

Cryptographic Privacy-Preserving
Enhancement Method for
Investigative Data Acquisition



Zbigniew Kwecka

April 2011

A THESIS SUBMITTED TO EDINBURGH NAPIER UNIVERSITY FOR THE
DEGREE OF DOCTOR OF PHILOSOPHY IN THE FACULTY OF
ENGINEERING, COMPUTING AND CREATIVE INDUSTRIES

Contents

CONTENTS	2
LIST OF TABLES	6
LIST OF FIGURES	7
ABSTRACT	9
DECLARATION	12
COPYRIGHT	13
ACKNOWLEDGMENTS	14
LIST OF ACRONYMS	15
CHAPTER 1 INTRODUCTION	17
1.1 INTRODUCTION	17
1.2 BACKGROUND AND CONTEXT	18
1.3 MOTIVATING SCENARIOS	19
1.4 RESEARCH QUESTION, AIM AND OBJECTIVES.....	22
1.5 CONTRIBUTION TO KNOWLEDGE	24
1.6 PUBLICATIONS	24
1.7 THESIS STRUCTURE	25
CHAPTER 2 BACKGROUND AND THEORY	27
2.1 INTRODUCTION	27
2.2 DIGITAL FORENSIC	27
2.3 INVESTIGATIONS USING THIRD PARTY DATA	31
2.3.1 <i>Data Protection Act</i>	31
2.3.2 <i>Regulation of Investigatory Powers Act</i>	33
2.3.3 <i>Data retention</i>	35
2.3.4 <i>Commission and diligence for recovery of documents</i>	37

2.4	PRIVACY, AND ITS WIDER CONTEXT	37
2.4.1	<i>Privacy</i>	38
2.4.2	<i>Measuring privacy</i>	39
2.4.3	<i>Privacy in surveillance systems</i>	40
2.5	CRYPTOGRAPHY	42
2.5.1	<i>Classification of cryptographic protocols</i>	43
2.5.2	<i>Authentication, Integrity and Non-repudiation – Public Key Infrastructure (PKI)</i>	45
2.5.3	<i>Operations on encrypted plaintext</i>	46
2.5.4	<i>Commutative Cryptography</i>	46
2.5.5	<i>Cryptanalysis</i>	49
2.6	CONCLUSION	50
CHAPTER 3 LITERATURE REVIEW		53
3.1	INTRODUCTION	53
3.2	PRIVACY-RESPECTING INVESTIGATIVE DATA ACQUISITION.....	54
3.2.1	<i>Privacy self-defence</i>	54
3.2.2	<i>Privacy Controls in an Investigative Scenario</i>	56
3.3	PRIVACY-PRESERVING PRIMITIVES	59
3.3.1	<i>Privacy-Enhancing Technologies</i>	59
3.3.2	<i>Multi-Party Computation</i>	62
3.3.3	<i>Sharing a secret</i>	63
3.3.4	<i>Retrieving records in a private manner</i>	66
3.3.5	<i>Private Value Comparison – Locating interesting records</i>	71
3.3.6	<i>Combined approaches to selective information retrieval</i>	74
3.3.7	<i>Security Considerations</i>	78
3.4	CONCLUSION	81
CHAPTER 4 IMPROVING THE ACQUISITION PROCESS		84
4.1	INTRODUCTION	84
4.2	METHODOLOGY	85
4.3	INITIAL REQUIREMENTS	86
4.4	OVERALL DESIGN	88
4.5	APPROACH 1: COMBINATION OF PET PRIMITIVES.....	91
4.6	APPROACH 2: COMBINED PET PRIMITIVES	95
4.7	EVALUATION	96
4.7.1	<i>Experiment Design and Implementation</i>	96
4.7.2	<i>Empirical Evaluation</i>	100
4.7.3	<i>Feedback from practitioner</i>	105
4.8	REQUIREMENTS REVIEW	107
4.9	CONCLUSIONS	107

CHAPTER 5	NOVEL DATA ACQUISITION PLATFORM.....	109
5.1	INTRODUCTION	109
5.2	METHODOLOGY	110
5.3	IDAP DESIGN	111
5.3.1	<i>Lowering Processing Time</i>	<i>111</i>
5.3.2	<i>Allow multiple selection criteria</i>	<i>118</i>
5.3.3	<i>Reassuring the Public</i>	<i>119</i>
5.4	IMPLEMENTATION.....	123
5.5	PROPOSED QUANTITATIVE EVALUATION.....	125
5.5.1	<i>Overall design of experimental environment</i>	<i>126</i>
5.5.2	<i>Experiments</i>	<i>128</i>
5.5.3	<i>Proposed Qualitative Evaluation</i>	<i>129</i>
5.6	CONCLUSION	130
CHAPTER 6	EVALUATION	132
6.1	INTRODUCTION	132
6.2	PRESENTATION OF PERFORMANCE IMPACT	133
6.2.1	<i>Evaluation of the complexity model</i>	<i>133</i>
6.2.2	<i>Varying the number of interesting records</i>	<i>137</i>
6.2.3	<i>Varying number of enquiries</i>	<i>140</i>
6.2.4	<i>Evaluation of IDAP with directory of identities</i>	<i>142</i>
6.2.5	<i>Use of dilution factor with different protocols</i>	<i>143</i>
6.2.6	<i>Controlling the balance between privacy and feasibility</i>	<i>144</i>
6.3	PRESENTATION OF QUALITATIVE EVALUATION.....	145
6.3.1	<i>Correctness and Security</i>	<i>145</i>
6.3.2	<i>Survey Results</i>	<i>148</i>
6.4	CONCLUSION	149
CHAPTER 7	CONCLUSIONS AND FUTURE WORK.....	151
7.1	INTRODUCTION	151
7.2	ACHIEVEMENT OF OBJECTIVES	152
7.3	MOTIVATING SCENARIOS WITH SOLUTIONS	156
7.4	CONTRIBUTION TO KNOWLEDGE	158
7.5	CRITICAL ANALYSIS.....	160
7.6	MAIN FINDINGS	163
7.7	FUTURE WORK	166
CHAPTER 8	REFERENCES	168
APPENDIX A	SIMPLIFIED OPERATION OF IDAP	178
APPENDIX B	EMPIRICAL EVALUATION RESULTS.....	180

APPENDIX C SURVEY..... 182

List of Tables

Table 4-1 Cryptographic operation performance measurements in nanoseconds (<i>ns</i>)	101
Table 4-2 Computational complexity of the OT-based approach.....	103
Table 4-3 Computational complexity of the PE-based approach	103
Table 6-1 Initial definition of IDAP's complexity.	135
Table 6-2 IDAP's complexity.	136
Table 6-3 IDAP's complexity for datasets with publically available dictionaries.	137

List of Figures

Figure 2-1 Analogy to the operation of the three-pass protocol.....	47
Figure 2-2 Three-pass protocol operation	48
Figure 3-1 Locking a secret under Khayat’s secret sharing scheme.....	64
Figure 3-2 Illustration of unlocking a secret under Khayat’s secret sharing scheme	65
Figure 3-3 Shamir’s simple secret sharing scheme	66
Figure 3-4 Basic Oblivious Transfer Protocol by Schneier [41]	69
Figure 3-5 Private Equality Test based on Commutative Cryptography.....	73
Figure 3-6 Operation of the Private Equi-join protocol	76
Figure 3-7 Efficient PIR based on Secure Coprocessor	76
Figure 3-8 Graphical representation of the Private Equi-join protocol.	77
Figure 4-1 Typical request for investigative data mapped into SQL.....	89
Figure 4-2 Request enabling privacy-preserving queries mapped into SQL.....	89
Figure 4-3 Retrieval Phase	95
Figure 4-4 Test Dataset	98
Figure 4-5 Total running time for both approaches including preparation time.....	104
Figure 4-6 Data Acquisition processing time excluding preparation time.	104
Figure 4-7 Performance of both approaches for varied number of records ($1 < m < n$).....	104
Figure 5-1 Process flow of the protocol incorporating the <i>dilution factor</i>	114
Figure 5-2 Lowering Processing Time Phase A – Preparation.....	115
Figure 5-3 Lowering Processing Time Phase B – Searching	116
Figure 5-4 Lowering Processing Time Phase C – Retrieval	118
Figure 5-5 mapped into SQL.....	118
Figure 5-6 Reassuring the public by introducing semi-trusted third parties.....	121
Figure 5-7 IDAP.....	125
Figure 5-8 Measuring processing time	126
Figure 5-9 Measuring bandwidth used during a protocol run	127
Figure 6-1 Complexity table reading vs. actual measurements.....	134

Figure 6-2 IDAP complexity in detail for varied m ; $k=1$; $o=1,000$ and $n=1,000,000$	136
Figure 6-3 Computational complexity of IDAP and PE for increasing m (logarithmic scale).....	138
Figure 6-4 Computational complexity of IDAP and PE for increasing m	139
Figure 6-5 IDAP's processing time for varying m and different values of n	139
Figure 6-6 IDAP's processing time for varying m and different values of o	139
Figure 6-7 Comparison of processing time for IDAP and PE.....	140
Figure 6-8 Detailed comparison for IDAP and PE with.....	141
Figure 6-9 IDAP's performance for different values of γ , as compared to PE.....	141
Figure 6-10 Performance gain for IDAP run on dataset with a directory.....	142
Figure 6-11 Performance gain for IDAP run on dataset with a directory.....	143
Figure 6-12 Comparison of IDAP to a modification of OT-based approach,	144
Figure 6-13 Processing time against the dilution factor o , for IDAP with directory of records.....	145
Figure 7-1 Hiding of request originators can improve privacy	165

Abstract

The current processes involved in the acquisition of investigative data from third parties, such as banks, Internet Service Providers (ISPs) and employers, by the public authorities can breach the rights of the individuals under investigation. This is mainly caused by the necessity to identify the records of interest, and thus the potential suspects, to the dataholders. Conversely, the public authorities often put pressure on legislators to provide a more direct access to the third party data, mainly in order to improve on turnaround times for enquiries and to limit the likelihood of compromising the investigations. This thesis presents a novel methodology for improving privacy and the performance of the investigative data acquisition process. The thesis shows that it is possible to adapt Symmetric Private Information Retrieval (SPIR) protocols for use in the acquisition process, and that it is possible to dynamically adjust the balance between the privacy and performance based on the notion of k -anonymity. In order to evaluate the findings an Investigative Data Acquisition Platform (IDAP) is formalised, as a cryptographic privacy-preserving enhancement to the current data acquisition process.

SPIR protocols are often computationally intensive, and therefore, they are generally unsuitable to retrieve records from large datasets, such as the ISP databases containing records of the network traffic data. This thesis shows that, despite the fact that many potential sources of investigative data exist, in most cases the data acquisition process can be treated as a single-database SPIR. Thanks to this observation, the notion of k -anonymity, developed for privacy-preserving statistical data-mining protocols, can be applied to the investigative scenarios, and used to narrow down the number of records that need to be processed by a SPIR protocol.

This novel approach makes the application of SPIR protocols in the retrieval of investigative data feasible.

The *dilution factor* is defined, by this thesis, as a parameter that expresses the range of records used to hide a single identity of a suspect. Interestingly, the value of this parameter does not need to be large in order to protect privacy, if the enquiries to a given dataholder are frequent. Therefore, IDAP is capable of retrieving an interesting record from a dataholder in a matter of seconds, while an ordinary SPIR protocol could take days to complete retrieval of a record from a large dataset.

This thesis introduces into the investigative scenario a semi-trusted third party, which is a watchdog organisation that could proxy the requests for investigative data from all public authorities. This party verifies the requests for data and hides the requesting party from the dataholder. This limits the dataholders ability to judge the nature of the enquiry. Moreover, the semi-trusted party would filter the SPIR responses from the dataholders, by securely discarding the records unrelated to enquiries. This would prevent the requesting party from using a large computational power to decrypt the diluting records in the future, and would allow the watchdog organisation to verify retrieved data in court, if such a need arises. Therefore, this thesis demonstrates a new use for the semi-trusted third parties in SPIR protocols. Traditionally used to improve on the complexity of SPIR protocols, such party can potentially improve the perception of the cryptographic trapdoor-based privacy-preserving information retrieval systems, by introducing policy-based controls.

The final contribution to knowledge of this thesis is definition of the process of privacy-preserving matching records from different datasets based on multiple selection criteria. This allows for the retrieval of records based on parameters other than the identifier of the interesting record. Thus, it is capable of adding a degree of fuzzy matching to the SPIR protocols that traditionally require a perfect match of the request to the records being retrieved. This allows for searching datasets based on circumstantial knowledge and suspect profiles, thus, extends the notion of SPIR to more complex scenarios.

