

FIA: An Open Forensic Integration Architecture for Composing Digital Evidence

Sriram Raghavan, Andrew Clark, and George Mohay

Information Security Institute, Queensland University of Technology,
Brisbane 4001, Australia
{s.raghavan,a.clark,g.mohay}@qut.edu.au

Abstract. The analysis and value of digital evidence in an investigation has been the domain of discourse in the digital forensic community for several years. While many works have considered different approaches to model digital evidence, a comprehensive understanding of the process of merging different evidence items recovered during a forensic analysis is still a distant dream. With the advent of modern technologies, pro-active measures are integral to keeping abreast of all forms of cyber crimes and attacks. This paper motivates the need to formalize the process of analyzing digital evidence from multiple sources simultaneously. In this paper, we present the forensic integration architecture (*FIA*) which provides a framework for abstracting the evidence source and storage format information from digital evidence and explores the concept of *integrating* evidence information from multiple sources. The *FIA* architecture identifies evidence information from multiple sources that enables an investigator to build theories to reconstruct the past. *FIA* is hierarchically composed of multiple layers and adopts a technology independent approach. *FIA* is also open and extensible making it simple to adapt to technological changes. We present a case study using a hypothetical car theft case to demonstrate the concepts and illustrate the value it brings into the field.

1 Introduction

In a digital investigation, investigators deal with acquiring digital data for examination. Digital records can vary in forms and types. Documents on a computer, telephone contact list, list of all phone calls made, trace of signal strengths from base station of a mobile phone, recorded voice and video files, email conversations, network traffic patterns and virus intrusions and detections are all examples of different types of digital records. Digital investigations must also contend with new challenges introduced by electronic equipment such as different devices, processor types, operating systems, storage formats and processing mechanisms that are used to store records in numerous formats. For the sake of this discussion, we restrict the classification of digital evidence to its source, data semantics and storage formats. We classify digital evidence based on its source, such as hard disks, volatile memory, or network traffic, its logical representation that defines its storage format and the type of information

that can be extracted from the source which determines the evidence semantics. No digital investigation is complete without an elaborate and systematic analysis along all three dimensions identified above.

A variety of new digital devices are being introduced with rapid advances in digital technology. Coping with such advances has become challenging owing to the use of proprietary data structures and protocols in most devices rendering them difficult for interpretation without relevant documentation, let alone, in a forensically sound manner. The large volumes of data collected in typical cases can be attributed to this variety and sifting through them can be enormously time consuming. Although digital forensics is in its early stages, there is a definite need to categorize digital evidence. This categorization is expected to limit the investigation space and minimize the effort spent on examining a variety of digital evidence.

From a forensic standpoint, there is too much entropy in the forensic examination process to capture all data and process it in one go. There is a need for capturing, understanding and analyzing information from disparate digital sources uniformly. Cohen [7] describes the PyFlag network forensic architecture, which is an open-source effort in providing a common framework for integrating forensic analysis from diverse digital sources. While PyFlag does support multiple image types and formats, it can only mount and examine one image at a time. PyFlag, thus, sorely lacks an architecture such as the one described in this paper to make the analysis more cohesive. As a first step to providing a common forensic analysis framework, this paper presents the architecture for integrating evidence information from different sources irrespective of the logical type of its contents.

Turner [19] states that as devices become more specialized, forensic examiners will require acquaintance with as many different processing tools to interpret the data they contain. This is attributed to the fact that forensic tools can only process digital devices as independent monolithic entities. The problem that this paper addresses is the multifarious interpretation and analysis of such evidentiary data in a uniform manner independent of origination source and storage formats. A preliminary validation of the concepts has been carried out on a hypothetical case involving a single disk image. The development of a prototype is planned as the next logical step to carry out a more comprehensive examination. This paper presents a conceptualization to how evidence integration can be achieved using content information from diverse evidence sources.

To illustrate the significance of evidence integration, consider a hypothetical case where investigators seize a personal computer and a mobile phone from a suspect. In the context of the investigation, it is essential to analyze the data contained in these sources uniformly, irrespective of semantics and storage formats. It is imperative that such a forensic framework be developed to support data interpretation from multiple sources.

