# The Big Data Paradigm Shift

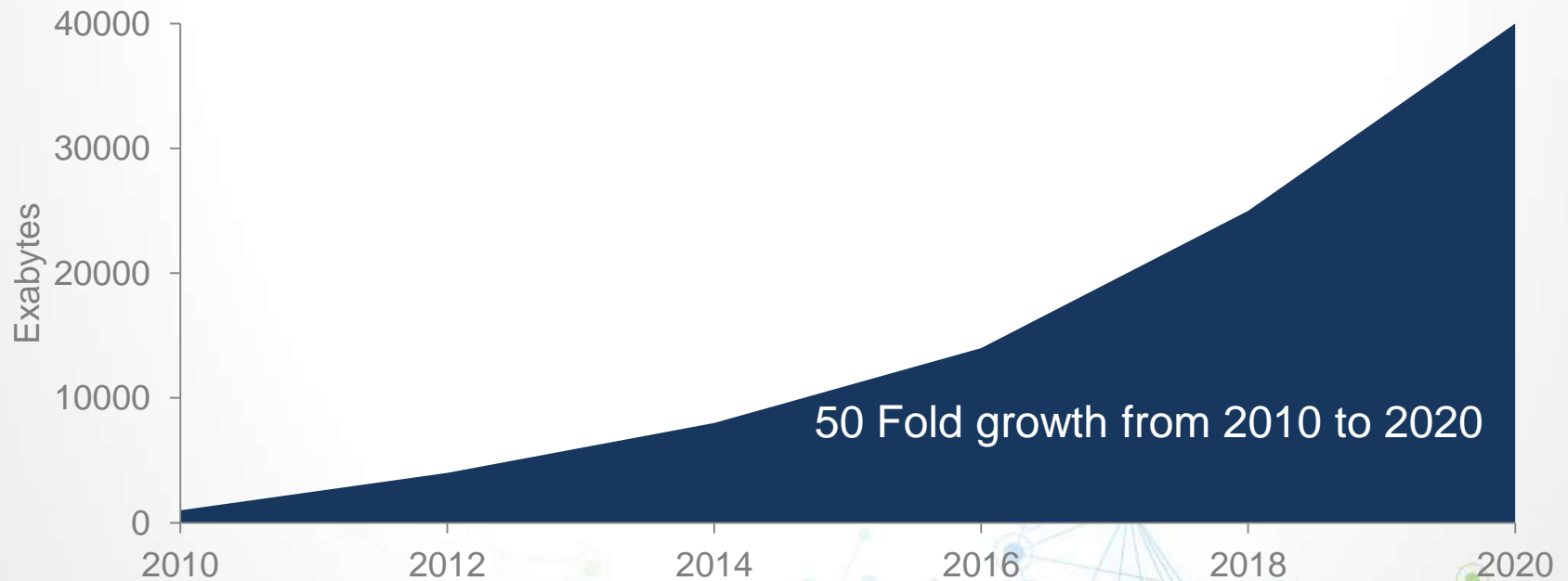## Insight Through Automation

emcien™

# Agenda

- The Problem
- Emcien's Solution:
  - Algorithms solve data related business problems
- How Does the Technology Work?
- Case Studies

# The Problem

- Data is growing at an unprecedented rate
- **Less than 1%** of data is analyzed

50 Fold growth from 2010 to 2020

Exabytes

40000
30000
20000
10000
0

2010    2012    2014    2016    2018    2020

Source: IDC's Digital Universe Study, sponsored by EMC, December 2012

emcien™

www.emcien.com | 3

# Old Paradigm: Manually Intensive Analysis



Unpredictable

Slow

Expensive

Collect → Analyze → Report

www.emcien.com |

emcien™

# New Paradigm: Automation of Analysis



**PREDICTABLE**

**FAST**

**ECONOMICAL**

Collect → Solve → Review & Act

emcien™

**www.emcien.com** | 5
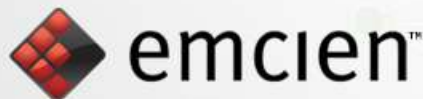
# Emcien's Unique Value Proposition

Emcien's **automatic pattern-detection platform** delivers timely mission critical insights from data

- Automated analysis for fast, predictable, accurate insight

- Applicable across all data types:

