

# The Network Factor in Proactive Digital Evidence Acquisition

Irvin Homem, Spyridon Dosis, Oliver Popov  
*Stockholm University, Sweden*

## Abstract

*While cybercrime proliferates – becoming more complex and surreptitious – the tools and techniques used in performing digital investigations are still seriously lagging behind, effectively slowing down law enforcement agencies at large. In this paper we briefly describe the architecture of a comprehensive proactive digital investigation system that is termed as the Live Evidence Information Aggregator (LEIA). Particular focus is made on the importance of the efficiency of network communication in such a system in order to retrieve potentially evidentiary information remotely and with immediacy. We demonstrate the live evidence capturing capabilities of such a system highlighting the necessity for better throughput through the use of Peer-to-Peer Systems.*

**Keywords:** *Digital Forensics, Cybercrime, Digital Evidence Proactive Forensics, , Big Data, Hadoop, Hypervisors, P2P, Collaborative Live Investigation*

## 1. Introduction

Malevolent activities, previously a characteristic only of the physical world, quickly adapt and evolve according to their environments. Similarly in the “Digital Realm they have also been seen to readily adapt, taking up their niche markedly on the inter-webs. They are no longer simply rare occurrences with mild consequences. They have permanently set up camp in intricate and surreptitious forms, taking unjust advantage over unsuspecting users going about commonplace activities on the Internet. Such activities as learning, socializing, shopping, seeking medical treatment and even just performing basic communication are now at the risk of being compromised by malicious actors aiming to steal, sabotage or paralyze others’ activities for personal benefit.

The consequences of such malicious activities for the unsuspecting user have also become more detrimental, persistent and having far reaching effects in that they are largely untraceable and easily invisible to the untrained eye. Developing novel and innovative methods that enable them to remain effectively undetected and untraceable is the hallmark of these pseudo-anonymous evildoers. They are almost always one step ahead of the pursuers. Furthermore it is relatively easy to hide among the deluge of data that is created among the communication devices that support the basic

network communication on the internet, thus malevolent activity in the “Digital Realm” can easily become rampant and uncontrollable if there are no equally innovative methods to counter the offending actors and their activities. The rate of creation of and uptake of novel techniques by law enforcement agencies must at the very least be equivalent to that of their criminal counterparts, if they are to keep up with the proliferation of crime on the Internet.

One of the foremost areas in digital crime investigations where innovative means of combatting crime are highly necessary but largely lacking is the evidence capture process. This is the initial stage of an investigation where artifacts from the scene of the crime need to be retrieved in their original form, or, in the case of digital investigations, in some form of a copy of the original artifact that can be proven to be devoid of any tampering. This process needs to be performed meticulously, carefully and in many cases slowly in order to ensure that there is no potentially crucial piece of evidence left behind. This is the state of affairs in the real physical world.

However, today’s crime scene is rapidly edging away from a physical reality into a more virtual one. The forms of evidence found in these “Digital Crime Scenes” have also moved from the traditional fingerprints, footprints, hair samples, blood samples or other DNA related evidence, into more digital artifacts.. Such digital forms of evidence commonly include hard-disk drives, live (RAM) memory, network traffic captures, mobile devices, RAID sets [1], and virtually any other form of technology that records past events of its actions; that can be captured and can be analyzed during or after the criminal event and whose integrity can be verified.

This opens the floor to almost any form of computer appliance (physical or virtual) that can be thought of. Thus arises the heterogeneity problem among devices – or simply put the seeming lack of standardization among vendors of devices that perform related tasks. Different devices may have different physical connectors, operating systems, software applications, storage formats, encoding formats and communication protocols [2]. This heterogeneity makes the job of a Digital Investigator a lot more difficult because of the wide variety in which evidence could manifest itself in the wild. This greatly hampers any manual efforts of collecting evidence even with the assistance of semi-automated tools of today such as disk imagers.

In addition to this, Electronic Crime cases today often involve more than just a single device. Several computer-like appliances including tablets, mobile phones, digital cameras, GPS devices, smart-TV's and even embedded devices such as onboard vehicle computer systems (from trucks, cars and even ships) could be seized for a single case, in order to be subjected to further investigative analysis. If we also bring in the vast realm of the Internet also into play, such evidence sources could include web application accounts, online email accounts, cloud storage facilities, network traffic captures and logs [3]. It is not difficult to imagine that all these evidence forms could easily be part of a single case in today's world and even more so in the imminent realm of the Internet of Things. The sheer volume of data that one would have to sift through in order to investigate a single case could be in the order of Terabytes and can be a more than daunting task to perform. [4]

Furthermore, in the realm of the Internet, composed massively interconnected devices sharing vast amounts of highly varying data, crossing paths at high velocities, the speed of the capture of potentially evidentiary information is of essence. The same levels of meticulousness and carefulness of physical evidence acquisition may as well be sacrificed to some extent for the agility that is needed in reacting to crime in the digital world. This is because potentially evidentiary information that is not captured almost instantaneously, is likely to be lost forever in just a matter of seconds. However, this does not mean that all accuracy and care in collection of digital evidence artifacts is ignored, rather it is traded-off and reduced in favour of speed. Nevertheless, the maintenance of the chain of custody is always very important in any digital investigation. New methods of achieving similar standards of the preservation of digital evidence to those of physical evidence also need to be sought after and integrated into legal standards.