Assume that on initial examination, investigators recover a set of suspicious documents which leads to the extraction of email messages exchanged between the suspect and suspect's contacts. Irrespective of the location and type of storage (either on mail servers or on a personal hard drive as user client profile), the data derived reinforces support to existing evidence and hence must be added to the framework under the same case.

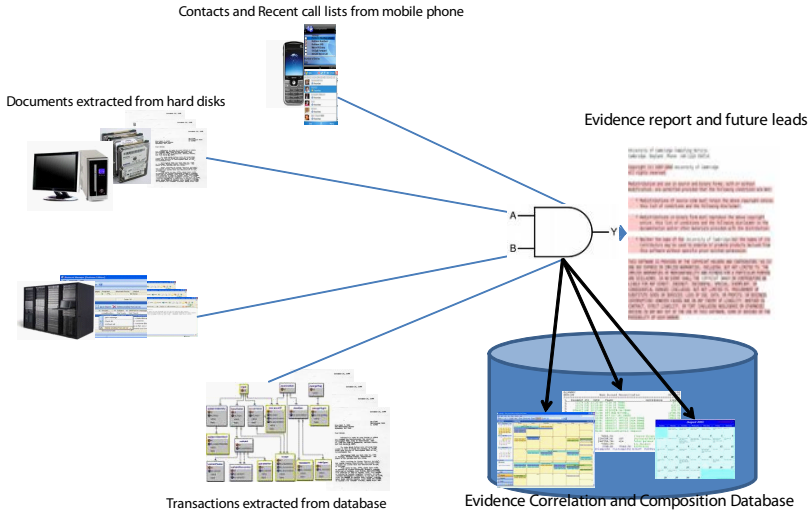


Fig. 1. An evidence composition example

The examination of email messages is expected to reveal some contact information and certain dates and times that might then be correlated with the current case to develop a *social calendar* of events and timelines. In addition, if the framework sources criminal records from a pre-existing repository that is indexed, then the correlation of extracted evidence with the repository can potentially reveal the underlying theme for the case, names and details of individuals involved, dates and times of activities reported or discussed relevant to the case. Such an extensive examination framework is illustrated in Figure 1. The framework aims to correlate (and hence compose) all reported information with the extracted evidence in an attempt to reconstruct the past. The rest of the paper is organized as follows. In Section 2, we review recent work in digital forensics and motivate the need for a common framework. In Section 3, we introduce the *forensic integration architecture*. In Section 4, we present a sample case study using a hypothetical case to demonstrate how FIA will operate. We conclude in Section 5 with a brief summary of the work reported and provide scope for future work.

2 State of the Art in Digital Investigations

Turner introduced the *digital evidence bags* (DEB) model [19] aimed at simplifying human interpretation. However, it is not intended to provide a methodology for combining different sources of evidence acquired at the various stages of an investigation process. Further, the model provides no scope for collecting and collating evidence from multiple digital sources which forms the crux of our work.

Schatz and Clark proposed a representation model to integrate metadata with evidence information in the sealed digital evidence bags (SDEB) [16]. The SDEB assumes the pre-existence of a forensic domain ontology model to support the representation of digital evidence. Such a model is yet to be developed.

Case et al. [6] introduce the FACE evidence correlation engine that parses data from different sources and correlates them. However, FACE assumes availability of all evidence sources at the start of analysis and the presence of known correlation in data. Besides, the engine does not integrate semantic information from evidence sources or provide for developing and validating assertions based on evidence analysis.

Alink et al. describe the XIRAF architecture [1] for indexing and retrieving stored digital evidence. The architecture indexes raw disk images and stores the content as annotated XML. However, XIRAF lays emphasis on feature extraction (indexing) and retrieval of digital evidence rather than on integrating evidence information that enables comprehensive forensic examination.

There are several other works in the literature that have reiterated the need for a common forensic analysis framework. Garfinkel highlights the problems associated with forensic analysis of raw computer disk images [9] and calls for the need to maintain an open and extendable standard for forensic analysis. Garfinkel introduced the Advanced Forensic Format (AFF) which is a two-layered forensic file system providing abstraction and extended functionality. However, the AFF is tailored to suit hard disk images and doesn't provide mechanisms to integrate evidence from multiple sources.