    Structured & Unstructured data, Text or Numeric

- Algorithms designed to <u>solve business problems</u>

emcien™

# Types of Data: Structured, Unstructured, Static, Streaming…

**Social, Blogs, Newsfeeds**

**Email Data**

**Click stream data**

**Machine  Data Network log files**

**Marketing Data**

**Sales Data**

**Corporate Data**

**www.emcien.com**  | 7

# Limitations of Current Solutions

**Manually Intensive**

- Very slow and unreliable
- Search or query based
- Visualization as a means for discovery → High error

**Only certain data types**

- Numerical analysis only
- Text only, NLP methods, very high set up cost

**Data staging**

- Streaming data and recent analysis
- At-rest data and historic analysis

**Lack of Scalability**

Current approaches focus too much on storage methods

**www.emcien.com**

# Another View of the Big Data Stack

Our Focus

| Value Layer |
| :-: |

**Use Cases**
**(Industry specific or cross industry)**

| Sectors |
| :-: |

| Banking | Insurance | Manufacturing | Retail |
| :-: | :-: | :-: | :-: |
| Internet | Telco | Intel/Security | Healthcare |

| Analysis Layer |
| :-: |

Algorithms and Analytics

| Infrastructure |
| :-: |

Processor

Storage

| Data |
| :-: |

Click stream

Machine Data

Corporate Data

# How Does Our Solution Work?

- Big Data problems need graph analysis
- Framework for analyzing relationships
- Highly scalable representation

> Data values → Tokens → Graph
> Tokens are linked if they occur together

Unstructured



Structured

# Algorithms Solve to Extract <u>Patterns</u>

- Algorithms surface the highly relevant dependencies
  - Defocus the redundant/noise to surface the signal

# Data Patterns That Reveal "The Insight"

Algorithms designed to reveal graph constructs that solve a business problem

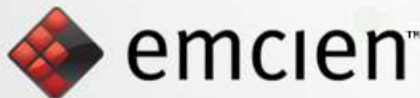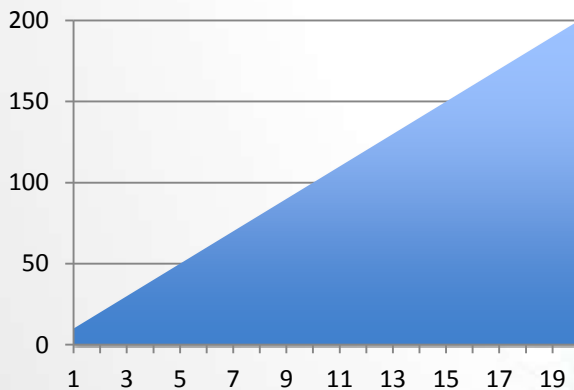| Solving a Graph Problem ➡ | Results in Solving a Business Problem |
|---|---|
| **Loosely Federated Communities** | - Reveals <u>groups that behave similarly</u><br>- Reveals dimensions that bind the group<br>- Impossible to detect in a typical querying system |
| **Cliques** | - Highly correlated elements<br>- Optimal query that would lead to insight |
| **People Network** | - Reveals influence network of individual<br>- Highly predictive for adoption behaviors |
| **Substitute Nodes** | - Nodes that behave very similarly<br>- ID theft or product substitutes |

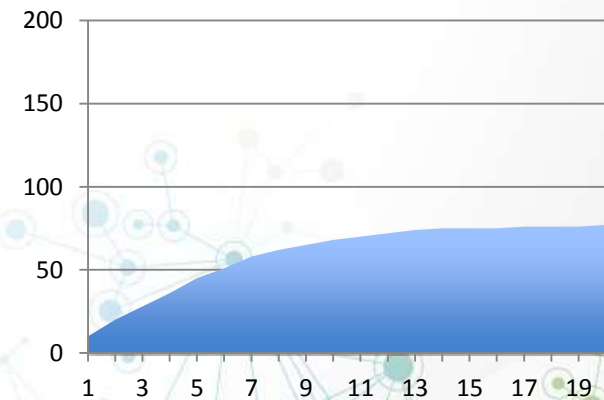# Algorithms Are Highly Scalable For Big Data

## Traditional Data Storage:

- Linear growth with transactions
- Very large storage requirements are
- Increases response time