Finally, at present, investigators grapple with the problem of the relatively immature forensic tools that they are presented with. Current industry standard forensic tools such as EnCase, FTK, XRY, Volatility and Wireshark<sup>1</sup>, at the moment of writing, do not cater for the highly divergent nature of digital evidence sources. Most, if not all tools, focus on a single niche area such as Filesystem Data, Live Memory, Network Traffic, Mobile Devices or Log data. None of the tools provide a comprehensive method to interface with all the variety present to provide a uniform investigation platform. In addition to this, current tools have rather limited capabilities for capturing potentially evidentiary data on demand over networks as well as dealing with extremely

<sup>1</sup> EnCase® Forensic and FTK® (Forensic ToolKit) are traditional filesystem forensics tools. MicroSystemation XRY is for mobile forensics, Volatility for live memory forensics and Wireshark for network forensics.

large datasets. Most of the tools would struggle and would quickly become problematic when presented with Internet-Scale crime scenes.

In this paper, we present the architecture of a scalable, distributed, multi-component incident response and digital investigation platform aimed at dealing with large scale distributed cybercrime investigations. We name this system the Live Evidence Information Aggregator, or LEIA, in short. The LEIA architecture aims at curbing cybercrime through assisting law enforcement agencies in improving their digital crime reaction capabilities. This is to be done through addressing several of the aforementioned problems such as the innate and growing complexity of the fast growing "Internet-of-Things" types of cases as well as dealing with the constantly growing amounts of heterogeneous data vis-a-vis the present shortage of physical resources and technical capacity within Law Enforcement. We also address the need for proactive collection of evidence from potential evidence sources on-demand over networks, and further show the need for faster throughput network transfers such as those seen in Peer to Peer technologies. The rest of this paper is organized as follows: In Section 2, we review related work outlining shortcomings of previous similar solutions. Section 3 describes the requirements for a comprehensive digital investigation platform. The functionality of the LEIA system with particular focus on the networking component is described in Section 4. The network-focused proof of concept implementation and results are outlined in Section 5. In Section 6 and 7, we summarize the work done in this study and propose further work that may be done in this area, respectively..

## 2. Background and Related Work

Several progressive efforts have been made towards improving the efficiency of the Digital Investigation process. The motivations behind these have spawned mainly from the changing requirements of national and international legal systems, the evolution in the digital crime scene, the visible backlogs of cases overburdening law enforcement agencies and advances in technological capabilities.

Some of these efforts include: Delegation and Collaboration among teams; Reduction of evidence sizes through filtering out known files; and simple automation of important but mundane, repetitive tasks (such as indexing data for subsequent searches, file carving, parsing running process in memory or TCP flows in network captures). Most of these capabilities have been implemented in current industry standard forensic tools, however, investigators and analysts still remain overburdened. This is because of the presently abundant and steadily growing amounts of heterogeneous and

disjointed datasets from multiple sources that they are tasked to collect and analyze. Methods to alleviate this problem through fully automating the remote collection and pre-processing of such data are so far either lacking in efficiency or in scalability.

Several unidirectional solutions have been proposed in a bid to solve this multi-faceted problem, however, none has been unequivocally successful. In recent times there have been initiatives to centralize processing and data collection onto powerful centralized mainframes and data storage centres, respectively. There has also been a push towards having the different parties, involved in solving a case to work together, even from geographically separate locations [5], particularly among technical staff in niche areas (filesystem forensics, network forensics, live memory forensics or mobile forensics) and the legal experts. Collaboration has been the mainstay of the attempt to get cases solved faster.

Reducing the amount of data that is needed to be collected is also a means of reducing the amount of time needed to analyze the data. This has previously been done through “Known File Filtering” as well as through heuristic analysis of the entire captured data. Network Security Monitoring has also been an avenue for gathering data through the assistance of Intrusion Detection Systems (IDS’s) assisted through Artificial Intelligence. However, this has been the specific mandate of the IDS, centralized or distributed, as the case may be, with terminating (end) devices or intermediary devices generally playing very minor roles in this task.

As far as is known to the author, there has not been much done, through any single initiative, in terms of expanding the scope of data captured to be the mandate of all possible devices of reasonable capability. Enabling individual devices to natively act as part of the Incidence Response System, towards the aim of collecting potential evidentiary data, has not been widely studied. Additionally, collaboration on the human processing level has been emphasized, but it has not been introduced among unrelated networked devices. These devices could possibly be harnessed to work together towards aiding in intelligent real-time capturing, filtering and processing in order to attain and retain that which could be considered as possible evidentiary data, antecedent to the event of a crime being detected. It is for these reasons that we delve into this area to explore it further.