The Common Digital Evidence Storage Format Working Group has re-iterated the drawbacks with current forensic analysis tools [8] in terms of not being able to cope with multiple proprietary image formats. The authors emphasize the need for introducing a common digital evidence storage format that is common to a variety of evidence sources.

Beebe and Clark [2] argue the need for an objective based framework for digital forensics owing to the uniqueness of every forensic investigation. They divide the investigation process into 7 stages and propose a 2-tier hierarchical objectives framework. However, the focus of this framework is to maintain evidence integrity at all stages of an investigation which merely complements our focus in integrating evidence information and enabling further investigative leads.

Hosmer calls for the need to standardize the concept of digital evidence [11] to provide a common platform for investigators to perform forensic analysis. Since digital evidences can be altered, copied or erased, he proposed the 4-point principles of authentication, integrity, access control and non-repudiation for handling digital evidence.

Besides these efforts, several efforts in advancing the state of the art in techniques for data acquisition from electronic devices [5] have been reported. Some recent works have addressed challenges in the effective acquisition of volatile memory [14, 15, 17] and specifically in Windows based memory analysis in a computer [13, 18], while Buchholz and Spafford have studied the role of file system metadata in digital forensics [4]. Since digital forensics has predominantly been reactionary, some research contributions have been reported in formal methods for event reconstruction [10] and building theoretical foundations [12] to digital forensics. Turner has applied the DEB model to selective imaging of hard disk drives [20] and Beebe and Clark [3] introduce a text string search engine in for thematic searching in digital evidence.

The models and techniques described above have independently viewed the challenges in evidence analysis but are only stepping stones to integrate collected evidence from different sources. We require a framework that enables the development of new tools for interpretation of diverse data. Our work derives motivation from work reported in [8] and presents the FIA architecture as a means for abstracting

technology dependence of evidentiary data and integrating and composing information from different sources.

3 FIA for Composing Digital Evidence

We introduce a new architecture called the forensic integration architecture (FIA) that consists of 4 layers. The architecture is illustrated in Figure 2. The layers that constitute the FIA are:

1. *evidence storage and access layer;*
2. *representation and interpretation layer;*
3. *meta-information layer;* and
4. *evidence composition and visualization layer.*

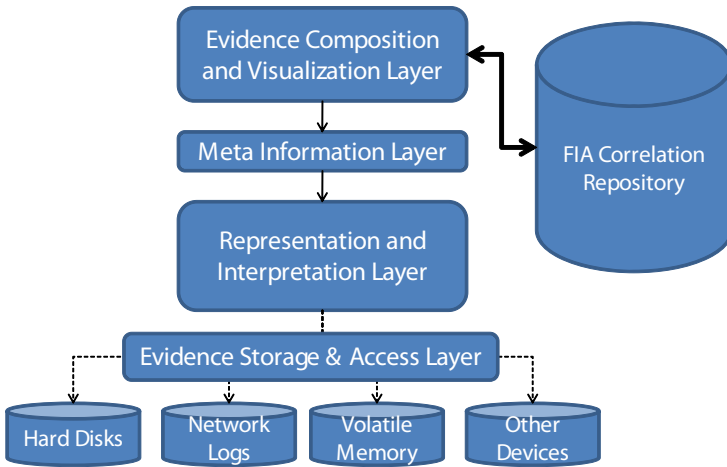


Fig. 2. Illustration of FIA evidence composition architecture

The FIA architecture is consistent with forensic principles. Based on a preliminary version of Turners DEB model, the FIA architecture pads the layers with added functionality that enhances its features and promises a natural transition towards automation. The layers are designed to allow scope for future extensions and selective modifications during analysis.

3.1 Evidence Storage and Access Layer

The *evidence storage and access layer* provides a binary abstraction to all data seized during an investigation. Acquisition of digital evidence is outside the scope of this work and the layer assumes that the evidence sources are forensically imaged copies stored on persistent media. All media must comply with *read only* semantics to maintain integrity of the data at all stages of an investigation. The layer supports registration interfaces for the acquired sources and their interpreters to be registered with FIA.