## Graph Data Storage:

- Size of total number of entities
- E.g. Store has 500,000 items → graph has 500,000 nodes
- Weights updated with transactions
- Delivers a global view of the data

# Speed of Data → Answers
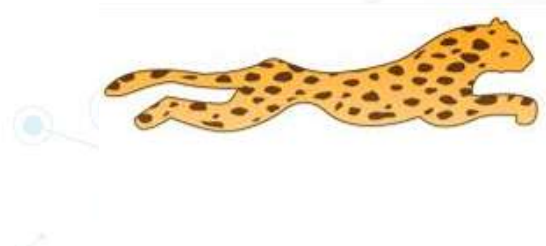# Access Time vs. Processing Speed

## Traditional Data Storage:

- Limit is query speed
- In-memory, hadoop cluster approaches
- Highly dimensional data is a problem
- Unstructured content is a problem

## Graph Data Storage:

- All results are Pre-computed (like Google)
- Pre-computing speed: 50K trans/sec compute on 1-core 8GB RAM system
- Speed of response is "access time"

emcien™

**www.emcien.com** | 14

# The Business Problems We Solve Across Different Types of Data

## Automatically Extract Dependencies

- Web click-stream - Reveal click patterns & market segments
- Sales data - Reveal consumption patterns and propensity
- Clinical trials - extract hypothesis for testing

## Social Patterns & People Network

- Marketing - Reveal conversation patterns, people communities
- For Intel - Reveal bad actors based on conversation patterns

## Surprising Streaming Content

- Machine Network traffic – Reveal network intrusion
- Sales transactions – Reveal fraud based on unusual patterns

## Entity Resolution / Cleansing

- Patterns automatically clusters similar entities.
- Example – credit card transactions, insurance claims with varying merchant names

# Cyber Threat Monitoring with Open Source Data
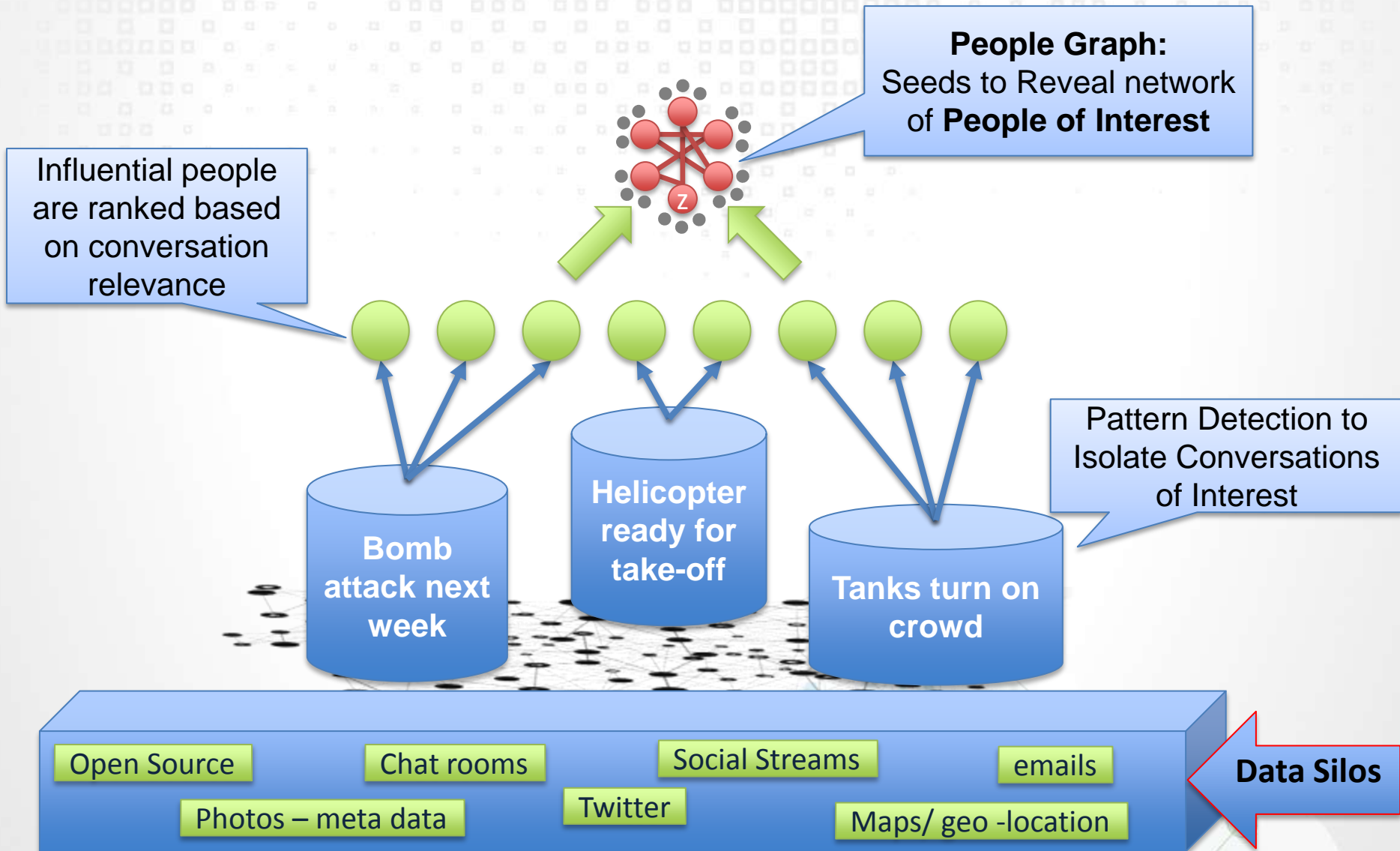
**Customer Overview And Current Situation**
- Federal agency is failing to keep up with the activity and data in open source
- Open Source (social, IRC, blogs, etc.) are a key source of communication for underworld
- Link analysis leads to **people of interest network** – which is key for intelligence