Notable related studies include [6], that describes a live network forensics system that provisions varying Intrusion Detection Systems on host machines based on their respective resource costs. It works in a virtualized environment where snapshots are taken periodically and used to revert the system back to the point before an attack began. Each system rollback results in varying IDS’s being deployed to collect new and possibly better

information. This presupposes that the attacker re-enacts their malicious behavior repetitively each time their efforts are thwarted by the system. Storage of the potential evidentiary information in a forensically sound manner is not particularly dealt with in this study. The aim was to understand attacks better in order to make better decisions on what kind of preventive measures to deploy.

In [7]–[9] distributed system architectures for proactive collection and summarization of evidence, with centralized data storage and processing, are described. They are, however, particularly directed at closed domain enterprise systems, where there is some form of control and order instigated by system administrators. Participation of computer systems outside the control of the enterprise is not considered. The system being proposed in this study is aimed at being universal – applying to the entire Internet.

The work done by Redding in [10] is the most closely related study done in the area of pro-active and collaborative computer forensic analysis among heterogeneous systems. Redding proposes a peer-to-peer framework for network monitoring and forensics through which network security events can be collected and shared among the peers. “Analysis, forensic preservation and reporting of related information can be performed using spare CPU cycles,” [10] together with other spare, under-utilized, or unused resources. This system however seems to be designed to collect only network security events and not any other forms of evidence from individual host devices Furthermore it seems to be aimed towards an “administratively closed environment” under the control of some systems administrator within an enterprise. An open system that has the Internet as its domain of operation assisting in the collection of any form of computer based evidence is what is not dealt with in Redding’s work. Thus, it is this that is sought after in the current study as will be described later in this paper.

In order to facilitate uniform, seamless exchange of forensic artifacts between heterogeneous entities, some form of standardization of the transmitted evidence formats is necessary. One of the bodies that has made proposals related to this is the Common Digital Evidence Storage Format Working Group [2]. Other notable efforts include [11] which makes use of the Resource Description Framework (RDF) from Semantic Web technologies as a common data representation layer for digital evidence related metadata, using ontologies for describing the vocabulary related to this data, and [12] where a detailed ontology of Windows Registry artifacts of interests is introduced. The Open Forensic Integration Architecture (FIA) in [3] and FACE [4] describe methods for the integration of digital evidence from multiple evidence sources in a bid to facilitate more efficient analysis. The Advanced

Forensic Format [13], AFF4 [1] and XIRAF [14] describe annotated evidence storage formats that allow for addition of arbitrary metadata as well as interoperability among different tools.

All these evidence integration ideas assume that acquisition of the relevant data from the respective sources has already been performed. This is a stumbling block in the process of making digital investigations more efficient. This is because these methods have not been integrated into the evidence acquisition process. They are seen as a subsequent step rather than part of the process. As part of this study we integrate these hitherto unlinked processes.

The proposed idea that this study covers is composed of several areas of specialization, namely: The Internet of Things (IoT), Intrusion Detection Systems, Peer to Peer Networks, Virtualization infrastructures, Large Scale Cloud storage and Semantic Web technologies. Most of these technologies have been previously harnessed in different capacities, singularly or in small clusters, towards the benefit of digital forensics for today's complex internetworked and intertwined cyber realm. However, to the author's knowledge, there has so far not been any work done that aims to merge all these technologies together in order to provide a singular scalable solution that solves the recurring problems of large amounts of data, several sources of evidence, inability of collecting evidence efficiently over networks, heterogeneity among systems, insufficient processing power, security and privacy – that are constantly troubling digital forensic analysts and law enforcement agencies worldwide.

### 3. Characteristics of the Desired Solution

In light of the current state of electronic crime, the present state of forensic tools, and from experience, we describe below a wish-list of characteristics that one would like to have in a Cyber-Law Enforcement solution:

- *Distribution*: The ability to deal with massive amounts of distribution in terms of participants, data storage, processing and dissemination. The system needs to be able to handle the heterogeneity that may come with distributed systems as well.
- *Scalability*: Large scale interconnectivity, as well as the possibility of new entities joining, as well as others leaving the system dynamically and gracefully without drastic negative effects on the system. The ability to easily improve or extend the capabilities of the system through new modules is also desired.
- *Availability*: Providing suitable levels of functionality as and when required.
- *Universality*: Among the heterogeneity and lack of standardization among vendors of different systems, there needs to be some standardization

and common understanding between the systems on the level of communication and storage of potential evidentiary information.

