

# Recurrent Patterns Detection Technology

## White Paper

*November, 2007*



## **RPD™ Technology**

### **Network Based Protection Against Email-Borne Threats**

Spam, Phishing and email-borne Malware such as Viruses and Worms are most often released in vast quantities in a relatively short period of time, causing global damage of tens of billions of dollars annually. Although many methods to combat these threats have been developed and tested throughout the past several years, a common deficiency of most of these methods is that they lack the ability to adapt quickly enough to the rapid change of distribution and infiltration techniques invented by Spammers and Virus authors.

Recurrent Patterns Detection (RPD) technology, based on the identification and classification of message patterns, delivers the highest threat detection capabilities. The objective of this document is to discuss the characteristics of these threats, the challenges facing technologies that aim to mitigate these often malicious attacks, and describe how the RPD solution protects against all types of email-borne threats.

The RPD approach is based on the understanding that all threat outbreaks share some common characteristics, including:

- Most email messages within the outbreak have been altered to make it difficult to set blocking rules based on lexical analysis.
- Most outbreaks include millions of email messages to maximize the highest possible response rate and the greatest ROI for the attacker.
- Most outbreaks are released within a short period of time, requiring a real-time solution to detect the outbreak to limit or avoid the damage that can be incurred.
- The originators of the attacks invest heavily in disguising their origin to make it difficult to track the message back to them.

In more detail, the following challenges have been identified:

#### **Spam Detection**

Spammers typically produce Spam for commercial gain rather than malicious intent. Although the purpose may be to sell a product or service to recipients, Spammers often use unethical methods to acquire email addresses of recipients and to distribute mass mailings. In composing Spam messages, Spammers use sophisticated tactics to evade existing Spam detection applications. This includes covering the tracks to the Spammers, manipulating or hiding the commercial URLs, use of non-English words and phrases and a host of other methods.

Typically a massive Spam outbreak only lasts a few hours and is launched from a network of 'Zombie' machines. To complicate the detection process, each message within the massive Spam outbreak can be composed differently and employ more than one evasion technique.

#### **Phishing Detection**

Password Harvesting or Phishing messages are typically sent for the single purpose of identity theft. There is an underlying intent within each message to violate the privacy of the recipient and commit fraud. The goal of a Phishing message is to fool the user to go to a website which

requests them to enter private information. By using highly effective social engineering methods, these messages target users, often with a sense of urgency, to believe they have arrived at a legitimate site such as their bank or recognized online vendor. The sender is then able to gather personal information such as credit card numbers, passwords, social security or identification numbers and other private information.

Phishing messages appear to be from genuine or credible sources. Like Spam, Phishing messages can be sent in any language or format in attacks that typically last only a few hours and are usually launched from an army of 'zombie' machines on the internet.

### **Virus Outbreak Detection**

Typically, email-borne Virus or Worm outbreaks are created and released for the malicious purpose of either committing fraud, industrial espionage, or in the hope of gaining financial benefits by gathering private information that can be sold to Spammers or those with criminal intent. Like Spam and Phishing messages, each Virus message can be packed differently in terms of its content and the characteristics of the executable files that include the Virus. However, email-borne Viruses and in particular, Worms, can be received from legitimate and trusted email sources that might have been previously infected and were unintentionally distributing the Virus to others.

Like Spam and Phishing, email-borne Virus attacks often last for very short durations. In the case of Viruses, users are exposed and unprotected during the first hours of the attack because most Anti-Virus defenses depend heavily on the use of a database of signatures that identify the threat by matching it with already-known characteristics. Recently, Virus writers have become even more sophisticated and have begun distributing multiple modified instances of the same Virus within a series of outbreaks to evade heuristic systems and to maximize the impact before new signatures are propagated.

### **Message Patterns**

Massive outbreaks which distribute Spam, Phishing, and email-borne Viruses or Worms, consist of many millions of messages intentionally composed differently in order to evade commonly-used filters. Nonetheless, all messages within the same outbreak share at least one and often more than one unique, identifiable value which can be used to distinguish the outbreak.