**Customer Objective**
- Federal agency requires fast methods to process high volume open source data
- Need automated methods to highlight **conversations of interest**
- Need automated link analysis to focus on **people of interest**
- Fast and continuous data processing to **keep up with the speed of crime**

## emcien™

**www.emcien.com**

**People Graph:**
Seeds to Reveal network
of **People of Interest**

Influential people
are ranked based
on conversation
relevance

Pattern Detection to
Isolate Conversations
of Interest

**Bomb
attack next
week**

**Helicopter
ready for
take-off**

**Tanks turn on
crowd**

Open Source

Chat rooms

Social Streams

emails

Photos – meta data

Twitter

Maps/ geo -location

**Data Silos**

emcien

www.emcien.com | 17

Influential people ranked based on conversations



Overview
**250,000** Accounts
**Over 1,000,000** Connections

**Highly relevant" People of Interest" network**

# Cyber Threat Monitoring with Open Source Data

**The silent signal – Automatically detecting a sleeper cell**



Overview
**250,000** Accounts Analyzed
**Over 1,000,000** Connections
**1** Account of Interest

emcien™

©2013 Emcien, Inc. All

# *Network Traffic Log Files (1/6)*
# Revealing Patterns In Machine-to-Machine Data

**Customer Overview**

- Research Institute has thousands of users on their network
- Must provide controlled safe access for the internal working labs and the outside network
- Control illegal intrusions, malicious malware and illegal data transmissions

**Customer Objective – Automate Process of Intrusion Detection**

- Scan streaming machine-to-machine log file output
- Detect surprising/interesting anomalies/beacons
- Automatically send short list of top ranked "questions" to ask of the data into existing tools (such as CA, Sumologic, Splunk, etc.)



**emcien**™

**www.emcien.com**

# Example Use Cases

**Example Use Cases**

1. Summarize and Rank Log File data based on "Surprising flow patterns"

2. Determine Machine network based on flow patterns.
   - Rank Machines based on their "influence" in the network

3. Detect "communities of machines" based on how they "talk to each other"

www.emcien.com | 21

# Reveal Surprising Patterns In Network flow Data



Network Log Data

**EmcienScout** Home

## Acme Network Traffic Clusters

0 seconds captured about 26 hours ago
Mon, 3 Dec 2012 9:45 am - 9:45 am

Sorted by relevance

Comments for this Search (0)

**1121** Messages | 11/8/11 | Src 10.2.197.241 | Start Hour:13 | Sig Id:2007757
Signature Name Et Scan W3af User Agent | Priority:2
Classification Attempted Information Leak | Dst 154.241.88.201 | Dst 80