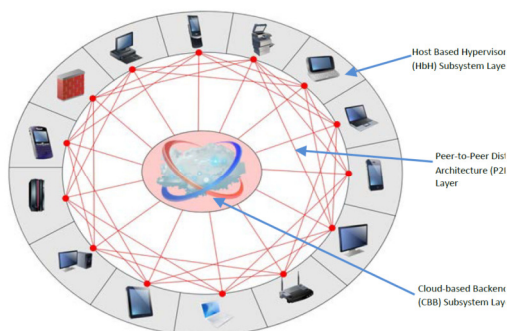
- *Responsiveness*: The system should be able to aptly detect when a security policy has been irrecoverably violated, thus collecting information in order to pursue the perpetrators of the criminal actions. This also improves on efficiency and privacy in that the system does not have to perpetually be collecting all possible information from all possible systems.
- *Resource Sharing*: Today, large complex problems that are being solved through collaboration and sharing of resources as seen in Crowdsourcing, P2P networks, and cloud infrastructures. They provide on demand rapid availability of large amounts of resources from collective resource pools providing speed, efficiency and the benefits from “the wisdom of the crowd”.
- *Integrity (Trust, Reliability & Accuracy)*: As a system facilitating law enforcement in digital crimes, the levels of trust, reliability, accuracy and integrity of the information needs to be high enough to be accepted as a veritable source of evidentiary information for a court of law. The Daubert standards and the chain of custody need to be adhered to.
- *Privacy & Confidentiality*: Personally identifiable and secret information must be maintained as anonymous and confidential as is reasonably acceptable, unless incriminated. Unauthorized access to such information is not to be allowed.
- *Security*: In addition to ensuring the security of the potential evidentiary information that it aims to collect and process, it must also take its own security into consideration – especially in terms of authentication, authorization, accountability and non-repudiation of activities undertaken.

### 4. The Live Evidence Information Aggregator

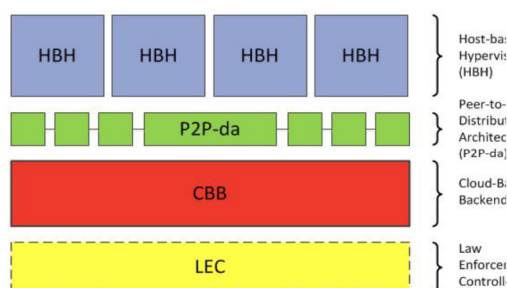
The LEIA is a 4-tiered system architecture that may be described as a combination of hypervisors, intrusion detection systems, peer to peer systems and cloud storage. It is made up of the following components:

- a) The Host-based Hypervisor (HbH)
- b) The Peer-to-Peer Distribution Architecture (P2P-da)
- c) The Cloud-based Backend (CBB)
- d) The Law Enforcement Controller (LEC)

The diagrams that follow portray the architecture of the LEIA. Figure 1 shows 3 layers of the architecture in a figurative manner, while Figure 2 displays the layered nature of the entire architecture.



**Figure 1: High level view of the LEIA architecture**



**Figure 2: Layered/Tiered representation of the LEIA architecture**

The functionality of each of the layers of the LEIA system is briefly described in the following sections.

**4.1 The Host-based Hypervisor (HbH)**

The Host-based Hypervisor (HbH) system is composed of a virtualization layer managed by a hypervisor – a privileged secure platform managing the guest operating system (OS). The hypervisor contains an inbuilt host-based intrusion detection system also termed as the embedded intrusion detection system (em-IDS). Security utilities within the guest OS such as anti-malware tools and intrusion detection systems maintain their own data and logs that are accessible to the HbH. The HbH collects and assimilates the information that it gets from its own inbuilt intrusion detection system together with other information collected from the other security utilities that may exist within the guest OS. This helps in getting a better perspective of current malicious activity that may be underway.

Further to this sharing of information within a single HbH system, individual HbH systems also share their information about malicious activity they may have discovered with each other. This

communication is facilitated through the Peer-to-Peer Distribution Architecture (P2P-da). This collaborative effort among the HbH systems further helps improve the accuracy of IDSs and eventually forensic data acquisition.

In order to reduce the amount of data that may need to be collected for analysis, each HbH maintains a hash list of the local files on its guest operating system (Local - Known Data Hash-List, L-KDHL). This L-KDHL is periodically cross-checked and updated against a Master – Known Data Hash-List (M-KDHL) stored at the Cloud-based Backend (CBB). This is managed by the Cloud-based Backend Differencing Engine (CBB-DE) component of the CBB. The aim of this is to quickly filter out known data or files through matching the files on a HbH against hashes of files that are known to be benign and have not been modified in any way.

A user data profile with its corresponding hash-lists is also created. The user-data hash-list is also maintained in a dual format – with a local copy residing on the HbH and a remote master copy being maintained at the CBB. Further to this the actual user data is backed up at the CBB. Thus, the user data hash lists are used to check which files have changed and may need to be backed up to the CBB.

With respect to “Known Malicious Files” which are files that have been previously identified as having malicious content within them, a “Known Malicious File” hash list is to be maintained only on the CBB. It is not held on individual HbH systems as it can easily become large and unmanageable for a single each HbH to maintain.

The hypervisor is the critical element when it comes to the collection of potential evidentiary data. Having privileged access, the hypervisor is able to directly interact with the file system, network interfaces, memory caches and other low-level resources, which are all primary sources of prime evidentiary data in digital investigations. The embedded IDS’s (em-IDS) also collect information mostly in the form of logs which are parsed to result in alerts. When evidentiary data from the local HbH is collected, it is transmitted towards the CBB via neighbouring HbH systems through the action of the P2P-da system (described in the next section). While such data is in transit through a neighbouring HbH system, and travelling onward to the CBB, it is always held in an encrypted form and only within temporary storage.