The constructed IDAP is thus a platform for investigative data acquisition employing the Private Equi-join (PE) protocol – a commutative cryptography SPIR protocol.

The thesis shows that the use of commutative cryptography in enquiries where multiple records need to be matched and then retrieved (m -out-of- n enquiries) is beneficial to the computational performance. However, the above customisations can be applied to other SPIR protocols in order to make them suitable for the investigative data acquisition process. These customisations, together with the findings of the literature review and the analysis of the field presented in this thesis, contribute to knowledge and can improve privacy in the investigative enquiries.

Declaration

No portion of the work referred to in this thesis has been submitted in support of an application for another degree or qualification of this or any other university or other institution of learning.

Copyright

Copyright in text of this thesis rests with the Author. Copies (by any process) either in full, or of extracts, may be made only in accordance with instructions given by the Author and lodged in the Edinburgh Napier University. Details may be obtained from the Librarian. This page must form part of any such copies made. Further copies (by any process) of copies made in accordance with such instructions may not be made without the permission (in writing) of the Author. The ownership of any intellectual property rights which may be described in this thesis is vested in Edinburgh Napier University, subject to any prior agreement to the contrary, and may not be made available for use by third parties without the written permission of the University, which will prescribe the terms and conditions of any such agreement. Further information on the conditions under which disclosures and exploitation may take place is available from the Head of School of Computing.

Acknowledgments

First, and foremost, I would like to thank my family for never doubting in me. So to my wife Wiola, to my parents Maria and Roman, and to my brother Andrzej, I thank you for always being there for me. Also, thanks to my daughter Amelia for keeping me sane near the end of this lengthy process, as her smile can always take weight off my shoulders.

My supervision team has always provided support and suggested plenty of avenues to explore. Prof. William Buchanan, my supervisor, was always an inspiration and never-ending source of research ideas. To Duncan Spiers, the second supervisor, many thanks for discussing the legal aspects of this thesis with me and for helping me to validate my ideas. Dr Michael Smyth has been a great panel chair, as keen listener with an eye for detail he has kept me on the right track throughout the process.

Finally, huge 'thank you' to my friends and colleagues. With special thanks to the C44 research team: to Jamie, Kevin, Lionel and Tim for getting me interested in doing research and helping me with my first steps. Also, I would like to thank Stuart and Allan for reading things at the last minute and providing ideas for improvements.

I know I did not manage to mention you all in here, but I do remember your efforts in motivating me. So thanks to all that throughout this long process have asked me a countless number of times "So when are you getting your PhD?". You know what? I just wouldn't make it without you!

Thank You!

List of Acronyms

CSP	Content Service Provider
FIPS	Federal Information Processing Standard
IDAP	Investigative Data Acquisition Platform
MPC	Multi-Party Computation
OT	Oblivious Transfer
PE	Private Equi-Join
PET	Privacy Enhancing Technology
PIR	Private Information Retrieval
SCOP	Secure Coprocessor
SIPR	Scottish Institute for Policing Research
SPIR	Symmetric PIR
SPoC	Single Point of Contact
TC	Trusted Computing
TOR	The Onion Router
ZKP	Zero-Knowledge Proof

Those who would give up essential Liberty, to purchase a little temporary safety, deserve neither Liberty nor Safety.

-Benjamin Franklin, 11 Nov 1755

Chapter 1

Introduction

1.1 Introduction

This thesis provides a novel method for obtaining investigative data from third-parties. Current processes often breach the human rights of the individuals being investigated, and they may expose the investigation. Consequently, in this thesis, a Privacy Enhancing Technologies (PET) solution is proposed based on the Symmetric Private Information Retrieval (SPIR) primitive customised to perform the specific task of data acquisition. The proposed customisation is the main contribution to knowledge of this thesis, and it can be applied to most single database SPIR protocols.

This research commenced a few years after the tragic events of 9/11, which started the 21st Century's *War on Terror*. It was this event that caused the UK to alter the laws and procedures for obtaining intelligence data towards more intrusive solutions. Thus, despite the protests of civil liberty activists, the UK decision-makers are still

prepared to widen the warrantless access to data, in order to satisfy the needs of investigators in the ever-changing communication environment. The possibility of direct access by the public authorities to information held by Communication Service Providers (CSPs) is one of the fiercely discussed topics. Consequently, the aim of this thesis is to propose a privacy-respecting mechanism that could be used by investigators in order to access data held by third-parties.

This thesis, thus, identifies that a major issue in the way that investigations are conducted is the need for public authorities to disclose the identities of their suspects to the data-holders. Such a disclosure may negatively impact the suspects, and also jeopardise the investigations. In order to protect the interest of investigations the public authorities may be forced to wait until they have a number of similar cases, or to widen their enquiry in order to dilute suspect's identity. This, depending on which technique is used, can lead to delays, or be treated as *fishing-for-evidence*. Therefore, this thesis sets out to provide an efficient way to automate the investigative data acquisition process, if the rights of the data subjects and the secrecy of the investigations are to be protected. The solution should be able to perform tasks regulated by a number of, often contradicting, legislative requirements. It should also allow for the fine tuning of the correct balance between privacy, security and performance.

1.2 Background and Context

Since 11 September 2001 many Western governments have passed laws empowering public authorities with wider rights to gather operational data [1, 2]. For many years public opinion accepted the invasion of personal privacy rights as the sacrifice needed to fight terrorism [3]. However, slowly, public opinion is shifting back to a state where such measures are often considered unacceptable. This is shown by public opinion surveys, such as the one conducted in US by Washington Post [4], where 32% of respondents agreed that they would prefer the federal government to ensure that privacy rights are respected rather than to investigate possible terrorist threats. This was an 11% increase from the similar survey conducted in 2003. The trend continues, since in [5] 63% of respondents stated that they are worried about the government's surveillance into personal lives.

In the UK, the public authorities, including the Police, request investigative data from third-parties on regular basis [6] and the data protection legislation allows for such requests, even without warrants [7, 8]. Depending on the way these requests are performed, human and natural rights of the data-subject can be breached, and/or the investigation may be compromised [9]. A recent proposal by the UK government went further and recommended allowing the public authorities direct access to data held by Content Service Providers (CSPs), such as mobile telephony providers, and Internet Service Providers (ISPs) [2]. There are a few major motivating factors behind this proposal: increasing access speeds to records; allowing for covert enquiries by anti-terror and national security agencies; reducing collateral damage to potential suspects under investigation; and enabling the analysis of data to facilitate the profiling of terrorists activities. In response, concerns were raised that if the proposal was implemented, it would thwart the privacy of Internet users around the globe in order to increase the security of one nation. This thesis shows that most of the objectives set out in the proposal can still be achieved while maintaining a high level of privacy. It is shown that an investigative system can maintain the privacy of the data subjects and also preserve the confidentiality of investigations. However, both security and privacy must be built into the system at the design stage in order to achieve this [1].

This thesis, thus, gives an insight into the use of PETs in improving the current investigative data acquisition practices and defines the Investigative Data Acquisition Platform (IDAP). IDAP is a proposed novel approach to maintain secrecy; preserving the suspect's privacy and gaining the public's support for the PET technologies in digitalised investigative enquiries.

1.3 Motivating Scenarios

There are a number of possible scenarios that the data acquisition process needs to facilitate. Legislations permit different public authorities to request investigative data from third parties under a variety of circumstances. To highlight the possible breaches of privacy and/or human rights of the suspects, we consider the following scenarios:

Scenario 1: Request for ISP subscriber data

A forensic investigation, carried out by a law enforcement agency on a confiscated Personal Computer (PC), has identified 14 different Internet Protocol (IP) addresses linked to organised crime. The agency would like to identify the owners of these IP addresses, and their subscriber data including the postal address. However, it is key that the nature of the enquiry, and identities of sought after individuals, are not revealed to the ISP (directly, or indirectly) in order to protect the integrity of the investigation.

Scenario 2: Banking transaction details

A shopkeeper has notified a law enforcement agency about the purchase of an uncontrolled substance that, in the wrong hands, can be used to produce an exploding device. The credit card number used in the transaction is made available to the agency. This agency would like to find out the list of purchases made on this card for the previous month, as well as the name of the owner. Since banks are not obliged to provide such information to the public authorities, the nature of the enquiry will have to be communicated to the bank. However, the identity of the potential suspect should be kept secret from the bank, as not to affect this individual's relation with the bank. If the bank was aware of a given individual being a suspect in an investigation then, as an example, the individual could be placed on a list of high-risk borrowers. This may stop them from getting a loan, even though they have not been charged with a crime. Most importantly, this individual may be unable to find out why his application was refused, since the disclosure of matters affecting national security and crime prevention are exempt from many provisions of the UK's Data Protection Act 1998 (DPA) (sections 28 and 29)[10].

It is clear that in these scenarios, investigations can be compromised by revealing the identity of the suspects to the data-holders. However, the second scenario demonstrates more clearly an invasion of suspect's rights, and, in this case, the party that caused the violation is the security services, as their actions have made a third-party aware of the identity of a suspect in an investigation.

In the UK, according to the DPA, organisations may provide other organisations with personal and sensitive personal information about a data subject in some exceptional circumstances (see Part IV of the DPA [10]). For instance, emergency services may

request information on the allergies of a casualty, and of a casualty's relatives, from any organisation that they suspect may have this data, and such organisation may lawfully disclose the data. Accordingly, the police and other public authorities may also request data related to their suspects, based on the same reasoning. Thus, in Scenario 2, the security services and the data-controlling organisation would act lawfully in accordance with the above legislation. However, their actions could seriously impinge upon the data-subject's natural rights, and, quite possibly, their right to privacy. In this scenario, it could have a detrimental impact upon the data-subjects rights concerning the future relations with the data-controlling organisation.

This raises interesting issues about the legal remedies, if any, open to the suspect. In similar circumstances a case could be made that there has been a breach of Article 8 of the Council of Europe's Convention on Human Rights (now enforceable in the UK under the Human Rights Act 1998 [11]). This would be difficult to pursue for a number of reasons. Quite apart from the practical difficulty of knowing that there has been a breach of rights, how the breach has come about, who is responsible and how to prove it (what might be called *evidential difficulties*), there is also the question of the extent to which those responsible might be able to claim exemption from responsibility (which might be called *substantive difficulties*). The right of privacy under Article 8 (like most human rights) is a qualified right, meaning that a public authority is entitled to disregard the right where the interests, among others, of national security, or the prevention of crime and disorder, require. Such an exemption would normally exclude the possibility of the affected data subject being able to pursue damages against the public authority. However, perhaps the correct approach is to regard the exemptions as only coming into effect where they are proportionate. If there is a way to obtain the evidence they require without invasion of privacy and other rights of the suspect, and without the adverse impact the scenario predicts, it is arguable that the public authority should take into account the rights of the suspect, and so to choose the least disruptive method of obtaining the evidence they need. It therefore could be argued that if they chose a method, which invades protected rights, and is likely to cause adverse impacts, the public authority have used an exemption disproportionately, and so should be obliged to recompense the suspect for the harm perpetrated by their choice of method. It is interesting to conjecture to what extent a court would entertain such a claim.

1.4 Research Question, Aim and Objectives

The work presented in this thesis will therefore address the following research question:

What are the improvements to SPIR methods that can be made within an investigative framework, and how can this be evaluated against the current methods to show efficiency gains?

Consequently, the main aim of this research is as follows:

To define new methods for the investigative data acquisition that can preserve privacy of relevant data-subjects, and which have perceivable performance gains over existing methods, and to allow variable parameters to preserve the balance of privacy against performance.

The thesis addresses these considerations with the following objectives:

- 1) Construct a literature review within the PET sphere (Chapter 3).
- 2) Define a set of requirements that data acquisition process must meet (Chapter 4).
- 3) Construct a novel methodology for privacy-preserving investigative data acquisition (Section 5.3 and Section 5.4).
- 4) Propose an evaluation framework suitable to assess performance of novel cryptographic enhancements to retrieval of investigatory data (Section 5.5).
- 5) Investigate parameters that could be used to assess the balance between the privacy and feasibility (Section 6.2.6).

In theory, it is possible to use Private Information Retrieval (PIR) primitives to search databases belonging to third parties, without revealing the search criteria. Thus, in investigative data acquisition, it is feasible to keep the identity of the suspects secret. While a PIR protocol would reveal to the investigators records other than those classified as interesting (the records referring to the suspect), a SPIR

protocol would potentially protect the interests of all the parties involved, since it can:

- Protect an enquiry by hiding the identities of the interesting records and some of the search criteria.
- Protect records kept on database, but unrelated to the enquiry.