Once registered, the layer guarantees forensically secure access to the registered media. The layer also appends case specific metadata information prior to commencing analysis.

3.2 Representation and Interpretation Layer

The types of data that the *representation and interpretation layer* will be capable of supporting hard disk images and memory dumps from various operating systems (Windows OS, Linux, UNIX, Mac OS, etc.), network and system logs, and mobile devices with third party file systems (Nokia mobile with Symbian file system, iPod/iPhone with HFS+, etc.). The layer exploits interpreter semantics to extract logical blocks of data from the evidence sources for further analysis. For a file system image, this operation might correspond to extracting directories and files, for a memory dump it might correspond to extracting process control blocks (PCB) from the various processes resident in memory at the time of imaging and for network or system logs, it might correspond to extracting records of entries and their attributes from the log files. The extracted blocks are passed to the layer above. Figure 3 is indicative of some types of evidence semantics that different evidence sources require. The functionalities of this layer can be mapped to the file system support provided by most forensic tool suites which interpret the clusters and sectors of a disk system (e.g., FTK, Encase, PyFlag, etc.). However, this layer has the additional capability of interpreting the contents from other digital media and supporting memory and network forensics.

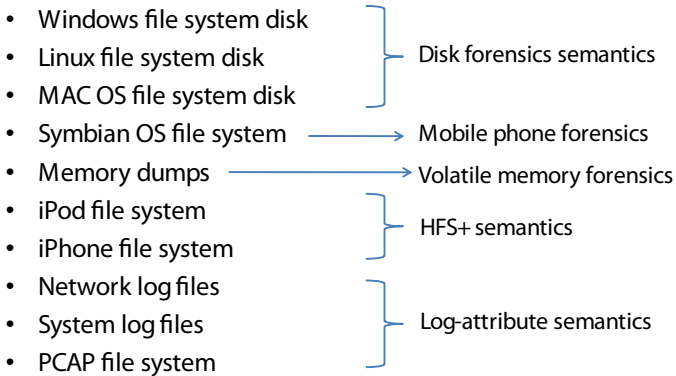


Fig. 3. Some types of evidentiary semantics used by different sources

3.3 Meta-information Layer

The *meta-information layer* supports application interfaces to extract metadata from objects present in the evidence sources. Every logical block of data extracted by the lower layer is represented as a file with properties that define its metadata. The meta-information layer uses a known file signature repository to filter metadata content from these blocks. For example, in hard disk images, files and metadata carry their usual

meaning. In memory dumps, PCB metadata might contain its allocated size in kilobytes, assembler type and process schedule information. In log file and packet capture sources, individual entry metadata might contain timestamps, process type, transaction source and destination and protocol information. Information such as the registered application executing a particular file is acquired while extracting metadata from file image. The functionality draws analogy to the file interpretation capabilities provided by existing forensic tool suites while supporting a larger variety of evidence sources.

3.4 Evidence Composition and Visualization Layer

The *evidence composition and visualization layer* is responsible for integrating information from various sources of evidence and composing the components into consistent and comprehensive evidentiary material for presentation to an investigator. This layer is composed of 3 sub-layers, *content indexing sub-layer*, *cross referencing sub-layer* and *knowledge representation and reasoning sub-layer*. The content indexing sub-layer is designed to index all syntactic content, such as keywords, locations, dates and timestamps, etc. in evidence sources and the cross referencing sub-layer cross references indexed data with entries in the FIA repository. This repository can be arbitrarily large and contain any external information that is deemed relevant to the case and be indexed in an identical manner. The knowledge representation and reasoning sub-layer is concerned with the truth value of information and logical inconsistencies in evidence data. The complete layer decomposition is illustrated in Figure 4. While the illustration shows the three sub-layers stacked one above the other, we acknowledge the presence of significant interplay between them and no particular order is pre-conceived in their representation.