**4.2 The Peer-to-Peer Distribution Architecture (P2P-da)**

The essence of the P2P-da is to provide reliability, scalability and rapid throughput of transmitted data even in the face of high rates of “churn”, that is, large numbers of nodes joining and leaving the network. In order to achieve this, a cocktail of P2P

protocols are put together in order to extract the best qualities of these and also allow for each to compensate for the other's shortcomings.. The particular P2P protocols that are put together in order to build the P2P-da are: Gradient (Hierarchical) overlay protocols [15] Epidemic (Gossiping) protocols [16], and the Bit-Torrent protocol[17].

There are 3 main functionalities of the P2P-da:

- I. Maintenance of the P2P Overlay
- II. Dissemination and aggregation of Malicious Behaviour Information and alerts.
- III. Incident response data collection

These functionalities generally correspond to the P2P protocols that make up the essence of the P2P-da. The function of the maintenance of the P2P overlay is facilitated mainly through gradient (hierarchical) overlays assisted through epidemic (gossip-based) overlays. The dissemination and aggregation of malicious behavior information is mainly facilitated by epidemic (gossip-based) overlays. Incident response data collection is mainly facilitated through an adaptation of the Bit-Torrent protocol. The details behind these individual functionalities are dealt with in the following sections.

#### 4.2.1 Maintenance of the P2P Overlay

This is the need for the overall P2P network to maintain connectivity among neighbouring nodes as well as the larger HbH node population. Further to this, the aim here is to link HbH nodes in such a way that they are most beneficial to each other as well as to the overall communication of security events and evidence transmission aims.

In order to do this a hierarchy is to be created among the peer nodes such that those less endowed with resources are lower in the hierarchy and those that are better endowed are higher in the hierarchy. The aim of this is to ensure that nodes that lack resources generally communicate security event information, or transmit potentially large evidence files towards more reliable peers. It is assumed that nodes with more resources are more likely to be better equipped to deal with larger amounts of information and are also more likely to be online and available to be communicated with.

A gradient overlay network is suited to ensure this form of a network structure. It is built in such a way that a utility metric is used to determine which nodes are most suitable to connect to, and which nodes to avoid. This utility metric is determined from a combination of factors including the amount of resources available on a node, the current state of use of the node and the amount of time that it has been online. These utility metrics are shared through random node interactions, typical of "gossip-based" (epidemic) P2P protocols in order for nodes to get to know of other nodes that might be better to link to.

As gossip-based P2P protocols are known to eventually converge to a generally stable state, a hierarchy of the HbH systems is thus formed with the less endowed elements on the outer edges and the more capable elements closer towards the centre of the LEIA system (that is, the CBB).

#### 4.2.2 Dissemination and Aggregation of Malicious Behaviour Information & Alerts

This capability is necessary in order to facilitate the collaborative mechanisms needed to ensure that security event information is shared and eventually potentially useful evidence information is captured efficiently and transmitted securely. The sharing of security event information known by individual HbH peers is duly shared out to others in order for the overall system to have a more informed security posture as well as to be forewarned of imminent malicious events. This includes the distribution of malicious activity signatures as well as the detection of malicious activity originating from certain hosts. When such messages are received, only a set of the most common and recently active malicious activity signatures are maintained at the HbH. These kind of messages are termed as "Management messages" and can be shared out to any peers that a particular HbH has address information about and that has connectivity.

The other major type of messages that are involved in this functionality are termed as "Security Incident Control messages". These messages facilitate the reaction to the detection of a malicious event. This mainly includes the communication of procedures to initiate the evidence capture process on certain components of certain nodes as well as initiating initial pre-processing such as determining IP addresses of recently connected devices in order to extend the evidence capture process to other suspected devices.

There may be other forms of messages that might need to traverse the P2P-da, however, the 2 categories mentioned thus far are the major types.

#### 4.2.3 Incident response data collection

This functionality is triggered by the detection of malicious events via the collective knowledge gained through collaborating HbH systems, the em-IDS and guest OS security mechanisms. For more volatile data such as network traffic and live memory, a fixed time period is chosen for which to perform the capture process (or a fixed number of snapshots of the data over a short period of time particularly for live memory) after which a decision is to be made whether subsequent captures need to be made, or whether what has been collected so far suffices. Correspondence with the Cloud-Based Backend-Differencing Engine (CBB-DE) filters out known system files through facilitating the hash

comparisons. Primary analysis for IP addresses and hostnames on the data collected may result in triggering of other HbH systems to capture data also.