**1441** Messages | 11/10/11 | Src 10.2.198.245 | Start Hour:9 | Sig Id:2003099
Signature Name Et Web Misc Poison Null Byte | Priority:2
Classification Access To A Potentially Vulnerable Web Application
Dst 154.241.88.201 | Dst 80

**6029** Messages | 11/10/11 | Sig Id:2002677
Signature Name Et Scan Nikto Web App Scan In Progress | Priority:1
Classification Web Application Attack | Dst 154.241.88.201 | Dst 80

**4950** Messages | 11/10/11 | Src 10.2.197.245 | Start Hour:9 | Priority:1
Classification Web Application Attack | Dst 154.241.88.201 | Dst 80

### Overview
**239,320** Messages to **42** Clusters
**100.0%** Compression

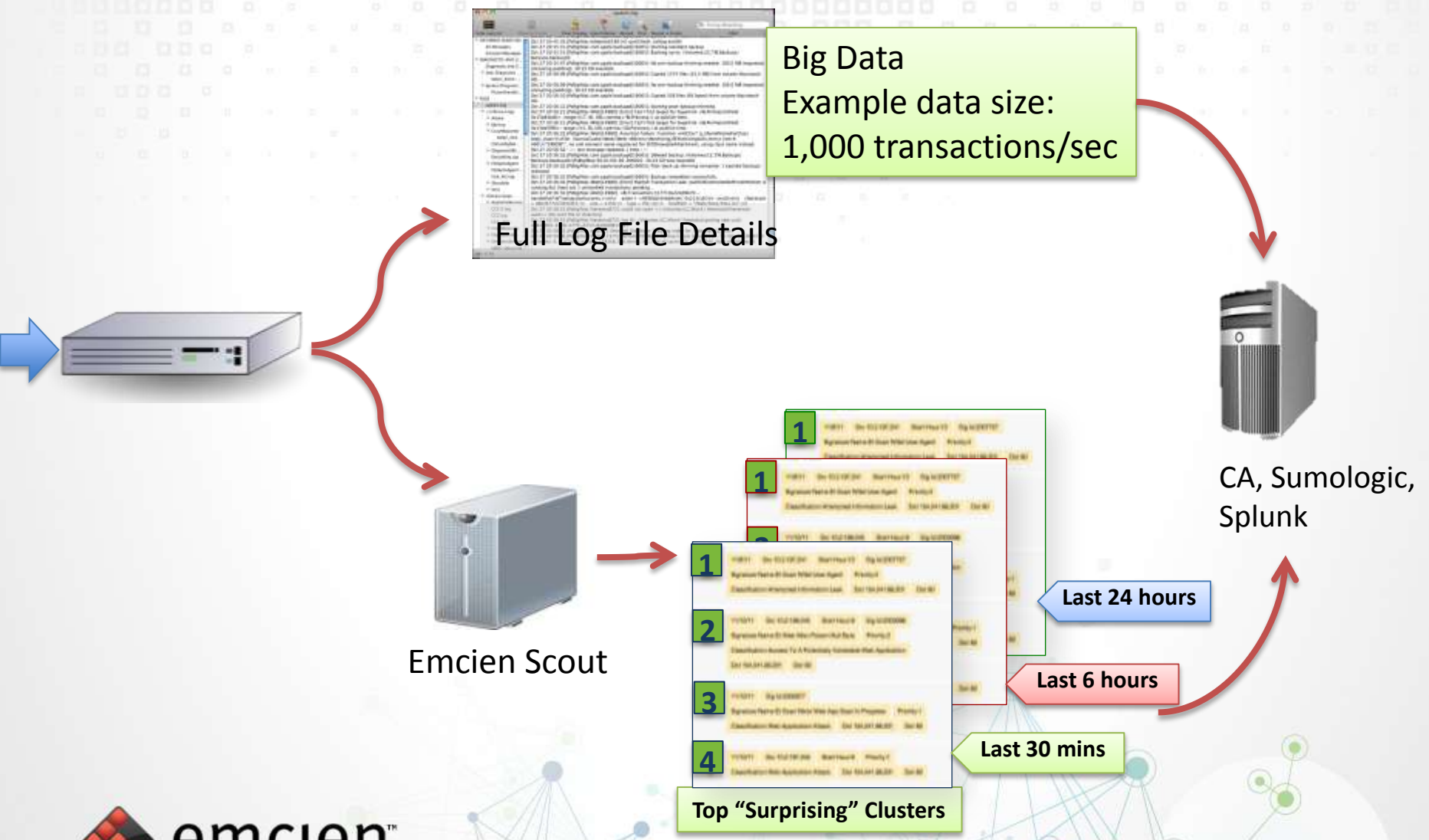This result has 239,320 messages (0 seconds) arranged into 42 Clusters.