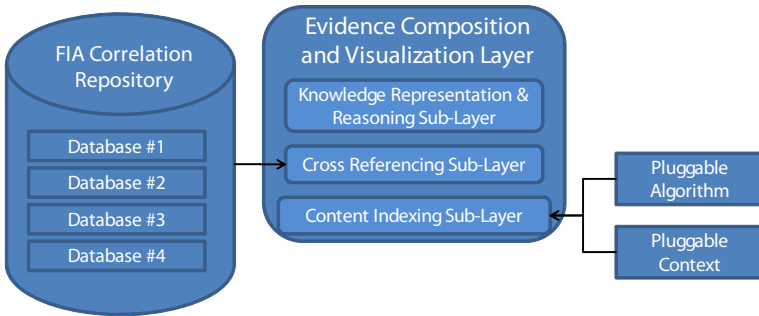


Fig. 4. Illustration of evidence composition and visualization layer

3.4.1 Content Indexing Sub-layer

The *content indexing sub-layer* supports mechanisms to index the logical blocks of data extracted by the lower layers. The indexing mechanism uses pluggable algorithms and pluggable contexts of keywords for indexing all syntactic content in data. The indexing process focuses on individualizing evidence, such as extracting and indexing names, locations, dates and events. We believe that such individualizations are crucial in any investigation and much sought after by investigators. This layer addresses the

challenges involved in syntactic indexing and correlation of digital evidence across different sources. Once the data is indexed for each source separately, the indices are integrated to create a comprehensive social calendar of names and events.

3.4.2 Cross Referencing Sub-layer

The *cross referencing sub-layer* is responsible for cross referencing indexed content with external databases to correlate evidence information with real world events. The sub-layer supports indexing a repository containing case relevant content databases using our evidence correlation model and cross referencing with data indexed from the evidence sources. Such a repository is built over a period of time from several investigations. For example, it might represent the collective knowledge of the investigation team learnt over that period. The list of content in the repository can include online dictionaries, automobile registration database, online map index database, calendar of dates and events and a database of social identification for individuals in a given area.

3.4.3 Knowledge Representation and Reasoning Sub-layer

The *knowledge representation and reasoning sub-layer* is concerned with the logical correctness of assertions and theories that are developed based on collected evidence. For example, consider a scenario where investigators discover simultaneous login attempts by a user from Brisbane and Perth into a corporate mail server. The information renders itself to the development of two independent assertions “*The user was in Brisbane at time X*” and “*The user was in Perth at time X*”. Clearly, information regarding user login attempts themselves, either from Brisbane or Perth, cannot be held against the user. However, correlating semantic information regarding the simultaneous attempts from two different cities provides a suspicious flavor to the actions and warrants further investigation. This sub-layer enables the development of such assertions and validating their truth value based on correlating information from multiple sources which is an integral part of any investigation. Any evidence to the contrary is flagged and presented to the investigator in an appropriate form which constitutes visualization.

4 Case Study – Car Theft Investigation

To demonstrate the concepts introduced in FIA, we present a case study using a hypothetical case concerning a car theft. The case was developed by Malcolm Corney at Queensland University of Technology as an assignment in a Computer Forensics course. The analysis of the case was carried out using existing forensic tools while following the FIA methodology. The true novelty of this architecture is described in Section 4.3. While the case contains only one disk image, we believe that the case involves sufficient diversity to demonstrate the utility of digital evidence integration. The actual value of this architecture, however, is perceived only when multiple such images are analyzed simultaneously.

4.1 About the Case

The case consists of a disk image containing multiple file system partitions. The case revolves around a chain of email messages recovered from Google Mail using the Thunderbird Client. The image contains several pictures of Australian wildlife with steganographic content containing pictures of car models. These car pictures represent cars recently reported stolen and currently under investigation. Each picture is password protected and the passwords are contained in an encrypted mail attachment. In addition, the disk slack space contains suspect's personal mail account details and a car model sequence that is traced back to the sequence of car thefts reported.

4.2 Extracting the Data

The disk was imaged using dd UNIX imager and the copy was hashed to preserve its integrity. The imaged disk was then registered under a new case with source and semantics information. This action reflects the functionality of the *evidence storage and access layer*. The image was then analyzed using FTK to detect and extract files, mail drafts and inbox messages. These actions reflect the extraction of logical blocks of data from the *representation and interpretation layer*. The same tool was also used to extract file properties or metadata information from the files and wildlife picture files. PRTK was used to crack the password of the encrypted file which contained passwords to the steganographic pictures which in turn provided more metadata. In FIA, this operation is performed at the *meta-information layer*. These operations were repeated with Encase and Sleuthkit to corroborate the results.