The actual data collection procedure involves 3 stages:

*a) Data Partitioning*

Different data formats (memory dumps, logs, files, packet captures, disk images) are compressed and stored temporarily on the HbH system in a modified AFF4 data structure that also contains simple RDF metadata describing the evidence. This data structure is termed as the Incident Data Archive (IDA). Each IDA data structure is partitioned in equal size pieces that will be referred to as shards. The shard is a signed and encrypted partition of the IDA analogous to the idea of a “piece” in the BitTorrent Protocol. A metadata file termed as the “reflection” (which corresponds to the BitTorrent Metadata file) is also created and sent directly to the CBB. In this way the CBB acts as the “tracker” and “leaches” IDAs from participating HbH systems in the P2P-da, thus benefiting from the high throughput of the BitTorrent protocol

*b) Shard Distribution*

Multiple copies of each individual shard are distributed to more capable neighbours (supporters), facilitated by the gradient overlay. Each time a shard is passed on it increases its “heat level”. After a certain “heat” threshold (the “melting point”) a HbH system is obliged to directly upload to the CBB (more specifically the HbH Master Peers of the CBB), else an election procedure is initiated to determine which previously supporting HbH should be delegated the uploading task. In order to avoid an individual node being the only “proxy” and thus a potential single point of failure, individual HbH systems are only allowed to partake in uploading a certain number of IDA shards governed by the “dependency value”. This improves the overall reliability of the larger system through reducing the possibility of having a single point of failure in the transmission process.

*c) Rapid fragment reconstruction*

For a particular shard, downloads are initiated from all their respective supporter locations. This is done for redundancy and bandwidth maximization purposes. Similar to the BitTorrent Protocol download, priority is given to the shards that are the least commonly available, that is, those that have the fewest recorded supporters.

In order to reconstitute the IDA, individual hashes of shards are verified as they are received, against that in the reflection. Several supporters upload at the same time, thus if a shard is in error, that from another supporter is taken. Once successfully transferred, shards are deleted from supporting HbH systems.

### 4.3 The Cloud-based Backend (CBB)

The CBB system is a highly available, scalable, responsive, centralized back end storage service capable of storing large amounts of data in a homogeneous form. It is subdivided into 3 major components: The Storage System (SS), the Differencing Engine (DE) and the HbH Master Peers.

The Storage System (SS) is built upon the Hadoop HDFS architecture [18] that provides not only the raw storage capabilities but also scalability, availability, reliability and responsiveness. The Differencing Engine (DE) filters out known files before having them stored on the CBB. This is provisioned through the MapReduce [19] capabilities supported by Hadoop. The DE also provides a query-response mechanism to the HBH systems with information on known benign data as part of the Master Known Data Hash-List (M-KDHL). The M-KDHL contains data about known files, memory processes, protocol flows, and log entries and thus enables their removal from IDAs being prepared. This reduces the size of IDAs before being stored on the Storage System (SS) of the CBB.

The HbH Master Peers are a particular set of well-endowed peers that are directly connected to the core CBB system (that is, the SS and DE) providing an interface to the rest of the LEIA system through the P2P-da. They do not have other core functionalities unrelated to their LEIA responsibilities and are essentially the backbone of the P2P-da and ultimately the provider of connectivity of the LEIA system outwards to the other HBH systems. The HBH Master Peers also serve as the central point through which system software updates and malicious event detection heuristics are originated from and disseminated outwards to the HBH systems in the wild.

### 4.4 The Law Enforcement Controller System

The Law Enforcement Controller is the main interface that law enforcement personnel interact with in order to perform their directed analysis for a particular digital investigation case. Through it, a Law Enforcement Agent can initiate specific queries to the data sets stored on the CBB, thus retrieving detailed, structured information as well as new knowledge inferred through correlation of data originating from different sources that may help in solving a case. The aim of this is to automate otherwise manual tasks of correlating data from different heterogeneous sources in order to pose valid assertions based on the data that could assist a forensic analyst in performing their duties of making sense of digital artifacts. This functionality is described in more detail by Dosis in [20].

Additionally, from the new found knowledge, patterns of malicious activities are to be learnt and stored. These Malicious Activity Patterns are to be used as feedback to the HbH systems in order to improve the detection capabilities of the inbuilt IDS's and thereby also improve the accuracy of collection of data of potential forensic evidentiary use.

### 5. Proof of Concept Evaluation and Results

In order to demonstrate the need for better throughput networks such as those exhibited in P2P overlays, an experiment was set up to simulate the conditions of the LEIA, however without the P2P-da component. This means that, the experiment was performed with the transmission of potentially evidentiary information from a HbH system to the CBB over a traditional client-server paradigm. The experiment itself focused on the remote extraction, compression and transmission over an encrypted channel of disk data evidence from small scale devices over the Internet and the subsequent reconstruction and storage of this data on a Hadoop cluster.

It should be mentioned that for the sake of simplicity of the experiment, the actual hypervisor of the HbH system was not built, however closely similar conditions – particularly in terms of the LEIA prototype application having privileged access – were met.

In order to test and measure the performance of the proof of concept application working over the client-server paradigm, four different small scale devices were used. The table below outlines the specifications of the devices being captured.