However, even efficient SPIR protocols may struggle to handle privacy-preserving requests for investigative data, as the databases involved usually contain a large number of records that are likely to change frequently. Consequently, it will be necessary to investigate possible modifications to existing SPIR protocols that would be suitable for enhancing the performance of these protocols in an investigative scenario. At the same time, the balance between the performance and the privacy must be kept at acceptable level. Therefore, the criteria for selecting this acceptable level will be also defined by this thesis.

The main goals of this work are to put forward a new problem, establishing a *practical feasibility* result for this problem, and, in the process, develop techniques that allow for the scaling-up current SPIR schemes to the size required by investigative data acquisition. This work does not attempt to fully optimise the platform, as this would complicate the presentation of the problem and the solution. The platform should, thus, be mainly viewed as a feasible framework, which may be the basis for further optimisations.

This thesis shows that DPA and The Regulation of Investigatory Powers Act 2000 (RIPA) [12] that was brought in to regulate the data acquisition process in the specific investigatory cases, define controls aimed at protecting individuals being investigated by the public authorities. One of such controls is the requirement for the following roles in the acquisition process [8]:

- **Applicant.** A person that requests the data needed for a specific investigation or operation within a relevant public authority.
- **Designated Person.** An individual responsible for assessing the application for data acquisition that ensures the request is *necessary and proportionate*.

- **SPoC.** This could be an individual or a group of accredited individuals trained in lawful acquisition of communications data and co-operating with third parties.
- **Senior Responsible Officer.** The officer oversees the whole data acquisition process to ensure compliance with the appropriate legislations.

These roles are referred to throughout the thesis.

1.5 Contribution to Knowledge

The work presented in this thesis contributes the following to knowledge:

- 1) Demonstrating the manner in which SPIR techniques can be used to assist public authorities in privacy-preserving retrieval of investigative data from third parties.
- 2) Reducing the problem of investigative data acquisition to a single-database SPIR, thus, allowing for the limiting of the number of records that need to be collected from a dataholder in order not to affect the privacy of a suspect in a considerable way.
- 3) Presentation of a novel methodology for the privacy-preserving investigative data acquisition, IDAP, which is suitable for real-life implementation.
- 4) Creation of a dilution factor that can be used to control the balance between the privacy and performance in a single-database SPIR system.
- 5) Definition of a technique for building complex privacy-preserving enquiries based on multiple selection criteria.
- 6) The novel use of semi-trusted third parties to gain the support of the public for SPIR-based data acquisition techniques.

1.6 Publications

The main publications conducted with this research include:

- Z. Kwecka and W. J. Buchanan, "Minimising Collateral Damage: Privacy-Preserving Investigative Data Acquisition Platform.," *International Journal of Information Technologies and Systems Approach (IJITSA): Special issue on Privacy and Security Issues in IT*, vol. 4, 2011.
- Z. Kwecka, W. J. Buchanan, and D. Spiers, "Privacy-Preserving Data Acquisition Protocol," Proceedings of the IEEE Region 8 International Conference on Computational Technologies in Electrical and Electronics Engineering (SIBIRCON), Irkutsk, vol. 1, 2010, pp. 131-136.
- Z. Kwecka, W. Buchanan, D. Spiers, and L. Saliou, "Validation of 1-N OT Algorithms in Privacy-Preserving Investigations," Proceedings of the 7th European Conference on Information Warfare and Security, University of Plymouth, 2008, pp. 119-128.
- Z. Kwecka, W. Buchanan, and D. Spiers, "Application and Analysis of Private Matching Schemes Based on Commutative Cryptosystems," Proceedings of the 8th European Conference on i-Warfare, Lisbon, 2009, pp. 154-163.
- Research poster *Privacy-Preserving Investigations - A technical solution to allow for legal and ethical data sharing* that was presented during the second annual conference of the Scottish Institute for Policing Research (SIPR) held in Edinburgh.

1.7 Thesis Structure

The structure of this manuscript closely follows the methodology used to draw the final conclusions. Therefore, along with appendices that contain supporting data, this thesis is organised as:

- **Chapter 1 – Introduction.** This chapter outlines a brief context of the research domain, and identifies the key issues. Finally, it presents the research question to be explored.
- **Chapter 2 – Background and Theory.** This chapter provides a more detailed presentation of the background of investigative data acquisition and

of privacy issues in information systems. It also provides an insight into cryptography, which forms the basis of the privacy-preserving data acquisition technique presented in this thesis.

- **Chapter 3 – Literature Review.** This chapter presents an analysis and overview of privacy-preserving techniques that can, in theory, improve the ethics of investigative data acquisition.
- **Chapter 4 – Improving the Acquisition Process.** This chapter defines the requirements for investigative data acquisition process and analyses existing PETs in contrast to these requirements. As a result, a single protocol is selected as a candidate protocol for the acquisition process and presented, along with a number of drawbacks of this protocol.
- **Chapter 5 – Novel Data Acquisition .** This chapter presents the framework developed during this research, in order to explore the research question presented in Chapter 1. This framework is based on the PET protocol chosen and customised within this thesis for the specific task of investigative data acquisition. The chapter also outlines the process of experimentation and simulation, and provides the necessary narrative to place it within a research methodology.
- **Chapter 6 – Evaluation.** This chapter presents the results of the process of experimentation and simulation of the framework performance. It also includes the details on the qualitative evaluation of the framework.
- **Chapter 7 – Conclusions and future work.** This chapter provides a discussion and a summary of the main findings of this thesis, within the main considerations of the research domain. It will also justify the contributions to knowledge, and suggest the future work.

Chapter 2

Background and Theory

2.1 Introduction

This chapter introduces the concepts that are crucial to understanding the field of investigative data acquisition, and associated issues. Consequently, matters relating to the use of electronic evidence and obtaining investigative data from third parties are discussed below. This is followed by an introduction to the concepts of security and privacy in the information systems. Finally, different cryptographic techniques that find use in information retrieval and storage are outlined.

2.2 Digital Forensic

This section provides the background on the field of digital forensics. The definition of the term is followed by the comparison of the digital forensic to forensic science. This is done in order to introduce the concepts of investigation, and digital evidence, which includes data collected during investigations.

The most comprehensive definition of digital forensic reads as follows:

(DF is) the use of scientifically derived and proved methods towards the preservation, collection, validation, identification, analysis, interpretation, documentation and presentation of digital evidence derived from digital sources for the purpose of facilitating or furthering the reconstruction of events found to be criminal, or helping to anticipate unauthorized actions shown to be disruptive to planned operations.

Palmer, [13], pp. 16

Computer-related offences started to emerge when personal computers became easily accessible. Nevertheless, DF was only recently recognised as a separate discipline. While technology is now advanced enough to assist in DF investigations, none of the DF techniques used are court-approved in the UK. This means that the DF evidence cannot be presented directly in court, unlike the physical evidence, such as the results of the fingerprints analysis. Instead, most DF evidence needs to be interpreted and presented to court by an expert witness in order to have any significance in a court case. Also, the lack of proper regulation of DF means that on many occasions electronic evidence is rejected by courts due to its illegitimate origin or alleged mishandling [14].

Little can be done to improve the quality of DF evidence collected from personal computing devices, as the environment is under the full control of the end-user. However, legislations for the corporate world, corporate information systems, as well as communications links, could allow for a greater access to data for investigators [14]. Currently, in order to protect themselves from costly legal suits, organisations often choose to monitor and log the flow of information that occurs under their governance. However, the methods that are used in these activities can render any information that is provided to the investigators, unusable in court. This is partially due to the fact that often the data is logged without the clear and explicit consent of the sender, or recipient, of the data, where by British and European law interception of information transfer without such consent is forbidden [12]. There are also more difficulties in using data from the corporate environment as evidence, as in cases

involving prosecution evidence provided by commercial organisations and their employees, the organisations can be selective about the evidence they provide [14].

During conventional forensics investigations, scientists are often capable of reconstructing events by examining physical evidence; therefore, it should also be possible in the computer environment. However, electronic evidence is quite different from the physical evidence. Traces of evidence in conventional forensic science are often difficult to forge; however, this is not the case in the world of digital forensics. For example, perpetrators can put fake evidence onto a machine belonging to someone else, or claim they become have victims of a Trojan horse [15]. This is partially due to the fact that conventional techniques of storing information are often more tamper-proof than their digital equivalents. For instance, if someone were to erase and overwrite a message written on paper, this process is likely to leave evidence that the information been tampered with. A perpetrator, thus, with moderate skills would be capable of a similar forgery in a digital world without leaving traces of this activity [16]. Within a computer operating system there are certain controls that can be used to trace such actions (such as meta-data, event logs, and so on), but these can also be fooled [15]. Some events, though, that take place in information systems often do not leave any long-lasting physical traces, as they only run in volatile memory [17, 18]. Thus, while monitoring computer systems, as well as in communication links, it can be treated as a potential breach of privacy and it can be used to prosecute the guilty, but it also has a potential to protect the innocent.

During a DF investigation, an examination of the data can also show signs of incidents different to the one being investigated. Thus, handling and examination of the evidence should be performed in a way defined by the appropriate regulatory documents [17]. In the corporate environment, these documents would most likely be a part of institutional security policy, written in accordance to the guidelines provided by national law enforcement. The UK authority, Association of Chief Police Officers (ACPO), provides general guidance for the recovery and handling of digital evidence. They suggest four following principles in working with computer-based electronic evidence (ACPO [17], pp. 4):

- *No action taken ... should change data held on a computer or storage media which may subsequently be relied upon in court.*

- *In exceptional circumstances, where a person finds it necessary to access original data ..., that person must be competent to do so and be able to give evidence explaining the relevance and the implications of their actions.*
- *An audit trail or other record of all processes applied to computer-based electronic evidence should be created and preserved. An independent third party should be able to examine those processes and achieve the same result.*
- *The person in charge of the investigation has overall responsibility for ensuring that the law and these principles are adhered to.*

The common practice of preserving digital evidence is taking complete copies of the source data bit-by-bit using *raw* formatting [19]. This, though, is often a costly and wasteful approach, as, in order to create forensically sound copy of a 200GB hard disk requires 200GB of storage for the evidence, even if there is only 100MB of data on the drive. The introduction of a vendor independent Common Digital Evidence Storage Format is discussed in [19], which allows data from various sources, including hard drives, network traffic monitoring, memory dumps and other files, or data acquired as evidence, to be preserved in a single format. This is the digital equivalent of the evidence bags that are used for physical evidence. Additionally, [20] suggests that evidence derived from server logs and network probes, such as the traffic data collected by the ISPs, should be split into different data formats, and data that is repeated should not be duplicated in the storage. In practice, this would mean that timestamps are stored in a date format, while IP addresses could be stored within 32 bits (four bytes), rather than in the *decimal-dot-delimited* format that requires up to 15 bytes when stored in *raw* format as a string. Also, the IP addresses could be stored for the initial packets in a given TCP connection logged, but then omitted from the remaining log entries. However, these considerations are mostly important to the dataholders. Since this thesis focuses on the mechanism for retrieving a small quantity of records from datasets held by third parties, this is not an issue, but it is worth noting that there are valid concerns related to the storage requirements for safekeeping any potential evidence.

2.3 Investigations using third party data

The public authorities are often required to carry out investigations based on data supplied by third parties. Such investigations may include: benefit fraud enquiries from the Her Majesty's Revenue and Customs (HMRC); solving a crime by the Police; investigating alleged terrorism cases by Scotland Yard; or, gathering health information about a patient at an Accident and Emergency (A&E) department. The process of obtaining third-party records is usually referred to as *data acquisition* [8].

In the UK there are two major legislations that the public authorities can use to justify their request for third-party data. These are The Data Protection Act 1998 (DPA) [10], and The Regulation of Investigatory Powers Act 2000 (RIPA) [12], including its Scottish counterpart The Regulation of Investigatory Powers (Scotland) Act 2002 (RIPSA). In this thesis RIPA and RIPSA are both referred to as RIPA, unless specified differently. The reason for the lack of the distinction is that the matters relating the data acquisition are mostly the same for the UK. This section discusses the aspects of DPA and RIPA that are important to the data acquisition process, along with current data retention practices by Content Service Providers (CSPs).