4.3 Evidence Composition

Once all the relevant data is extracted, the *evidence composition and visualization layer* takes over and indexes content in the extracted logical blocks. In our case study, the chain of email messages was used as the main source to generate evidence composition as illustrated in Figure 5. The contents were then cross referenced with multiple databases held in FIA repository to determine potential connections. Using directed keyword and metadata searches, an illustration of how FIA might piece the different sources of evidence together is illustrated in Figure 6. The dates of creation of picture files produced a pattern that traced back to the dates in the email chain indexed previously. Metadata analysis of the pictures further revealed the use of a particular camera that was recovered from the suspect's premises. The contact list from Thunderbird client revealed two persons with criminal record history, when their names presented hits in a simulated police database. The car registration numbers were cross referenced with simulated databases containing *automobile registration details* to determine the owners of the cars and *police complaint details* to verify if a stolen complaint has been registered since the theft and whether theft details fitted the description. Further, an address recovered from the email content was searched for a registration log (again, added to the FIA repository) to determine if the owner of the premises had collaborated with the suspect to store the stolen cars until they were shipped offshore.

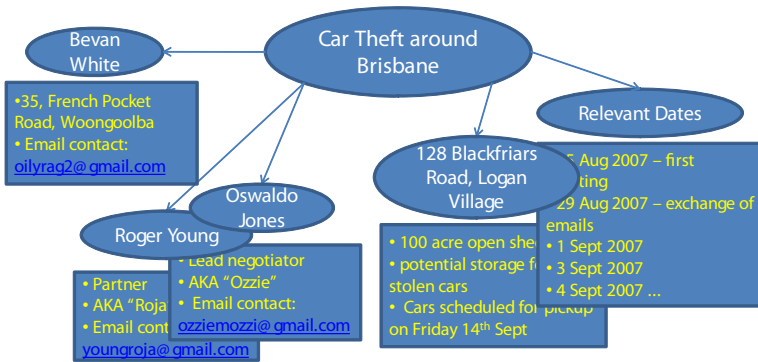


Fig. 5. Illustrating evidence composition for car theft case

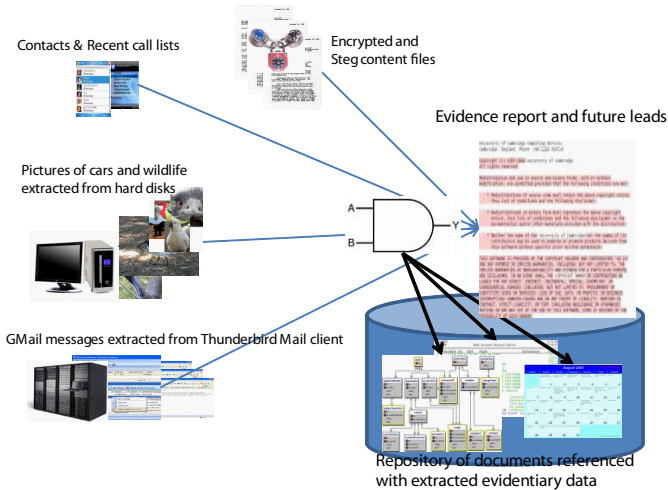


Fig. 6. Correlating data in car theft investigation using FIA

We have thus demonstrated the operation of the FIA architecture on a car theft case to perform evidence integration. FIA views the extracted data as conceptual sources, both at a syntactic and semantic level, and correlates information present in these sources in an attempt to reconstruct the past. In the process, FIA also aids the generation of further investigative leads that enable investigators to build a strong foundation for the case based on scientific evidence and facts.

5 Conclusions and Future Work

In this paper, we presented the FIA architecture for integrating digital evidence from multiple evidence sources in a technology independent manner and composing evidence information. The architecture supports all requirements of forensic security. In

addition, FIA also supports indexing content identified as conceptual sources and cross references them with a repository of internal and external databases relevant to a case. To the best of our knowledge, FIA is the only known work that attempts to integrate different sources of evidence and compose comprehensive evidence. The architecture is hierarchical and completely modular and extensible to keep pace with challenges that frequently crop up in this field. The model has been demonstrated with a hypothetical case study involving car theft.