**Table 1: Small scale device specifications**

Device	Platform	RAM	Disk
Chumby Classic	Busybox v1.6.1	64MB	64MB
HTC Incredible S	Android OS v2.3.3 (Gingerbread)	768MB	1.1GB
HTC MyTouch 4G Slide	CyanogenMod 10.2 Alpha	768MB	4GB
Samsung Galaxy Tab 2 (WiFi Only)	Android OS, v4.0.3 (Ice Cream Sandwich)	1GB	8GB

In order to perform the testing and the performance evaluation, partitions of the various devices were filled with known sizes of random data and subsequently captured over the network. The capture process was repeated 10 times for each individual partition size in order to get the average file transfer times that each size took. The repetition process and the use of the averaging was done in order to avoid random processes that could have

affected network transmission times. Such random processes could include network traffic from other users of the networks being used, phone calls coming in and interfering with the I/O processes of the devices, or applications being updated on the devices, among others.

The tables below show the partition sizes used and the average times (in milliseconds) taken to perform the transfer:

**Table 2: Results from Test Cases on "HTC Incredible S"**

Partition Amount used	# of Test Runs	Avg. File Transfer time (ms)
16MB	10	13664
133MB	10	84600.8
250MB	10	392323.9
507MB	10	553933.1
1000MB	10	978571.8
1500MB	10	1360375
2000MB	10	2932376.8
3000MB	10	3877676.8
4000MB	10	4814006.6

**Table 3: Results from Test Cases on "HTC MyTouch 4G Slide"**

Partition Amount Used	# of Test Runs	Avg. File Transfer time (ms)
21.4MB	10	8583
87.0MB	10	31467
255MB	10	230709
500MB	10	338180
1000MB	10	1174482
1550MB	10	1323845.90
2000MB	10	1673928
3000MB	10	2052952.40
4000MB	10	3015056.60

**Table 4: Results from Test Cases on "Samsung Galaxy Tab 2"**

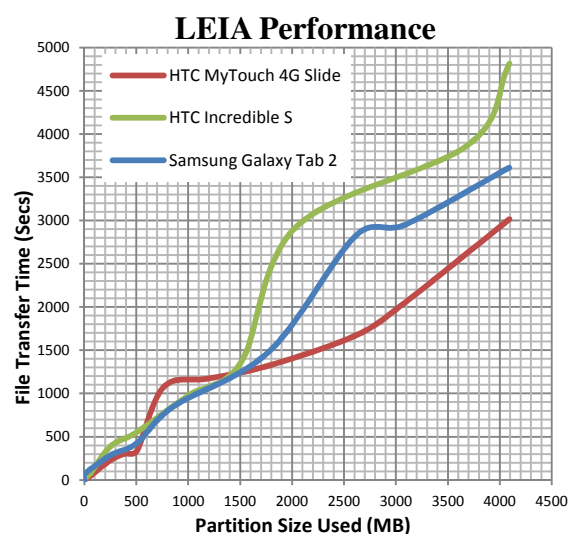
Partition Amount Used	# of Test Runs	Avg. File Transfer time (ms)
4MB	10	1235
11MB	10	67608
250MB	10	286947
500MB	10	426783
1000MB	10	960952
1500MB	10	1488236
2000MB	10	2829355
3000MB	10	2951551
4000MB	10	3707556

The data above from three of the four different specimen devices is plotted on a graph in order to visualize the general trend of the file transfer time against the partition size for the client server network



paradigm of remote evidence acquisition as explicated from the experiment performed.

The graph that follows depicts the graph that was attained:



**Figure 3: Performance of the LEIA Proof of Concept**

As has been noted only 3 out of the 4 specimen devices where graphed. This is because the fourth device (Chumby Classic) was not able to produce much valuable data. This is because it was equipped with a rather small sized disk (64MB) that did not permit much scaling upwards to the sizes that were tested on other devices. It should however be noted that it was successfully dumped several times and was the main device used while developing the application as its small disk size enabled quick disk captures that were needed while prototyping the initial system and testing for bugs.

From the graphs the curves seem to start off with a linear relationship which soon seems to turn into more of an exponential relationship. The "HTC MyTouch 4G Slide" clearly portrays this characteristic, with the rest of the devices also exhibiting this however not as vividly. Overall there seems to be a more exponential relationship between the Partition Size and the File Transfer Time with respect to the larger sizes of partitions. One could posit that as the partition sizes increase, even to sizes substantially larger than those in the graph, the relationship will become ever more exponential. This means that the times taken to acquire such partition sizes would be increase in exponential magnitude and thus shows that the client-server paradigm is likely not suitable enough for the task of performing remote evidence acquisition, especially in the type of environment that the LEIA system is aimed at. This is a fairly clear motivation for the need of a more efficient network transfer paradigm for this type of activity. Thus the suggested use of P2P networks between the evidence capture location and the

eventual storage location seems to be suitable as they generally provide higher network throughput.

## 6. Conclusion

In this study we outlined the numerous problems that blight the digital investigation process, and law enforcement agencies at large, rendering them slow and ultimately ineffective. We proposed a comprehensive architecture of a system – that makes use of hypervisors, P2P networks, the RDF framework and cloud storage – that could essentially revolutionize the digital investigation process through automation. Finally, through a small proof of concept we demonstrate part of this system, as well as make a clear motivation for need of P2P networks in order to improve the speed of remote evidence capture.