2.3.1 Data Protection Act

In the UK, the DPA regulates the processing of data on identifiable individuals. Similar regulations also exist in other countries of the European Economic Area, as the DPA was enacted in implement of the European Data Protection Directive 95/46/EC [7]. The DPA provides eight principles for handling personal data and ten conditions governing the processing of sensitive personal data. It should be noted that this legislation is not aimed at regulating all data about individuals, and the scope of principles provided is limited by the definitions of data and sensitive personal data provided in Section 1(1) of the Act. Consequently, the DPA regulates the processing of any data about an individual, which is:

- intended to be processed automatically;

- intended to be recorded in a structured manner allowing for the retrieval of information about an identifiable individual, i.e. as a part of *relevant filling system*;
- a relevant health, educational or public record.

Additionally, the key to understanding DPA are the terms *data subject* and *processing*. *Data subject* is a term widely used to describe an identifiable individual whose data is kept by the given dataholder (referred to as *data controller* in the Act), where *processing* is used to describe any operation on the *data*.

The eight principles of the DPA are:

- *Personal data shall be processed fairly and lawfully*. Thus, any operation on the *data* relating to a *data subject* is done with expressed or implied consent of the *data subject*, unless it is required to satisfy legal requirements of the *data controller*.
- Information should only be used for the original purpose for which it has been obtained.
- *Personal data* needs to be adequate, relevant, and not excessive to the purpose for which they are *processed*.
- *Data controller* needs to make sure that the data is accurate and kept up-to-date when necessary.
- *Data* should be kept on a system only when needed for the original purpose for which it has been obtained.
- *Data subject* should be assured that any *processing* of *personal data* is performed in accordance to the DPA.
- *Data controllers* should regularly evaluate the risk to data and implement appropriate countermeasures if required.
- The data cannot be transferred to a territory outside the European Economic Area (EEA) unless this territory can provide adequate level of protection for

the data. Since, DPA and similar legislations guarantee to protect the rights of the subject only in the EEA, care needs to be taken when data are transferred to territories outside of these controls.

The DPA provides a voluntary mechanism to enable the *data controller* to disclose information on *data subjects* to the public authorities, in circumstances that the *data controller* perceives as reasonable for such disclosure. Some may argue that such an exclusion in the DPA is unreasonable, since the public authorities can then obtain information about their suspects without any court warrants. On the other hand, such a provision is required for life-threatening situations where medical staff or police needs to gather information quickly to protect lives and critical infrastructure. Consequently, the valid uses of this provision possibly outweigh the abuses. Based on the DPA, the public authorities cannot enforce any disclosure of information without a warrant. However, the *data controller* can disclose any information to investigators, if investigators can demonstrate a valid reason for the disclosure.

An example of when the system was abused is found in [21] where an ex-policeman was able to gain access to the database of UK Driver and Vehicle Licensing Agency (DVLA) and obtain postal address of an individual based on registration of a car. The ex-serviceman performed this action in order to help in his private investigation for a missing dog. This was a plain breach of the DPA, however, the consequences of this breach were limited; if the missing dog enquiry was handled through the police channels, it would result in the same information being obtained by the ex-policeman. On the other hand, the DPA voluntary disclosure can be used to obtain a data subject's medical details by a hospital A&E department. This can allow the public authorities to act fast in life and death situations, where even a slight delay may cost the data subject, or another individual, their life.

2.3.2 Regulation of Investigatory Powers Act

CSPs are the most common third party sources of investigative data [22]. Historically CSPs stored communication data of all transmissions taking place for billing purposes. Such data, often referred to as traffic data, included telephone numbers, time of call, duration, and so on. ISPs, a subset of CSPs, used to record similar data for the Internet transaction, i.e. IP addresses; types of packets; some high-level

Internet addresses, such as Hyper-Text Transfer Protocol (HTTP) Uniform Resource Locators (URLs). Consequently, in the past, investigators were able to make enquires requesting details of communication data based on the voluntary provision mechanism of the DPA. However, ambiguity and abuse of this investigative technique showed the need for further regulation of this area [23]. Consequently, RIPA was introduced to regulate:

- amount of information collected by CSPs about their customers;
- amount of time CSPs were allowed to retain this data;
- who, and under what circumstances, can request to see this data without a subpoena.

The rules addressing these issues are only a part of the RIPA act that was introduced in order to satisfy the directives of the European Convention on Human Rights [24]. RIPA set out to control the interception of communications, acquisition and disclosure of: communications data; the use of covert surveillance and human intelligence sources; as well as access to electronic data protected by passwords. Soon after RIPA, its Scottish counterpart RIPSAs was announced, and regulates the general conduct of surveillance in Scotland [25]. According to the Home Office, these acts were supposed to strictly limit the use of covert surveillance techniques, and intrusive intelligence gathering to the most serious crimes. However, it was criticised by lawyers and privacy activists for loosely defining what was meant by the most serious crimes, and the exceptional circumstances that allowed the public authorities a warrantless intrusion of privacy [22]. Although, the history shows that, generally, the rights given to the public authorities by RIPA are not being abused [25].

Under RIPA, a public authority may send a data acquisition notice to a CSP requesting the disclosure of certain traffic data. Unlike, in the DPA data acquisition request, RIPA notices do not require justification being presented to the *data controller*. The *data controller* must then disclose the requested data within a reasonably practicable time, or face a penalty. Since, RIPA requires the collection of a relevant subset of data from a database of the CSP, the requesting party should make a financial contribution to cover the costs incurred by the CSP. All the notices

are then subject to the approval by the senior officers of the requesting party. In the case of police, a RIPA notice must be authorised by an officer of *Superintendent* or higher rank, and for ambulance services the *Director of Operations* must approve the request. Once the notice is authorised it is then processed by a Single Point of Contact (SPoC) within the public authority, who serves the notice to the CSP's SPoC. As the name suggests, the SPoC is an individual, or a group of individuals, that was/were appointed as the main contact points between the organisations. This ensures that an investigator cannot request any data under RIPA without having the appropriate approval, and that the full process must be followed in order to obtain the data. Thus, the process is self-enforcing, which shows that RIPA and related processes attempt to provide high degree of privacy protection.

Finally, RIPA states that both the notice and the communications data should be transferred in a secure manner according to the DPA, in order to protect the information in transit. However, there is a lack of clear guidelines on how such transfers should occur in the code of practice published by the Home Office [8]. On the other hand, if the investigative data retrieved from a third-party by the police is encrypted, under RIPA, the data subject may be required to disclose the information in an *intelligible form* or provide encryption key(s) and tools required to render the information intelligible to the investigators [12].

2.3.3 Data retention

Under RIPA and DPA, a CSP should not retain any communications data any longer than it is required for billing purposes, and settling any consequent billing disputes with its customers. Many CSPs do not require the storing of such information for prolonged periods of time, and, consequently, some communication data is disposed of shortly after the monthly bills are issued. However, after 9/11 the Anti-Terrorism, Crime & Security Act 2001 was introduced which allows CSPs to voluntarily retain communications data for periods of time that could allow the UK public authorities to have sufficient information available to them, in order to protect national security [26]. This act did not make CSPs store any more data than required for the billing purposes, and did not modify any data acquisition procedures stated in RIPA. Instead, the Act simply provided CSPs with an ability to store the data for the periods specified in Appendix A of [26], and maximum of 12 months, without breaching the

DPA, RIPA, and Human Rights of the data subject, even when the data is no longer required for business operations. Thus, the individual *data controllers* could consider storing communication data for prolonged periods of time as *necessary* in relation to the DPA. This act is often referred to as the *Voluntary Code of Practice*.

In April 2009, the UK Government issued a public consultation that was the first step of modernising the current approach to data retention [2]. The main reason for this is that many CSPs (especially ISPs) do not require any communications data for billing purposes, anymore. Nowadays, most ISPs charge a monthly subscription fee for the unlimited access to the Internet. Also, it is likely that this will also become the case for telephony providers as they are shifting their operation towards using Voice over Internet Protocol (VoIP) instead of the Public Switched Telephone Network (PSTN). This shift makes the cost of telephone calls negligible, and most calls would not be billed separately but would be included in a monthly fee. Consequently, in the future, CSPs will often have no need to store traffic data for the operational purposes, and the traffic data will no longer be available to the public authorities. In [2] the Government proposes a solution, where all the CSPs would monitor all the Internet transactions taking place over their networks, and would make the traffic data available, and intelligible, to a centralised search engine, referred to as a *query hub*. Such a search engine could be used by the public authorities in case of RIPA enquiries. If the system proposed by the consultation gets introduced into practice this would modify the rules for data retention and processing towards a more *intrusive* solution.

RIPA, and other legislations governing communications issues, often refer to the term *communications data*, which refers to all data about the communications apart of the content of communications. This approach was first introduced with telephone systems in mind. Thus, communication data would refer to: the number dialled and duration of a call (traffic data); the services paid for by the subscriber (services data); and subscriber address (subscriber data); but not the content data, which is the actual conversation, this is the information carried over the telephone circuit during the call. The division between communications data and content data was relatively easy at the time when all the conversations took place over circuit switched telephone networks. However, the difference between telephone systems and new means of

communication, such as Internet, is the fact that traffic data cannot be easily separated from content data. Thus, one can consider HTTP headers as traffic data, since it is used to control the request and response during Web browsing, as an analogy to the telephone service this would be traffic data, since it is used to establish the communication. On the other hand, HTTP headers contain information about the content being viewed, and consequently can allow the investigators to infer a good deal, if not all, of the content data. While most legislations differentiate between these two types of data, there is a lack of clear definitions [27], opening the way for precedence lawyers.

2.3.4 Commission and diligence for recovery of documents

In most of the data-acquisition scenarios considered in this thesis, the evidence is being requested by the public authorities. Other cases, such as when a private party requests data to be provided for a court case, are deliberately left out from the scope of the work as they require a subpoena. However, a particular concept used by the Scottish Court of Session in disputes between private parties finds use in this thesis. Chapter 35 Section 4 of the Rules of the Court of Session specifies that a commissioner may be appointed to fulfil a request for third-party data made by a party in the dispute. This is to ensure that only the relevant evidence are collected from the *haver* (dataholder) according to the *specification of documents* prepared by the requesting party, while the data not related to the case, especially *haver's* trade secrets, is filtered-out [28, 29]. Similar principles can be found in other legal jurisdictions, for example the *discovery escrow* in the US intellectual property law [30].

2.4 Privacy, and its wider context

This section provides a brief background on the concept of privacy, and the way that matters of privacy can be examined in an Information System. It is also shown that certain security, auditing and surveillance schemes that were introduced to protect a given population, often breach the privacy of individuals, while other systems with analogous aims contribute to the privacy of the population that they cover. Therefore, privacy levels in a given information system are not directly dependant on the purpose, but are related to the design decisions and implementation of the given

system. A previous section has mentioned that in UK there are plans to increase amount of data that CSPs need to collect in order to enable investigators from public authorities to protect the public and the nation, so this section explains the reasons why some of the proposed measures are justified.

2.4.1 Privacy

A given piece of personal information can be perceived as confidential by one individual, while others would not attempt to conceal it. For this reason it is difficult to define privacy. Stanford Encyclopaedia of Philosophy [31] and [32] provide extensive and neutral discussions on this term. According to these sources some lawyers and philosophers argue that privacy is merely a collection of rights available to an individual to protect the information considered as confidential, and as such, does not merit to be legislated on its own. William Parent defends a view of privacy in the domain of personal information that does not confuse the basic meanings of other fundamental terms. He defines privacy as the condition of not having undocumented personal information known or possessed by others, but he also stresses that privacy is a moral value, and not a legal right.

Perhaps, the best definition for use in this thesis is the slightly wider definition by Alan Westin. This definition describes privacy as *the ability to determine for ourselves when, how, and to what extent information about us is communicated to others* (as discussed in [31]). Or as Swire and Steinfeld put it *privacy is providing individuals some level of information and control regarding the uses and disclosures of their personal information* [1]. These definitions of privacy may look ambiguous, however, the best practice of handling private data is to allow the data subjects to have a certain amount of control, as no individual is the same as another. Nevertheless, for legislative reasons, there is a need to expand this definition and create laws that could ensure privacy is being maintained by organisations that have access to personal information. In the previous sections the DPA has been discussed as the UK legislation defined for this purpose, and the following section outlines the origins of the DPA, and different views on privacy by different social groups.

