating data**

- **This capability is applicable to any intelligence gathering and analyzing activity**

# Contact

## Thomas E. Potok, Ph.D

**Group Leader**
**Applied Software Engineering Research**
**potokte@ornl.gov**

OAK RIDGE
National Laboratory