## 7. Future Work

Though this architecture is promising, several parameters within the P2P communication protocols need further optimization. A PKI infrastructure can be infused in the system in order to improve the security of the communication and storage facilities. The concept of privacy also needs to be addressed within the scope of this solution. Finally, an experiment with a wider scope would be greatly desired in order to better drive this architecture towards becoming a reality.

## 8. References

- [1] M. Cohen, S. Garfinkel, and B. Schatz, "Extending the advanced forensic format to accommodate multiple data sources, logical evidence, arbitrary information and forensic workflow," *Digit. Investig.*, vol. 6, pp. S57–S68, Sep. 2009.
- [2] B. C. Frank Adelstein, J. K. Eoghan Casey, Simson L. Garfinkel, Chet Hosmer, and P. T. Jim Lyle, Marcus Rogers, "Standardizing Digital Evidence Storage," *Commun. ACM*, vol. 49, no. 2, pp. 67–68, 2006.
- [3] S. Raghavan, A. Clark, and G. Mohay, "FIA: an open forensic integration architecture for composing digital evidence," *Forensics Telecommun. Inf. Multimed. Lect. Notes Inst. Comput. Sci. Soc. Informatics Telecommun. Eng.*, vol. 8, pp. 83–94, 2009.
- [4] A. Case, A. Cristina, L. Marziale, G. G. Richard, and V. Roussev, "FACE: Automated digital evidence discovery and correlation," *Digit. Investig.*, vol. 5, pp. S65–S75, Sep. 2008.
- [5] M. Davis, G. Manes, and S. Sheno, "A network-based architecture for storing digital evidence," *Adv. Digit. Forensics IFIP Int. Conf. Digit. Forensics*, vol. 194, pp. 33–42, 2005.
- [6] S. Zonouz, K. Joshi, and W. Sanders, "Floguard: cost-aware systemwide intrusion defense via online

forensics and on-demand IDS deployment,” in *Computer Safety, Reliability, and ...*, 2011, pp. 338–354.

[7] C. Shields, O. Frieder, and M. Maloof, “A system for the proactive, continuous, and efficient collection of digital forensic evidence,” *Digit. Investig.*, vol. 8, pp. S3–S13, Aug. 2011.

[8] J. Yu, Y. V. Ramana Reddy, S. Selliah, S. Reddy, V. Bharadwaj, and S. Kankanahalli, “TRINETR: An architecture for collaborative intrusion detection and knowledge-based alert evaluation,” *Adv. Eng. Informatics*, vol. 19, no. 2, pp. 93–101, Apr. 2005.

[9] M. I. Cohen, D. Bilby, and G. Caronni, “Distributed forensics and incident response in the enterprise,” *Digit. Investig.*, vol. 8, pp. S101–S110, Aug. 2011.

[10] S. Redding, “Using Peer-to-Peer Technology for Network Forensics,” *Adv. Digit. Forensics IFIP Int. Fed. Inf. Process.*, vol. 194, pp. 141–152, 2005.

[11] B. Schatz and A. Clark, “An open architecture for digital evidence integration,” in *AusCERT Asia Pacific Information Technology Security Conference*, 2006, no. May, pp. 15–29.

[12] D. Kahvedžić and T. Kechadi, “DIALOG: A framework for modeling, analysis and reuse of digital forensic knowledge,” *Digit. Investig.*, vol. 6, pp. S23–S33, Sep. 2009.

[13] S. L. Garfinkel, “AFF: A New Format for Storing Hard Drive Images,” *Assoc. Comput. Mach. Commun. ACM*, vol. 49, no. 2, pp. 85–87, 2006.

[14] W. Alink, R. a. F. Bhoedjang, P. a. Boncz, and A. P. de Vries, “XIRAF – XML-based indexing and querying for digital forensics,” *Digit. Investig.*, vol. 3, pp. 50–58, Sep. 2006.

[15] J. Sacha, J. Dowling, R. Cunningham, and R. Meier, “Discovery of stable peers in a self-organising peer-to-peer gradient topology,” in *International Conference on Distributed Applications and Interoperable Systems (DAIS)*, 2006, pp. 70–83.

[16] M. Jelasity, S. Voulgaris, R. Guerraoui, A.-M. Kermarrec, and M. van Steen, “Gossip-based peer sampling,” *ACM Trans. Comput. Syst.*, vol. 25, no. 3, pp. 1–36, 2007.

[17] B. Cohen, “Incentives build robustness in BitTorrent,” *Work. Econ. Peer-to-Peer Syst.*, 2003.

[18] K. Shvachko, H. Kuang, S. Radia, and R. Chansler, “The Hadoop Distributed File System,” *2010 IEEE 26th Symp. Mass Storage Syst. Technol.*, pp. 1–10, May 2010.

[19] J. Dean and S. Ghemawat, “MapReduce: simplified data processing on large clusters,” *Commun. ACM*, pp. 1–13, 2008.

[20] S. Dosis, I. Homem, and O. Popov, “Semantic Representation and Integration of Digital Evidence,” *Procedia Comput. Sci.*, vol. 22, pp. 1266–1275, 2013.