2.4.2 Measuring privacy

One of the first comprehensive guidelines for creation of privacy laws was the US Code of Fair Information Practices, developed in 1973 by US Department of Health, Education and Welfare [33]. The document describes five key factors required to achieve privacy:

- **Openness.** Data subjects should be aware that a system keeping their personal data exists.
- **Disclosure.** There must be a mechanism for the data subject to access their own records, and to find out how these records are used.
- **Secondary usage.** Gathered data can only be used for the purpose it has been collected for, unless there is consent from the data subject to further process the data.
- **Record correction.** If the records are incorrect, the data subject should have the right to request a correction.
- **Security.** Where the reliability of the records for their intended use is required, the data controller must take precautions to prevent misuse of the data.

These, and other guidelines, were later adapted by the Organization of Economic Cooperation and Development (OECD) to form the Guidelines on the Protection of Privacy and Transborder Flows of Personal Data. The OECD is formed by 24 countries, where the creation of the guidelines was to harmonise the efforts of the member countries in creating privacy laws [33]. This, in turn, greatly influenced the European Data Protection Directive and, finally, the UK DPA. Consequently, the OECD guidelines, and the DPA, can be used as guidelines for designing information systems. However, they define only the bare minimum that a system must meet.

The DPA, like most guidelines before it, allows any kind of data collection and processing, as long as the user gives consent to such operation. However, the consent is often only implied, from the fact that the data subject uses a given service offered by an organisation. What is more is that an individual wanting to use a services of a communication service provider in the UK often must agree to the terms of use, that

likely include references to information interception. This also applies to other operations where consent is required. For example, most often banks have similar privacy policies, and for a user to be able to have a bank account, such a policy must be accepted. Consequently, it can be argued that an individual does not often have a choice, and therefore the requirement of consent does not necessarily improve privacy. This confirms claims of [34] that with the evolution of new technologies, the sets of rules proposed in US Code of Fair Information Practices has become outdated. In [34], Marx proposes 29 questions that may be used to help assess the ethics of a given surveillance process. In order to define the framework, Marx identified conditions which, when breached, could violate an individual's rational perception of privacy and dignity. These conditions call for the following in any surveillance system:

- avoiding harm;
- ensuring validity of reasons;
- building trust with the data-subjects;
- giving notice;
- obtaining permission when crossing personal borders.

Marx's framework is not technology specific, in order to keep the framework universal and lasting longer than the ever-changing technologies of surveillance. Marx explains such approach in the following words: *in matters so complex and varied we are better served by an imperfect compass than a detailed map* (Marx, [34], pp. 17).

2.4.3 Privacy in surveillance systems

Western society is subject to many forms of surveillance on a daily basis. Some surveillance activities are overt such as Closed-Circuit Television (CCTV) monitoring, while other are hidden beneath a cloak of customer rewards schemes, or they are being completely concealed from the public eye. One of the surveillance projects that for many years has been concealed from the public was Echelon, which is a surveillance operation monitoring international communications links [35]. Also,

the behavioural advertising scheme by British Telecom, the Phorm, has been kept secret for many months, despite consisting of technology that analyses content of the Internet communications [36]. The UK Government did not react to the breaches of privacy by Phorm, as the scheme formed a useful and readily accessible under RIPA source of surveillance information for the public authorities. This is because, according to RIPA, the public authorities cannot require the ISPs to collect content data for all Internet communications. However, a different situation arises when the ISPs has logs of the content data, as well as a list of interests for every user, collected for business purposes. In such a scenario an investigator could request such data without a warrant. It appears that a significant percentage of monitoring activities are performed in an unethical manner [37]. Often such unethical surveillance systems are created for *the greater good of society* (e.g. [38]), and whilst this is mostly true taking into consideration their aims, sometimes the ends do not justify the means, and society often does not benefit from these monitoring systems. In [39], researchers reported on the effectiveness of a number of CCTV deployments in UK. Their findings showed that only a few of the systems achieved their goal of reducing crime levels. The reason for this was that despite the reduction of crime being the initial objective, the design stage of the deployments was not bound to achieving this objective, neither were the way the systems were managed.

In the digital world, the surveillance measures, equivalent to physical CCTV, are logging, monitoring and auditing. Interestingly, in [1] Swire and Steinfeld showed that the implementation of surveillance measures in information systems does not have to go along with lowering privacy of the users. They argue that auditing and monitoring of information systems are standard procedures, and without these, the privacy of data, could be in greater danger, since in an unprotected system the data could be easily stolen or misused [1]. Even though it easy to understand the worries of the general public where their privacy is concerned [40], well-designed and configured surveillance systems may actually stop breaches of privacy from occurring. Swire and Steinfeld believe that security and privacy are complementary, as there are common goals between the two terms. They both are concerned with stopping unauthorised access, use and disclosure of personal information [1]. However, in order to protect the security and privacy of the data subjects, any surveillance system should adhere to firm privacy rules.

2.5 Cryptography

Cryptography is the science of keeping data secure from eavesdroppers during transit. Data in its unsecured form is often referred to as plaintext, or cleartext. One method of securing data is encryption and the encrypted data is referred to as ciphertext. A ciphertext may be transformed back into the original plaintext by decryption. When discussing cryptography a set of common symbols is used to denote these operations and states of data. In this thesis the following symbols are used:

- M – plaintext.
- C – ciphertext.
- E – encryption operation.
- D – decryption operation.

Using these symbols, cryptographic operations can be written in an algebraic form. Thus, encryption E of a plaintext M , that produces a ciphertext C , is shown in Eqn. 2-1. Similar equations may be used to describe decryption, and other cryptographic operations.

$$C = E(M) \qquad \text{Eqn. 2-1}$$

In the past many cryptographic protocols, also referred to as ciphers, were based on the secrecy of the mathematical functions (algorithms) that were used to encrypt and decrypt messages. This solution did not scale well, as the mass use of any single cipher was impossible, and every group requiring to communicate securely would need to develop a new mathematical function that could not be broken by the eavesdroppers. Nowadays, only the protocols that are open to scrutiny of the public and cryptanalysts are perceived as secure [41]. In these protocols, the secret is protected by the secrecy of the key used during the encryption process, rather than the secrecy of the cipher, itself. In arithmetical notation the key is denoted by K :

$$C = E_K(M) \qquad \text{Eqn. 2-2}$$

This section describes the different types of cryptographic protocols, together with their advantages and disadvantages. In cases where the specific source is not provided, the information is based on Schneier's comprehensive reference book [41].

2.5.1 Classification of cryptographic protocols

Traditionally, many cryptographic algorithms used the same key for both encryption and decryption, or one key could be simply derived from the other. These algorithms are referred to as symmetric algorithms, or secret-key algorithms, since, the key needs to remain secret from anybody outside the trusted domain. Consequently, in order for a number of parties to exchange secret messages they had to first exchange the encryption key. Some methods of dealing of this problem included out-of-band communication, or calculation of a common key by two remote parties based on the Diffie-Hellman (DH) algorithm [42] in a way that an eavesdropper cannot produce a valid secret key. The first solution was practical only for the parties that knew in advance that there will be a need to exchange information securely, and that then could use out-of-band communication means (such as secure post) to exchange cryptographic keys. The second, enabled by the DH algorithm, lacked means of authenticating the remote party, and although the session key for data exchange could be securely communicated between two parties using DH, the cryptographic techniques known at the time did not allow the parties to verify the identity of each other.

In the same document as the DH algorithm was first published, Diffie and Hellman discussed a mechanism that would allow an encryption key to be published to the world without jeopardising the secrecy of any message encrypted under such a key [42]. These resulted in a number of protocols that allowed for asymmetric public-key cryptography to surface [43] including the Rivest-Shamir-Adleman (RSA) protocol [44]. (A Mechanism of the RSA for practical use is described in [45].) Such protocols mitigated the need for exchanging the encryption keys using out-of-band techniques, since different keys were used for encryption and decryption, and they were difficult to be derived from each other by anyone else than the creator of the key. These protocols were classified as asymmetric as the encryption and decryption

keys could not be easily derived from each other, and named public key cryptography because one of the keys could be made public.

The term public key usually refers to the key used to encrypt a message, whereas the key used for decryption is often named the private key. The concept of asymmetric cryptography was exactly what was needed in the world of computer interaction at the time. Thus, asymmetric cryptography is mainly used to facilitate exchange of keys for symmetric ciphers, as well as performing functions such as allowing for authentication and integrity checks. It is this exact combination of the two cryptographic approaches that is used to protect most of the secure transactions taking place on the Internet today.

Cryptosystems can also be classified based on the level of security they offer. Consequently, some cryptosystems may be considered as information-theoretically secure. This term is derived from the information theory developed in 1949 [46], and in the context of cryptography and security would classify a cryptosystem as secure if the original cleartext message could not be recovered from a given piece of ciphertext by a cryptanalyst with no access to the appropriate decryption key. A few decades later Shamir and other researchers argued that *in practice the important distinction is not between doable and the undoable, but between the easy and the difficult* (Shamir, [47], pp. 583). The reason for this was that a number of useful cryptographic algorithms were then (and still are) based on mathematical problems that were (and most of them still are) hard to solve. A problem would be considered hard if the solution could be calculated but the time taken for this calculation would make the results unusable. Thus, the cryptosystems where it is hard to derive the original cleartext message from the ciphertext are considered as computationally secure [47]. To give an example: one-time-pad cryptosystem based on a exclusive-OR (EX-OR) operation between the cleartext message and the key, equal in length to the message, would be classified as information-theoretically secure, whereas the RSA cryptosystem is classified as being computationally secure.

2.5.2 Authentication, Integrity and Non-repudiation – Public Key Infrastructure (PKI)

The uses of cryptography are not limited to encryption and decryption of messages, and there are a number of additional functions that are needed to handle secure information processing on computers. Thus, the identification of the message sender can be performed using authentication mechanisms. Integrity checking can detect any changes to the message, after it has been formed by the sender, so that it can be verified that the received message is valid. Finally, thanks to the nonrepudiation mechanism, it can be proven that a given message was originated by a given sender. These functionalities allow the communicating parties to trust in the authenticity of information received, in similar way that signatures, seals and tamperproof envelopes used to do it in the physical world. Consequently, any system designed to transfer information between remote parties must be capable of performing such checks. However, these functions do not have to be limited to verifying encrypted communications, as they are also centric to watermarking and copyright protection of digital goods [48-50].

These functions are possible thanks to a mix of the public and private key cryptography, and also cryptographic hash functions. Cryptographic hash function is a deterministic procedure that converts input data into a fixed-length bit string (or array), referred to as a hash signature. Such a hash signature can be used to uniquely identify the data that was used as the input to the hash function since any commonly accepted cryptographic hash function have the following properties:

- A small change in the input data results in a large difference in the output.
- Hard to reverse.
- Hard to find two different sets of input data that produce the same output.
- Easy to compute.

Thanks to these properties, the hash functions can be used to verify correctness of the input data if the hash signature of the valid input is known. The two most commonly used hashing protocols are the 128-bit MD5 and 160-bit SHA-1. MD5 has been found vulnerable to a number of theoretical attacks, with a successful attack published in 2008 [51], but still it is widely used. SHA-1 has also been found to be

weaker than initially expected, and a collision, a different input data that produces the same results, can be found in 2^{63} operations [52]. For these reasons it is recommended that SHA-2 (or stronger) is used in the applications that are currently being developed [53].

2.5.3 Operations on encrypted plaintext

In certain cryptographic protocols the mathematical operation on the ciphertext has a regular effect on the plaintext. This property is referred to as homomorphism. For example, multiplication of a ciphertext created with the RSA protocol will result in a multiplication of the plaintext. Thus, if the ciphertext of RSA is multiplied by an encrypted number two, after it is decrypted, the value of the plaintext will be twice the original plaintext.

Another homomorphic cipher is ElGamal that also allows for multiplication of the plaintexts [54, 55]. Some other homomorphic ciphers can perform the addition of encrypted plaintext, such as Paillier [56]. These protocols, have already found use in verifiable electronic voting systems [55], and other applications which require privacy and security. However, only recently, a homomorphic cipher which can perform both addition and multiplication was invented [57]. This cipher has not matured yet, but once it passes the scrutiny of peer review, it should be capable to securely evaluate any function (or circuit) over a ciphertext. Consequently, many novel privacy and security solutions could be based on this cipher, and their scope can only be limited by a poor computational performance of the cipher.