For example, in the case of Spam the objective is to lead the recipient to the same commercial websites that can be classified as Spam. In doing so, different Spam attacks are often launched from the same network of zombie machines that can be blacklisted. In the case of Phishing, recipients are lured to voluntarily disclose personal and confidential information via clever social engineering methods and the objective is often to lead the victims to the same fake URLs. Email-borne Viruses always contain the same malicious code (otherwise it is a different Virus or another instance of the same Virus). All these are recurring values of typical outbreaks. These values are called the 'message patterns' of the outbreak. Any message containing one or more of these unique patterns can be assumed with a great deal of certainty to be part of the same outbreak.

Message patterns are extracted from the message envelope, headers, and body with no reference to the lexical meaning of the content. Thus pattern analysis can be used to identify outbreaks in any language, message format, and encoding type. Message patterns can be

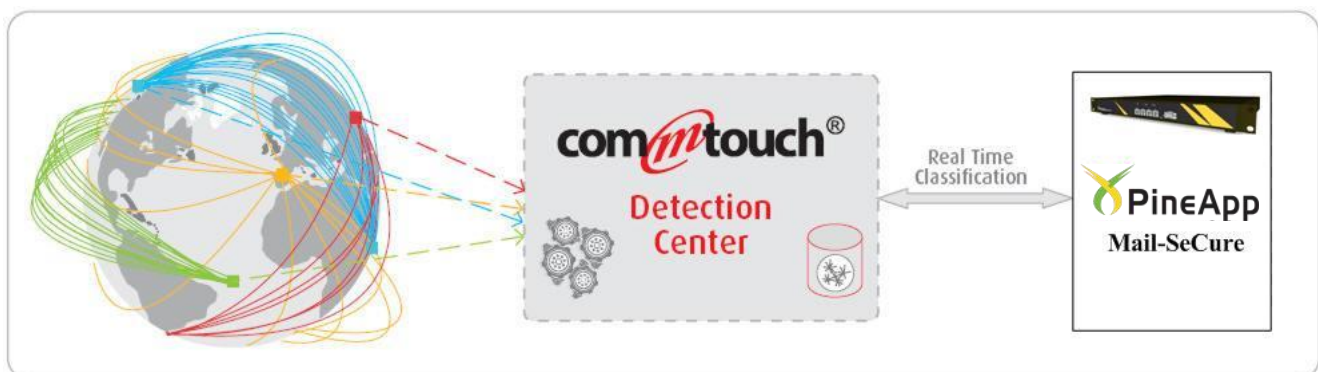
divided into distribution patterns, which determine if the message is 'good' or 'bad' by analyzing the way it is distributed to the recipients, and structure patterns, which determine the volume of the distribution.

The challenges of message pattern classification determining which message patterns identify outbreaks without generating cases of false positives, and how to extract and analyze these patterns before the outbreak wanes. Most outbreaks have a relatively short lifecycle measured in only a few hours. Therefore, any solution that does not detect and classify messages in real-time will only be effective towards the end of the outbreak, when most of the damage has already been done. All outbreaks attempt to disguise messages as legitimate email correspondence pretending to arrive from trusted sources and therefore, solutions that are based on pattern analysis must be able to tell the difference between 'good' and 'bad' patterns and avoid making mistakes.

The challenges are made more complex by the fact that each new outbreak usually introduces completely new patterns that were not previously analyzed and are therefore unknown to the pattern analyzer. Pattern detection represents a new and greater understanding of how email-borne threats are created and propagated. Because tactics for distributing Spam, Phishing, and email-borne Viruses and Worms are constantly evolving, it is necessary to proactively identify new and unique patterns in real-time in order to determine new outbreaks as they are released to the internet and begin targeting recipients.

## Recurrent-Pattern Detection (RPD™) Technology

Recurrent Pattern Detection (RPD) technology detects and classifies all types of email-borne threat patterns in real-time. RPD is hosted by the Commtouch® Detection Center, which proactively analyzes vast amounts of Internet traffic in real-time.