**10** Top Words | ‹ Hashtags | Topics ›

Classification Web Application Attack 7887 | Translate
Start_hour:9 7541
Signature_name Et Scan Nikto Web App Scan In Progress
Sig_id:2002677 6029
Src-10.2.197.245 5583
Src-10.2.186.254 2775
Src-0 2442
Dst-0 2442
Src- 2442
Classification Attempted Information Leak 2130

0    1272    2544    3816    5088    6360    7633

Surprising network activity
within the data flow

**emcien™**

www.emcien.com    |    22

# Ranked Summary of "Surprising" Events



Big Data
Example data size:
1,000 transactions/sec

Full Log File Details

CA, Sumologic, Splunk

Emcien Scout

Last 24 hours

Last 6 hours

Last 30 mins

**Top "Surprising" Clusters**

www.emcien.com | 23
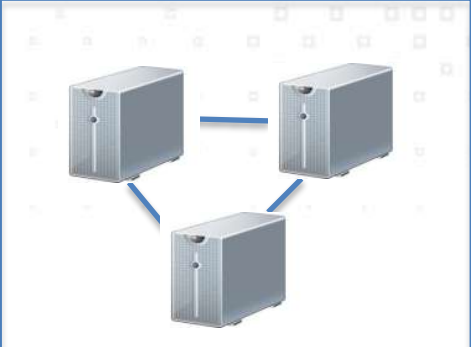
# Most "Influential" Nodes on network

Audience size for each machine



Emcien Scout

Most influential machines on the network based on communication patterns
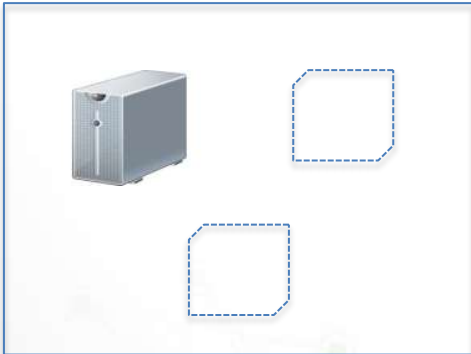
# Machine Communities based on "how they talk"