2.5.4 Commutative Cryptography

Many cryptographic applications employ sequential encryption and decryption operations under one or more underlying cryptosystems. The reasons to sequence (cascade) different cryptographic schemes together include: strengthening the resulting ciphertext; and achieving additional functionality, which is impossible under any given encryption scheme on its own [46, 55]. A basic cascable cryptosystem can consist of a number of encryption stages, where the output from one stage is treated as the input to another. In such a basic cascable cryptosystem it is necessary to decrypt in the reverse order of encryption operations. However, a special class of sequential cryptosystems – commutative cryptosystems – allows for

the decryption of a ciphertext in an arbitrary order. In conventional cryptography when a plaintext message is encrypted with two different cryptographic functions E_A and E_B , the resulting ciphertext will be different depending on the order of the key application (Eqn. 2-3).

$$E_A(E_B(M)) \neq E_B(E_A(M)) \quad \text{Eqn. 2-3}$$

For most cryptographic applications, this is a desirable behaviour, as it improves the security of the plaintext and the encryption keys. However, commutative algorithms are characterised by the opposite property:

$$E_A(E_B(M)) = E_B(E_A(M)) \quad \text{Eqn. 2-4}$$

Most implementations are computationally expensive. Some are on par in terms of performance with RSA. The commutative encryption protocols, in a similar fashion to the homomorphic encryption protocols, can bring a good deal of benefits to the areas of privacy and secrecy [47, 55]. The property shown in (Eqn. 2-4) makes these protocols an ideal choice for testing inputs for equality without revealing these inputs, which will be expanded upon in Chapter 3.

A typical example an application for the commutative cryptography is the Three-Pass (3Pass) protocol designed to enable two parties to share a secret without exchanging any private or public key. The 3Pass protocol can be described using the following physical analogy:

1. Alice places a secret message in a box and locks it with a padlock.
2. The box is sent to Bob, who adds his padlock to the latch, and sends the box back to Alice.
3. Alice removes her padlock and passes the box back to Bob.
4. Bob removes his padlock, and this enables him to read the message inside the box.

Figure 2-1 Analogy to the operation of the three-pass protocol

A more formal, graphical notation of this protocol is shown in Figure 2-2. Using this protocol Alice and Bob can share a secret without sharing a key first and without using a PKI infrastructure. This protocol is aimed at providing an alternative to public-key encryption and DH-like key negotiation protocols. 3Pass, though, has never been widely used in this way since it is susceptible to man-in-the-middle attacks [58] and is less efficient than RSA, a common choice public-key protocol [44]. However, related concepts are commonly referred to in the information sharing [59] and in the information retrieval [60] solutions.

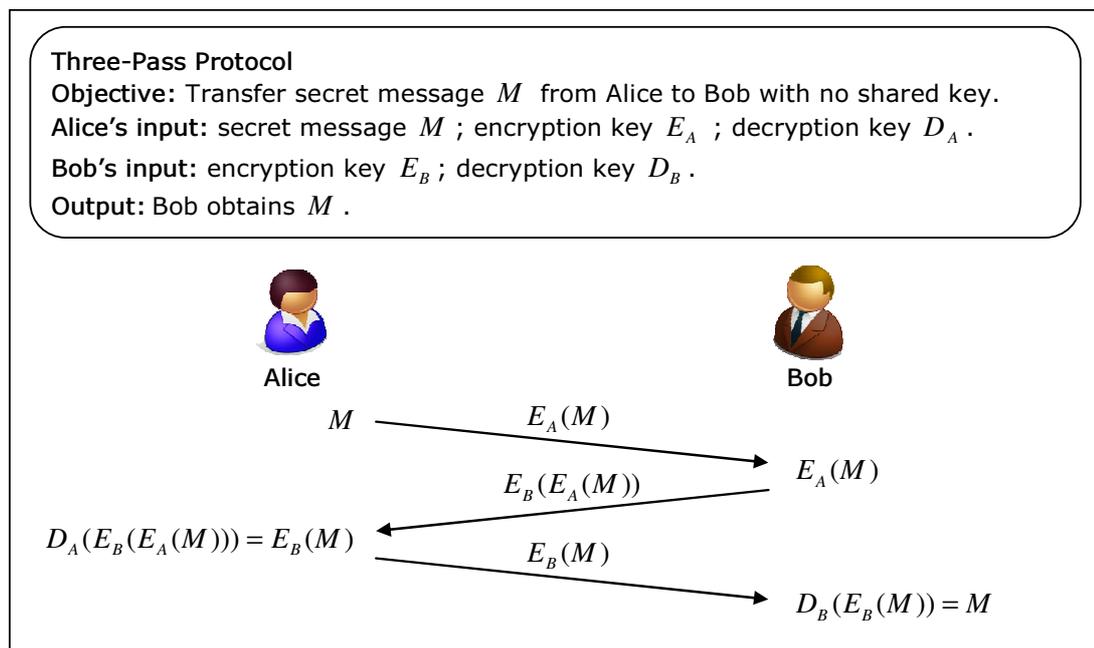


Figure 2-2 Three-pass protocol operation

Commutative algorithms include:

- Pohlig-Hellman is commutative for keys based on the common prime p [47]. Since Shamir was the first person to propose using this algorithm in this way it is often referred to as Shamir's commutative algorithm. Also, RSA can be modified to work in a similar manner due to the link between Pohlig-Hellman and RSA protocols.
- Massey-Omura has improved the above algorithm by performing operations in a specific Galois Field $GF(2^n)$ [61]. This allowed for faster realisation of the cryptographic operations than in case of Pohlig-Hellman algorithm where the operations are performed on the $GF(p)$ [62].

- ElGamal [54], can be used to form a semantically secure commutative algorithm [55] with modification of the universal re-encryption of the plaintext under this protocol [63]. Another way to build a commutative algorithm from ElGamal is discussed in [64].

2.5.5 Cryptanalysis

Cryptanalysis is the science of breaking the security of cryptographic ciphers in order to decrypt a specific ciphertext, or to obtain a cryptographic key used for a certain purpose. According to Swenson in [65] there are a few forms of possible attacks on any given cryptosystem. Ciphertext-Only Attack takes place when the cryptanalyst have access only to the ciphertext, and is looking to decrypt the message, and find the key that was used to encrypt it. Thus, all protocols must withstand this kind of attack as it assumed that the ciphertext can be made, or may become, public. In many systems it is possible for the cryptanalyst to obtain the plaintext associated with a given ciphertext. If this is the case such an attack is referred to as Known-Plaintext Attacks, and the objective is to derive the key that was used to encrypt the messages. A similar form of attack, listed separately by Swenson is the Probable-Plaintext Attack, is where the cryptanalyst has got a fairly good idea what certain parts of the ciphertext contain, which allows for easier deciphering of the message and deriving the decryption key. For example if the cryptanalyst tries to analyse a piece of a source code the ciphertext would contain a large amount of text that are a part of the programming language used. Known-Plaintext Attack helped in decoding the Enigma code during the World War II [66].

Chosen-Plaintext Attack is the most powerful attack. In this attack a special plaintext, prepared so that its ciphertext could reveal certain information about the key used, is fed into the encryption process. It is a powerful type of an attack; however, a careful design, implementation and exploitation of cryptographic technologies should prevent a possibility of such attack. A variation of this attack, suggested by Schneier in [41], is the Adaptive-Chosen-Plaintext attack, where the cryptographer may modify the plaintexts during a chosen-plaintext attack based on the results obtained in the previous attack. Finally, if the cryptanalyst have access to the decryption mechanism, i.e. in form of a decryption process on a computer, or is capable to eavesdrop a plaintext obtained from ciphertext submitted to a given

process, a Chosen-Ciphertext Attack may be rolled-out with an aim of obtaining the decryption key, or reverse engineering the decryption process.

As an addition to this list [41] provides two forms of attack not discussed earlier. Chosen-Key Attack is an unusual attack, in which the cryptanalyst has some understanding of the relationship between different keys. In [41] Schneier discusses a chosen-key attack against modification of Data Encryption Standard (DES) protocol, where the key is rotated two bits after every round. Such attack proved to be efficient; however, impractical since the variability in rotations of the key in DES mitigates any chances of rolling-out this kind of an attack on this protocol. Rubber-hose cryptanalysis (or Purchase-Key Attack) is often ignored by the system designers, but it is one of the most powerful attacks. It is a form of social engineering attack, where the cryptanalyst forces, or bribes, someone to deliver the key.

The classification of attacks on systems using cryptography alone shows the significance of the system design and implementation. The first form of attack can be rolled out against virtually any system, however, the other attacks are often a result of implementation error or procedural error in the way the protocol is being used [66].

2.6 Conclusion

In UK there are two major legislations that regulate collection, retention and release of personal data: DPA and RIPA. DPA regulates the *processing* of personal data that is stored electronically, or in an organised fashion that allows for retrieval of data about a specific individual. This legislation aims to guarantee lawful and fair processing of the *data*. Under DPA the *data controller* can voluntarily release personal data if required by certain public authorities. The provision can be used by police investigators, and the intelligence services, as well as accident and emergency staff in health care.

RIPA has been introduced to regulate investigatory powers that were frequently abused. It permits CSPs to collect and retain data, other than transaction content data, if it is needed for the system maintenance or billing. Under RIPA the public authorities have the right to request data from CSPs without subpoena or any

justification, and the data controller must provide the requested information or face penalties. Thus, RIPA, in comparison to DPA, gives investigators better access to data.

Once always-on Internet connections started to emerge, telecommunication companies introduced all-inclusive call tariffs, and the logging of traffic data was no longer required. Therefore, CSPs had no longer an incentive, nor legal right, to retain this data. This could cause UK investigators to lose one of the commonly-used investigative data sources, and the UK government introduced the *Voluntary Code of Practice* that allowed CSPs to retain this data for up to a year without breaching the DPA. Later, in April 2009 another modification to the laws governing the personal data was proposed by the Government. The proposal suggested the necessity to create a distributed database for all the traffic information from the UK-based CSPs, including data collected from transit connections. Consequently, the consultation document confirmed that the assumptions made during formulating the aims of this thesis were right, and the UK Government is looking into introducing new more intrusive legislation to replace and/or extend RIPA, and to widen the investigative rights.

This chapter also discussed the concept of privacy. Unfortunately, there is no clear definition of privacy, and some argue that, in legal terms, privacy is just a collection of rights, and not a legal term, while others recognise privacy as a moral value, and not a legal right. However, most definitions have a common factor in the notion of control that an individual should have over the way their personal data is used. One of the first instances of documented privacy guidelines was the US Code of Fair Information Practices developed in 1973 by the US Department of Health. This code identified five key areas in fair information processing: openness; disclosure; secondary usage; record correction; and security. Despite the code being prepared more than 30 years ago, the legislations in the UK and Europe, including the DPA, widely inherit from it. However, the code of practice, and the DPA itself, are not specific enough to assist system designers in the building and evaluating systems that respect privacy. The DPA can thus merely assess whether the user is provided with an appropriate level of control over their own personal data. For this reason another technique of assessing privacy in information systems is necessary. In the case of

surveillance systems, some researchers proposed a series of questions that can be used to verify the aim, and the means of the system in respect to privacy.

Any process that is used to collect data that may be required to be presented in front of court of law should follow evidence handling guidelines. The guidelines respected by the UK public authorities, and especially police operations, is the Good Practice Guide for Computer based Electronic Evidence put together by the ACPO. There are four principles presented in this guide: no action of the process should alter the evidence; when the evidence needs to be accessed directly or destructively tested such operation needs to be performed by qualified person that can later give evidence in court; strict audit trail needs to be maintained; and the law and these principles must be adhered to at any time when handling the data.

Finally, the chapter provided background to cryptography and basic cryptographic techniques. Thus, the symmetric protocols that employ the same key for encryption and decryption operations, and asymmetric protocols that use different keys for these operations were compared. The differences between symmetric cryptography and asymmetric cryptography were outlined together with the description of the way in which they allowed ecommerce and modern secure communications to exist. This chapter also detailed different types of attacks that could be employed against a cryptographic protocol.

Chapter 3

Literature Review

3.1 Introduction

This chapter reviews the research related to the most important concepts in the field of privacy-preserving investigative data acquisition. Section 3.2 identifies and discusses the research closely tied with the subject of this thesis, while Section 3.3 provides an insight into the primitives referred to by the research presented in Section 3.2.

The literature related to technologies and techniques allowing individuals to protect their privacy and control the data that relates to them is discussed. It is shown that despite the relevant privacy legislation, mainly the DPA, giving the consumers an option to opt-out from being a data subject in a system; financial and convenience factors may force the consumers to opt-in and use a system that they consider as privacy intrusive. Therefore, individuals that wish to preserve their privacy reach for PETs. Anonymous Internet browsing can be achieved by employing onion routing

techniques, while the identities of buyers (and, in some circumstances, the type of goods purchased) can be hidden from the seller with PIR and private comparison protocols.