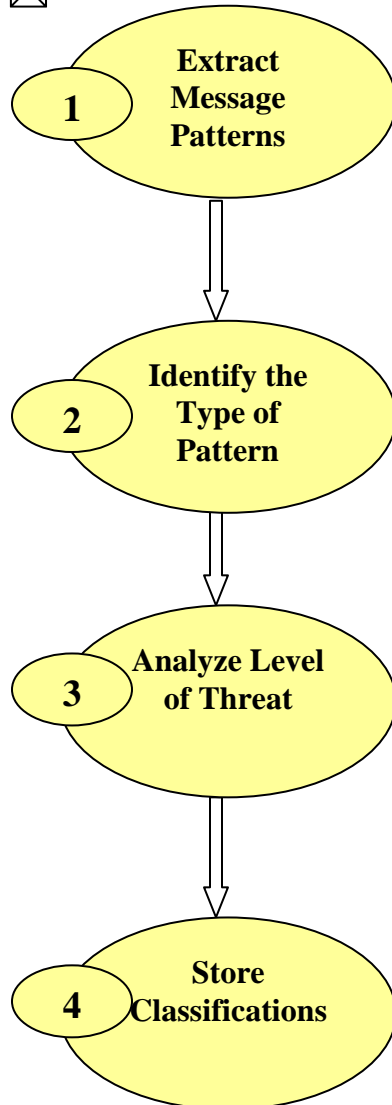
RPD, a patent-pending technology based on Commtouch's U.S. patent, extracts and then analyzes relevant message patterns, which are used to identify massive email borne outbreaks. RPD classifies both distribution patterns and structure patterns and the analysis results are stored in a vast warehouse of classifications. In addition to identifying new threat patterns, RPD is also used to modify or enhance classifications of already identified message patterns.

RPD is designed to distinguish between the distribution patterns of solicited bulk emails which represent legitimate business correspondence, from those of unsolicited bulk emails by applying a reverse analysis. The results of this analysis are 'bleached' message patterns belonging to 'good' messages such as popular newsletters, mailing lists, etc.

The RPD technology is used in a highly scalable environment to deliver extremely high performance rates by analyzing many millions of new patterns each day, (24x7x365). New outbreaks are identified within minutes since they are launched on the internet globally.

The RPD technology was designed to be fully automated and requires no human intervention. To ensure maximum privacy and business confidentiality, RPD was designed to analyze hashed values of message patterns and not the 'open' values nor the message content.

## Email Messages to Analyze



Unsorted flow of millions of messages each day consisting of legitimate emails, Spam messages, Phishing, Email-borne Viruses and Worms.

- **Distribution Patterns**
- **Structure Patterns**

- **"Good" patterns of solicited bulk emails**
- **"Bad" patterns of unsolicited bulk emails**

- **No threat**
- **Confirmed threat**
- **Bulk or suspected threat**
- **Unknown**

### Used for:

- **Spam detection**
- **Phishing detection**
- **New Virus outbreak detection**

RPD identifies nearly 100% of incoming threat messages with almost no cases of false positives. It is language-agnostic and is equally effective for all message formats and encoding types.

## Summary

To effectively combat email-borne threats, a successful solution must address a growing number of challenges. RPD is a proactive detection technology that continues to outwit those who continue to invent new methods to propagate email-borne threats because it does not rely on the contents of the email and therefore, is able to detect Spam in any language and in every message format (including images, HTML, etc.), non-English characters, single and double byte, etc.

RPD is also unaffected by the lack of signatures for new email-borne Viruses and Worms and is capable of detecting these Malwares within minutes of the attack's launch.

RPD technology offers:

- High Spam detection rate with almost no cases of false positives
- Early detection of Virus threats
- Protection against Phishing attempts
- Content-agnostic threat protection
- Multi-language threat detection
- Multi-format threat detection

For these reasons, the RPD technology is the best companion to multi-prone messaging and security applications. RPD is available for fast integration in various ways, the most common of which is to embed the cross-platform shared object managed with a simple-to-use API.

Because RPD uses pattern analysis, it provides the best protection of investment for enterprises and vendors of messaging and security applications.