Lab A          Lab B

**Physical Connections**
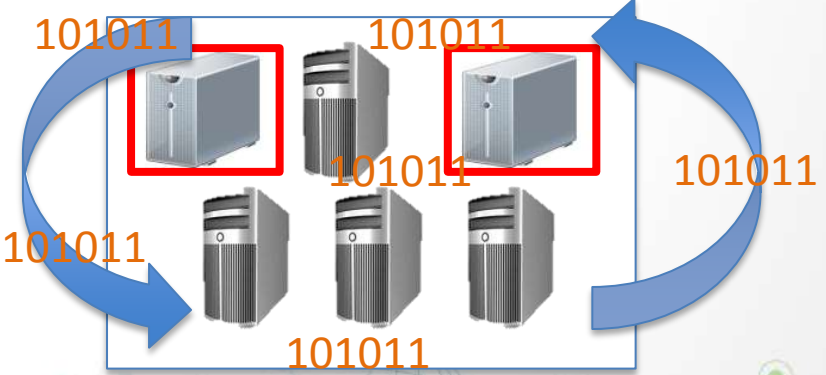


Lab A          Lab B

**Machine Communication Communities**



101011      101011

101011      101011

101011

101011

# Intel Case Study (1/6)
## Reveal Conversation Patterns & Network of Actors in Email Data

**Customer Overview And Current Situation**
- Federal agency is failing to keep up with the activity and data in email
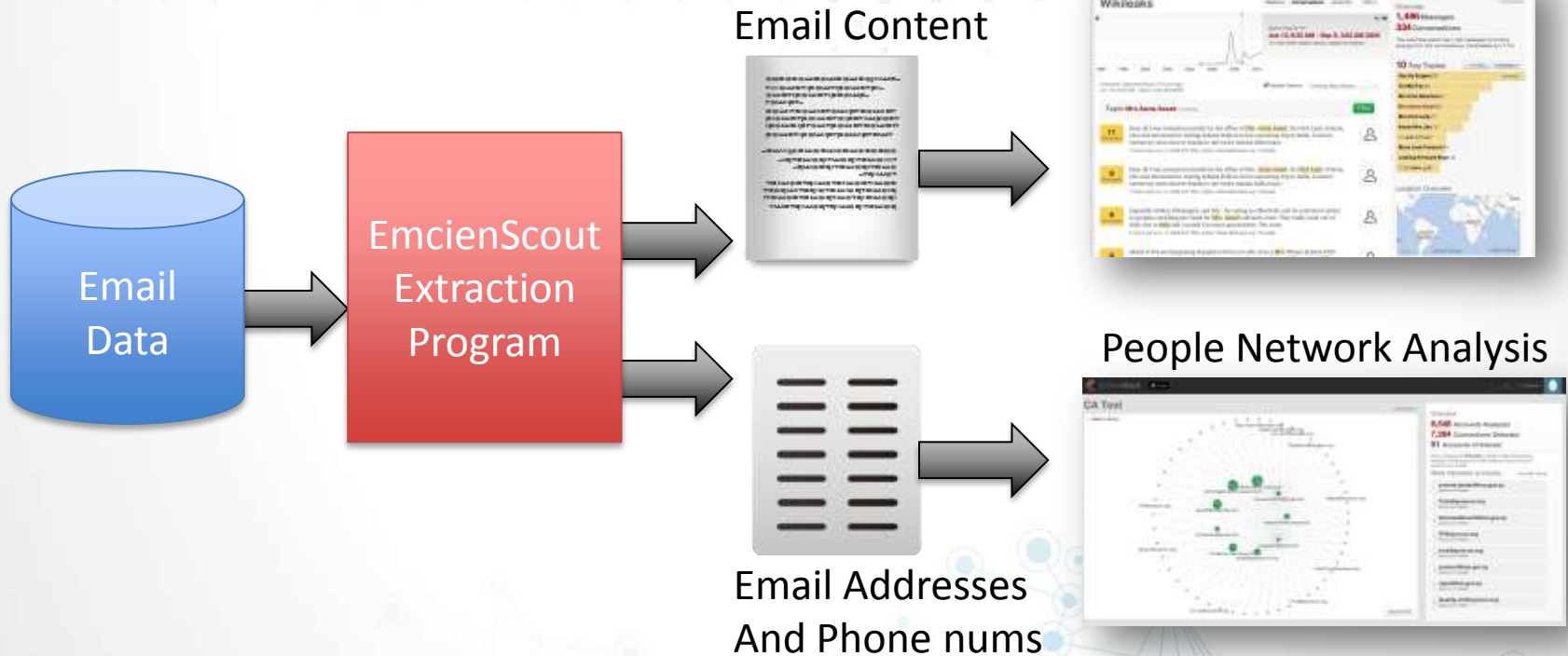- Too much data and current tools are manually intensive

**Customer Objective**
- Federal agency requires fast methods to process high volume of email data
- Need automated methods to highlight **conversations of interest**
- Need automated link analysis to focus on **people of interest based on emails**
- Fast and continuous data processing to **keep up with the speed of crime**

**www.emcien.com**

## Automatic Data Collectors

- Content extracted from emails
- Addresses extracted and linked

Text Summary

Email Content

Email
Data

EmcienScout
Extraction
Program

People Network Analysis

Email Addresses
And Phone nums

emcien™

**www.emcien.com** | 27

## Automatic Email Extraction

From: daniel.brown@enron.com
To: dan.leff@enron.com, david.delaney@enron.com
Subject: FW: EES Employee Issues
Cc: kalen.pieper@enron.com, judy.gray@enron.com
Bcc: kalen.pieper@enron.com, judy.gray@enron.com
Date: Wed, 12 Dec 2001 09:28:51 -0800 (PST)

Dan/Dave,

We are working to gather as much information as possible on our exposure to relocated former and current domestic and international employees impacted by Enron's bankruptcy filing. Lloyd has outlined our position on the urgent issues below. Please keep in mind that regardless of our obligation, the courts have only approved $15K per employee for all expenses less the $4500 payment if applicable.