The current privacy preserving measures during investigations conducted by the public authorities are shown to be limited to policy-based controls. As discussed in this thesis, these controls are not capable of hiding the identity of the data subject being investigated. Related research shows that it is possible to create surveillance systems with a set of hidden criteria that triggers and alerts if a monitored individual performs a forbidden action. It is also possible to create a privacy-preserved blacklist that informs the investigators about transactions performed by listed suspects.

Data collected during investigations needs to be stored with the same controls as digital evidence. It is also important to ensure that the source of the data is valid and has not been tampered with. For this reason literature relating to preserving of digital evidence is also reviewed in this chapter. Techniques for minimising amount of storage required for the evidence, and allowing data to be verified for authenticity are discussed.

3.2 Privacy-Respecting Investigative Data Acquisition

Chapter 2 introduced the investigative data acquisition process in UK. This section provides information on on-going research in this field.

3.2.1 Privacy self-defence

Research suggests that individuals have different attitudes to privacy depending on the service that they use, and depending on the organisation that is collecting the data [67]. Surprisingly, some individuals that report high levels of concern about privacy, do not consider giving away their personal data in exchange for services and goods as a serious privacy issue [67, 68]. As discussed in Section 2.4 privacy is a complex matter, thus, privacy decisions should be left to the individuals they concern. Thus, the consent of the data subject is important to privacy guidelines [33] and legislation [7, 10]. This would suggest that the best privacy defence measure available to individuals would be an opt-out mechanism that all organisations must provide under the relevant legislations. However, it should be noted that even after opting-out,

some organisations, and, especially public institutions, will retain the right to keep certain private data indefinitely [10]. Research shows that there are situations where the user must opt-in and cannot opt-out. Such situations are usually enforced by convenience and economical, rather than legal factors [20]. An example is the prepaid travel card for London public transport – the Oyster card, where the movements of a person paying with the card can be easily tracked, and as [20] states, the data collected by the back-end of the Oyster system is already being used by investigators from public authorities. Although it is possible to travel around London using cash, paying with the card works-out cheaper than with cash. There are ways to purchase this top-up card without registration, which requires paying a small deposit fee [69]. While this is inconvenient, once the serial number of the Oyster card belonging to an individual is identified using another means, such as CCTV footage from the time it has been topped-up or purchased then the movements of the data subject can be traced-back in the system. Similar considerations apply to access to mobile networks and Internet, as well as other services desired by consumers. Consequently, choosing not to opt-in (or choosing to opt-out) is sometimes not an option and, for this reason, individuals looking to have their privacy protected need to use other means to achieve privacy.

Onion routing networks, such as The Onion Router (TOR) allow for anonymous Internet browsing. This is an interesting approach for privacy that conscious users can use to protect their identity. The infrastructure is based on a series of relays that TOR clients can use to route their requests. These relays are simply network nodes belonging to clients that allow other clients to use their bandwidth in order to create a coherent anonymising network. Since the traffic between the client and relays within the network is encrypted, any request made through TOR to an Internet server can only be traced-back to a relay that executed the request, and not to any particular client. Even the participating nodes are unable to trace-back the request, since the requests are re-encrypted at each relay in the path. Thus, the IP address of a requestor is safely hidden in the population of active TOR clients [70]. However, there are transactions that require at least some form of authorization of the participants. As an example, a website may require all the visitors to be of a legal age, in order to enter. Microsoft researchers suggest that Identity Metasystem, that manages privacy during authentication and authorisation, could be the answer to this concern. The Identity

Metasystem can be based on the data minimisation principal (also discussed in Section 3.3.1) and limit the information released to the minimum required to perform a transaction. Such a system would work as a middleware for all identity-related interactions [71].

Some solutions from the area of Multi-Party Computation (MPC) could prove to be more effective in protecting the privacy of the Internet transactions (and in physical authentication as well). Thus, as an example, systems allowing for anonymous digital payments exist in enabling privacy-preserving purchases of electronic goods and services. The digital money, or *ecash*, has been designed in a way that makes the buyer untraceable, as long as there is no fraudulent attempt to reuse the *ecash* in another transaction. This ensures that the customers can make anonymous transactions, as long as they do not try to cheat the system [72]. Unfortunately, the *ecash* does not address the fact that the seller will know which goods have been purchased, and when. This can potentially help the seller to profile an anonymous buyer or inference some information about the buyer's identity, and thus, has been addressed by the protocols described in [73]. It is important to note that stopping the seller from finding out this information can possibly stop the seller from optimising sales. Consequently, the privacy of the buyer is protected at a cost to the seller. However, the buyers can then voluntarily provide some information, or feedback, to the seller. This approach is *in-synch* with the true spirit of privacy, where the individual described by the data is in control of the data.

3.2.2 Privacy Controls in an Investigative Scenario

Technology-based solutions to protecting privacy of potential suspects are not commonly used in practice. The most widely deployed controls in this area are processes enforced by data protection and human rights legislations. In [8] a code of practice for using the investigatory powers granted by RIPA is provided. It specifies the process for granting the authorisations and giving data acquisition notices, that involves four different roles:

- **Applicant.** A person that requests the data needed for a specific investigation or operation within a relevant public authority.

- **Designated Person.** An individual responsible for assessing the application for data acquisition that ensures the request is *necessary and proportionate*.
- **SPoC.** This could be an individual or a group of accredited individuals trained in lawful acquisition of communications data and co-operating with third parties.
- **Senior Responsible Officer.** The officer oversees the whole data acquisition process to ensure compliance with the appropriate legislations.

Out of these four roles, the designated person is delegated to protecting the rights of the data subjects. Consequently, the whole process is organised to ensure that the data acquisition request is *necessary and proportionate*. The key term here is *proportionate*, as it states that the amount of the potential collateral damage caused to the data subject is justified by the objectives of the request. Furthermore, as described in Chapter 1 the collateral damage can (and is likely to) occur, since in order to obtain data about an individual, the identity of this individual must be revealed. This is the case if only process-based privacy controls are used to safeguard privacy. However, the research discussed in [74] suggests that there are privacy-preserving primitives that can be used to provide greater privacy levels in the investigative scenarios. The suggested primitives are:

- *asymmetric equality* allowing two parties to compare their inputs without revealing them;
- *split comparison* providing the participants with the ability to compare inputs, once again keeping them secret from each-other;
- *equality with selection* that can be used to provide the requestor with a different output depending on the result of an *asymmetric equality* test.

These primitives are combined in [74] to form a system for privacy-preserving electronic surveillance, that allows the tracking of electronic transactions performed by individuals listed as potential suspects. This size of the set (referred to as n) containing identities of the suspected individuals is assumed to be much smaller than the population (denoted by N). The schemes described in [74] assume that it is tractable to perform $O(N)$ operations (270 million, the population of United States, is

the number used), but it is also noted that performing this many operations is impractical. Thus, trading suspect privacy for speed is discussed. Each asymmetric equality test performed to check whether a given identity belongs to one of the n suspects would typically use two unique sets of public keys. However, as [74] suggests, if these keys are reused in more comparison rounds, the communication cost, as well as the number of computations required, are greatly reduced. This is also the technique suggested in [60] where a system for selective sharing of information between parties is discussed. The main drawback of this technique is providing the potential perpetrator greater scope for a known-plaintext attack. Also, once the key is broken, all identities encrypted by this key are revealed.

Another approach to ensure that investigators or auditors can only review actions taken by certain individuals already considered as suspects, is proposed by Biskup and Flegel in [75, 76]. Their system can be used as a pseudonym-based privacy-preserving middleware for audit software. The identities of the users are hidden using pseudonyms unique per identity, or per transaction. While the first case allows auditors to identify suspicious activity of an individual hidden by a pseudonym, the second case makes it impossible to study (or profile) actions performed by individuals, thus providing full anonymity to the users. Auditors can however reveal identities of the pseudonyms after a warrant is given for the given pseudonym or individual, while the system reacts to the transactions taking place by automatically revealing identities of individuals that perform a number of prohibited transactions. The system is based on Shamir's approach to splitting a secret based on polynomial interpolation, so that the core of the system is not based on cryptography. Still, these systems are impractical for large-scale implementations, where it is not feasible to associate each transaction with a unique and tractable pseudonym. Certainly, such a system could not exist in scenario involving a population in the region of $N=270$ million, as proposed in [74]. As Kantarcioglu and Clifton describe it, *privacy is not free* [77] and keeping private information secret requires many computations and many communication rounds.

Researchers note that investigators can easily trick systems into providing data about innocent individuals by simply placing these individuals on the privacy-preserved lists of potential suspects. The solutions proposed in [74] and [78] suggest that a

warrant signature system is needed, where the party in-charge of data acquisition warrants signs a request to assure the participants that the investigation is authorised.

An interesting problem is considered in [77], where the government wants to split users of a certain third party system based on secret classification criteria, while the privacy and equality advocates want to ensure that the criteria is fair and that no data about the users is provided to the government agencies apart of the results of the classification. These requirements can be considered as contradictory. However, [77] shows that it is possible to achieve such a system by the use of commutative and asymmetric cryptography, and one-time pads. This illustrates the power of sequencing different encryption schemes to achieve a functionality that would be impossible to achieve using a single encryption scheme alone, as discussed in Section 2.5.4.

3.3 Privacy-Preserving Primitives

3.3.1 Privacy-Enhancing Technologies

PET is the common name for *a range of different technologies to protect sensitive personal data within information systems* (Koorn, [79], pp. 1). Such technologies find use in all types of information systems. The discussion in [79] deals with typical scenarios where the owner of the data has an incentive, such as a required legislative compliance, to provide privacy to the data subjects.

There are a number of conventional privacy controls that can satisfy the current legislative requirements, and although these are not as exciting to the academic community, they are still valid solutions to the privacy concerns. Overall, they have been tried and tested over the years, and for this reason they are often compliant with the security standards that dictates data processing in many organisations. In [79] the following types of conventional PETs have been identified:

- **General PET controls.** These are the controls that can be implemented with technologies similar to those used in data security. For this, the privacy is treated as the highest level of security, and in this way is a well-defined security policy where its controls protect the privacy. This type is further split into:

- *Data Minimisation.* Analogically to the principal of least privilege in security, a minimal access to data should be given to any requestors. If the requestor needs to know whether a given data subject is an adult, yes or no answers would suffice to fully answer this question, without the need to provide the requestor with neither the age of the data subject nor the date of birth.
- *Authentication and Authorisation.* These should really be treated as prerequisites for any system that carries data that is not publicly available. Without these, other controls, such as the data minimisation mentioned cannot be deployed.
- *Quality-Enhancing Technology.* Part of the requirements for fair information processing [33] and DPA is ensuring the correctness of the data. This can be done by improving the data collection mechanism, as well as allowing the data subjects to view and correct the data about them.
- **Separation of data.** This control splits a data-source into two or more domains, where the personal data that carries information such as a name and an address is stored in the identity domain, and the other personal data is stored in another domain against the pseudo-identity that was derived from the real identity. These domains are linked by *identity protector* software that enables only privileged users to restore the relationship between the data-records in the different domains. Consequently, with this control applied, the personal data can be analysed without revealing the identity of the data subject to the analyst.
- **Privacy management systems.** This type of controls is the least mature of the conventional methods presented in this section. It introduces software that ensures automated enforcement of the privacy policy. Such software intercepts any transactions that involve personal data, and tests these against the privacy regulations, which might include the privacy policy and privacy preferences of the data subjects that the transaction concerns.
- **Anonymisation.** This is a similar approach to the *separation of data*, where the difference is that the pseudo-identities cannot be linked back to the real identities of the data subjects, nor can the identity of the data subject be

inferred from the anonymised data. Thus, the process of anonymisation transforms personal data into data that can be freely processed without privacy controls.

It is worth noting that, despite the anonymisation techniques being valued the most by [79], these cannot be applied in systems where processing of personal data is necessary, unless the personal data is only processed on the input to the system, and it is automatically anonymised on writes to the database.

At the time of its publication [79] was an authoritative guide to PETs for decision makers, but it fails to mention PETs such as mixnets [55, 63], crowds [80] and PIR [81]. Thus, there is another type of PETs, which this thesis defines as:

- **Identity hiding.** Whether it is hiding the identity of the interesting record being retrieved from a database in a larger group of records (PIR) or hiding the identities of individuals that committed an action by a group of individuals (crowds and mixnets) it is possible to provide an additional level of privacy by hiding the target or an originator of an action in a larger group. Such solutions are usually computationally-expensive, thus, they are more likely to be utilised by individuals wanting to improve their privacy, rather than organisation seeking to protect the data subjects.