Future work will focus on the design and comprehensive validation of a prototype with real evidence data. Research is currently underway into developing data representation and effective indexing algorithms for content in FIA repository for evidence property identification in different evidence sources.

References

1. Alink, W., Bhoedjang, R.A.F., Boncz, P.A., de Vries, A.P.: XIRAF - XML-based indexing and querying for digital forensics. *Digital Investigation*. In: The Proceedings of the 6th Annual Digital Forensic Research Workshop (DFRWS 2006), vol. 3(suppl. 1), pp. 50–58 (2006)
2. Beebe, N.L., Clark, J.G.: A hierarchical, objectives-based framework for the digital investigations process. *Digital Investigation* 2(2), 147–167 (2005)
3. Beebe, N.L., Clark, J.G.: Digital forensic text string searching: Improving information retrieval effectiveness by thematically clustering search results. *Digital Investigation* 4(suppl. 1), 49–54 (2007)
4. Buchholz, F., Spafford, E.: On the role of file system metadata in digital forensics. *Digital Investigation* 1(4), 297–308 (2004)
5. Carrier, B.D., Grand, J.: A hardware-based memory acquisition procedure for digital investigations. *Digital Investigation* 1(1), 50–60 (2004)
6. Case, A., Cristina, A., Marziale, L., Richard, G.G., Roussev, V.: FACE: Automated digital evidence discovery and correlation, *Digital Investigation*. In: The Proceedings of the Eighth Annual DFRWS Conference, 5th edn., pp. S65–S75 (September 2008)
7. Cohen, M.I.: PyFlag - An advanced network forensic framework, *Digital Investigation*. In: The Proceedings of the Eighth Annual DFRWS Conference, vol. 5(suppl. 1), pp. S112–S120 (September 2008)
8. Common Digital Evidence Storage Format Working Group. Standardizing digital evidence storage. *Communications of the ACM* 49(2), 67–68 (February 2006)
9. Garfinkel, S.: AFF: a new format for storing hard drive images. *Communications of the ACM* 49(2), 85–87 (2006)
10. Gladyshev, P., Patel, A.: Finite state machine approach to digital event reconstruction. *Digital Investigation* 1(2), 130–149 (2004)
11. Hosmer, C.: Digital evidence bag. *Communications of the ACM* 49(2), 69–70 (2006)
12. Mocas, S.: Building theoretical underpinnings for digital forensics research. *Digital Investigation* 1(1), 61–68 (2004)
13. Mee, V., Tryfonas, T., Sutherland, I.: The Windows Registry as a forensic artefact: Illustrating evidence collection for Internet usage. *Digital Investigation* 3(3), 166–173 (2006)
14. Nikkel, B.J.: Improving evidence acquisition from live network sources. *Digital Investigation* 3(2), 89–96 (2006)

15. Petroni, J., Nick, L., Walters, A., Fraser, T., Arbaugh, W.A.: FATKit: A framework for the extraction and analysis of digital forensic data from volatile system memory. *Digital Investigation* 3(4), 197–210 (2006)
16. Schatz, B., Clark, A.: An Open architecture for digital evidence integration. In: *Proceedings of the 2006 AUSCERT R&D Stream*, pp. 15–29 (2006)
17. Schatz, B.: BodySnatcher: Towards reliable volatile memory acquisition by software. *Digital Investigation* 4(suppl. 1), 126–134 (2007)
18. Schuster, A.: Searching for processes and threads in Microsoft Windows memory dumps. *Digital Investigation*. In: *The Proceedings of the 6th Annual Digital Forensic Research Workshop (DFRWS 2006)*, vol. 3(suppl. 1), pp. 10–16 (2006)
19. Turner, P.: Unification of digital evidence from disparate sources (Digital Evidence Bags). *Digital Investigation* 2(3), 223–228 (2005)
20. Turner, P.: Selective and intelligent imaging using digital evidence bags. *Digital Investigation*. In: *The Proceedings of the 6th Annual Digital Forensic Research Workshop (DFRWS 2006)*, vol. 3(suppl. 1), pp. 59–64 (2006)