We will continue to work on getting a comprehensive listing over the next couple of days.

Daniel

**Extracts all Addresses in Header AND Body**

From: daniel.brown@enron.com
To: dan.leff@enron.com, david.delaney@enron.com
Subject: **FW: EES Employee Issues**
Cc: kalen.pieper@enron.com, judy.gray@enron.com
Bcc: kalen.pieper@enron.com, judy.gray@enron.com

Messages extracted, each word tokenized and connected into graph.

We are working to gather as much information as possible....

| We | are | working | to | gather |

# Intel Case Study (4/6)
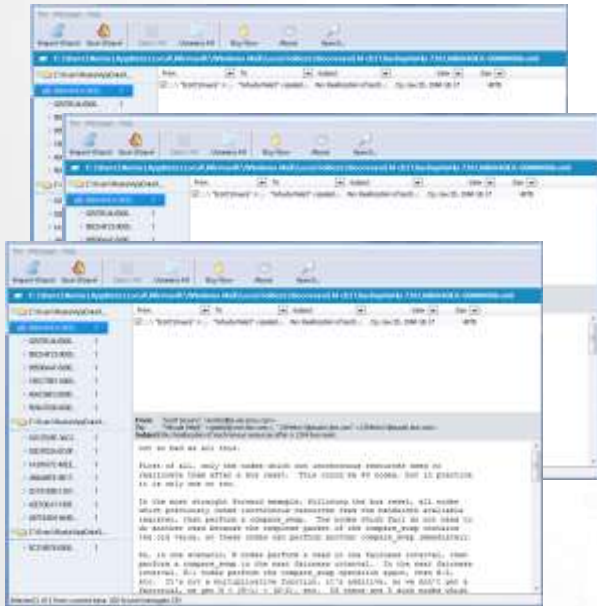## Automatic Email Summarization



Summarize content from emails to better understand group conversations

# Intel Case Study (5/6)
## People Graph (1/2)

Program extracts
To/From email addresses
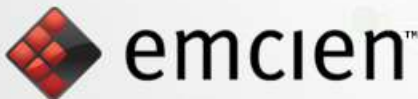__and__ phone numbers
from suspects email
account

Newly created contacts
file is loaded into Scout
People Graph

Large complex graph is
created using emails and
phone number
connections

joe@hotmail.com

jane@ofc.com

404 555-1212

emcien™

# Intel Case Study (6/6)
## Algorithm Computes People Graph (2/2)



Initial "bad actor" seed accounts (emails addresses or phone numbers) are selected or entered

Ranked List of other accounts that are likely involved with seed accounts based on their connections.

# How Emcien Fits Into Your Ecosystem

Feed downstream systems with data output

API

UI

UI for Analyst who wants to review results

Production Servers

**emcien™**

Pattern Detection Platform

Compressed Graph Data Representation

- Sales Trans
- Bank Trans
- Insurance Claims

- Social Media
- News / Blogs
- Emails / Chat

- Server Logs
- Web logs
- Security Logs

Structured Data

Unstructured Data

Unstructured
(Machine Data)

**emcien™**

**www.emcien.com** | 32

# Types of Data

- Many types of Data
  - Structured, Unstructured
  - Text, Numeric, Machine

- In many states
  - Static (slow batch)
  - Streaming or fast batch

**Marketing Data**  **Sales Data**  **Corporate Data**

**Social, Blogs, Newsfeeds**

**Email Data**

**Machine Data Network log files**

**Click stream data**

emcien™

# Limitations of Current Solutions

- <u>Manually Intensive</u>: Very slow and unreliable
  - Search or query based
  - Visualization as a means for discovery → High error

- <u>Limitation based on data types</u>
  - Numerical analysis only
  - Text only, NLP methods, very high set up cost

- <u>Limitation based on data staging</u>
  - Streaming data and recent analysis
  - At-rest data and Historic analysis

- <u>Scalability</u>
  - Current approaches focus too much on storage methods

emcien™

**www.emcien.com** | 35