The *Identity Hiding* techniques can, most likely, perform on par with the *Anonymisation* techniques in terms of effectiveness in privacy protection. Since, the privacy-preserving solution to investigative data acquisition is unlikely to be found in the classical PET technologies described in [79], as these would be in use by now, the remainder of this chapter focuses on the *Identity Hiding* techniques that can be used in this domain.

When discussing PETs and privacy-preserving operations on data, it is important to note the distinction between Private Data-Mining and PIR. Both are well researched subjects, however, the first term – private data-mining – is usually used in relation to obtaining anonymised, statistical data rather than retrieval of individual records as it is the case with PIR. While some techniques used in various available approaches to private data-mining can be modified and reused in information retrieval and vice-versa, these two primitives are dissimilar in objectives.

In the field of statistical data-mining, researchers have developed a number of techniques that permit operations on a subset, or cross-section, of datasets [82-84]. In [84], Agrawal and Srikant suggest a technique based on perturbations, where the larger the perturbations, the greater the level of privacy in the system, but such a technique can result in a loss of information. However, Agrawal and Yu ([83]) show that this is a natural trade-off between accuracy and privacy, similar to those caused by adding noise to data that is then approximately removed from the output [82]. An interesting approach is based on k -anonymity models [85] that ensure that any attempts to link a given record to the data subjects it describes result in at least k different identities being returned. Thus, contrary to the security where any leak of information may be unacceptable [81], privacy can be achieved by hiding the data subject in a larger group of individuals. Finally, a system that does not lose precision can be achieved by employing primitives from the area of MPC (Multi-Party Computation), but, in order for such schemes to be feasible, they often need to make use of an extra party in the protocol – semi-honest party trusted not to collude with other participants – otherwise the computational complexity of the protocol may be too high [82].

3.3.2 Multi-Party Computation

MPC (Multi-Party Computation) allows a number of parties to engage in a protocol that enables them to compare their secret inputs, or to compute a function, without revealing these inputs. It is used in scenarios where no trusted third party exists. Therefore, it can be used to solve a function, such as $f(a,b)$, where a is Alice's input data and b is Bob's input data without the need to reveal these inputs [74, 86]. A classic case is Yao's millionaires' problem, where two millionaires seek to compare their fortunes without revealing the exact figures involved [86]. Yao provides three different solutions to the problem, giving the basis for the multiparty computation of two different parties. He also specifies a method to scale up his techniques in order to allow computation for n different parties. However, schemes designed specifically to handle computation between n different parties were later introduced by Goldreich, Micali and Wigderson [87].

MPC protocols are largely based on the same functions as common encryption schemes, and therefore most have strong theoretical underpinning [88]. However, it

should be noted that the security of any MPC protocol strongly depends on the function that a given protocol is designed to evaluate. Thus, if Alice and Bob engage in a protocol that can evaluate a function $f(a, b) = a \times b$, the party that learns the result can also calculate the input from the second party without breaking the protocol. Consequently, some functions cannot be evaluated privately. For these an MPC protocol can obfuscate the process and hide a security vulnerability rather than solve the underlying problem.

Many MPC protocols, though, are characterised by an exponential growth in the computational complexity with linear growth of the number of records to be processed. These protocols are often impractical while working on large datasets, such as those containing ISP data or health records. The protocols that are characterised by a linear increase in processing time as a response to increased number of records or their size are referred to as *efficient* protocols.

3.3.3 Sharing a secret

An interesting primitive that is commonly used in the field of privacy-preserving information retrieval deals with sharing a secret with a number of parties. In this kind of information sharing schemes, a party (Alice) wants to share a secret with another party (Bob), only if the board of trustees agrees that Bob should have access to the secret [59]. An alternative description of the problem is that a number of parties want to lock their secret, so that it can be retrieved only when they co-operate [89]. For performance reasons, the secret is usually a small piece of data, however, it can be used to store a secret key to a larger dataset, and thus, such schemes can be used to pass control of any asset from Alice to Bob. Khayat in [59] employs the 3Pass primitive described in Section 2.5.4, to propose a secret-sharing with a board of trustees scheme. Figure 3-1 illustrates the encryption operation of this scheme, while the decryption is done in an analogical way (see Figure 3-2). A trustee may leave the scheme by removing his encryption from the ciphertext held by Bob. It is also possible for a new trustee to be added to the scheme, if, for example another trustee is leaving the scheme and a new trustee needs to be appointed. Unfortunately, Khayat has overlooked issues that could arise from an implementation of this secret sharing scheme in real-life scenario, such as:

- Death of a trustee.
- Betrayal.
- Corruption of the ciphertext.

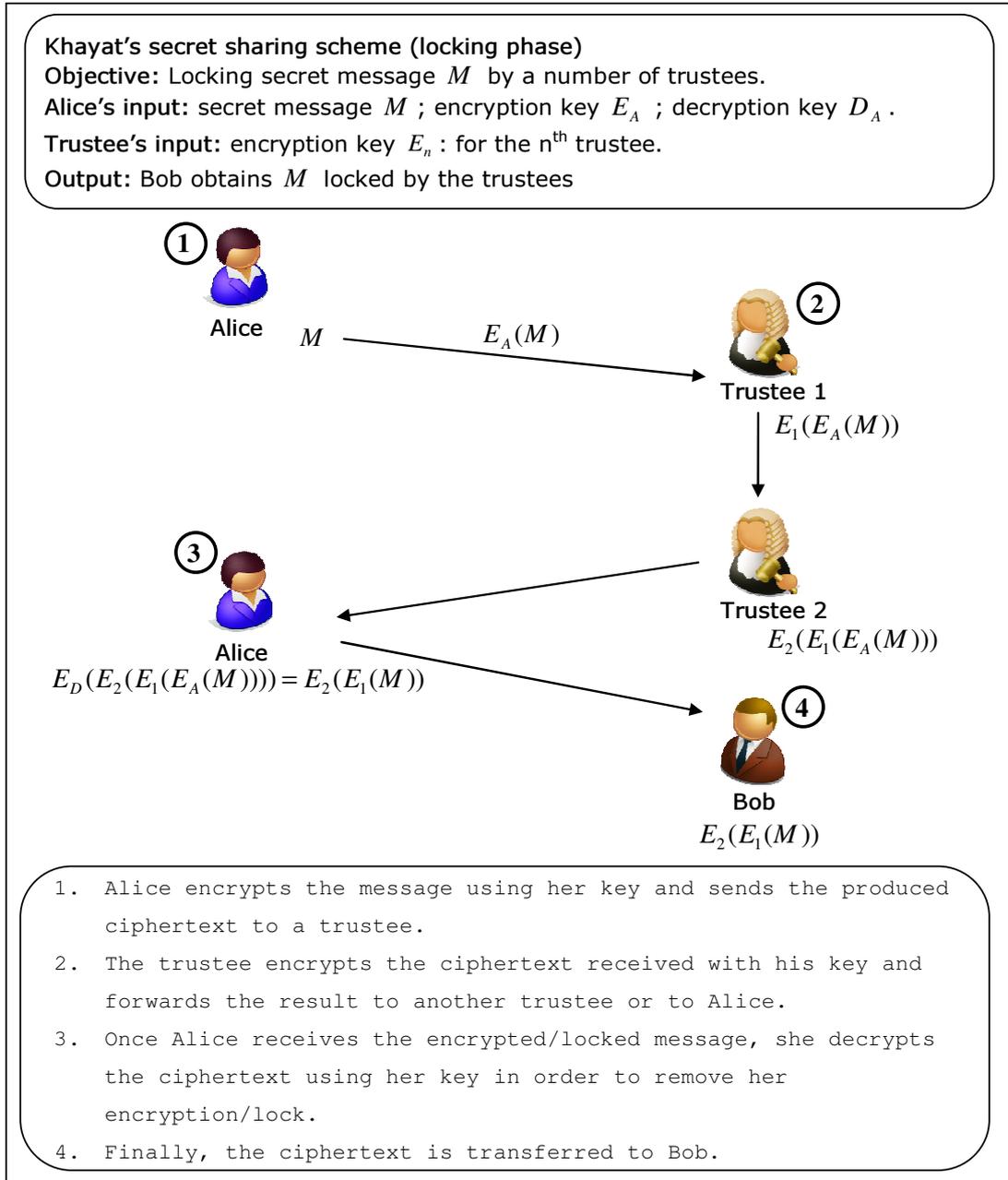


Figure 3-1 Locking a secret under Khayat's secret sharing scheme

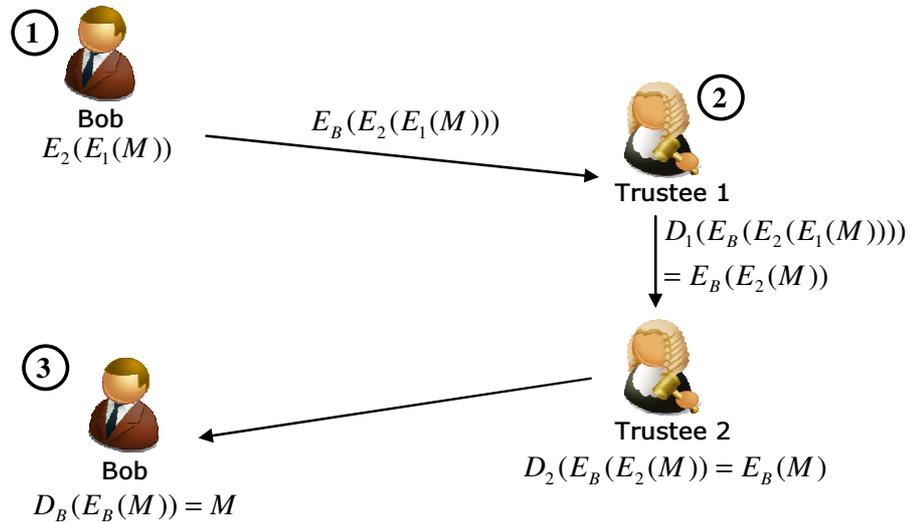
Khayat's secret sharing scheme (unlocking phase)

Objective: Unlocking secret message M locked by a number of trustees.

Bob's input: encryption key E_B , decryption key D_B .

Trustee's input: decryption key D_n : for the n^{th} trustee.

Output: Bob obtains M



1. Bob encrypts the ciphertext with his key and sends it to any trustee in the scheme.
2. A trustee removes his encryption from the ciphertext and passes it to another trustee, or back to Bob. The order of decryption (lock removal) by trustees is arbitrary.
3. Once all the trustees removed their locks, Bob can decrypt the ciphertext and can obtain the plaintext message.

Figure 3-2 Illustration of unlocking a secret under Khayat's secret sharing scheme

In the event of such issues arising, the above secret sharing scheme would be broken. However, these issues were considered by Shamir in a description of his own secret sharing scheme [89], some 25 years before publication of Khayat's scheme. In Shamir's secret sharing scheme, a secret is divided into n pieces, however, only k pieces are needed to use the secret. Thus, up to half of the pieces may be missing (or corrupted) and the operation of protocol would not be affected, and the secret could be retrieved in the most optimal case: $n = 2k - 1$. Consequently, Shamir's scheme (described in Figure 3-3), based on polynomial interpolation, and not an encryption protocol, provides an efficient and an adaptable solution for secret sharing. Also, this scheme is information theoretically secure, unlike the trapdoor-based solution proposed by Khayat.

Many different secret sharing schemes exist, and a good overview of these is provided by Schneier in [41], however, most popular protocols are a variation of the polynomial interpolation concept.

Despite Khayat's scheme not being suitable for the intended purpose, the concept of locking and unlocking a secret in an arbitrary order by a number of parties forms a useful digital equivalent to a *safety lockout hasp*. Such hasps are used by engineers to lock-out an area or a resource, while work is being carried out. Thus, each engineer places a padlock on the hasp on commissioning a task, and removes it once the work is done. The hasp can only be opened once all the padlocks have been removed, meaning that all engineers have finished their allocated tasks. This is an important primitive behind a number of information retrieval schemes described later in this chapter.

Secret sharing scheme based on polynomial interpolation

Objective: Split a secret into n different pieces so that only k pieces are required to read the secret.

Participants: n parties/trustees

Inputs: secret D , a prime number p .

1. Choose a prime number p that is larger than n and D .
2. Pick at random $k - 1$ degree polynomial:
$$q(x) = a_0 + a_1x + \dots + a_{k-1}x^{k-1}$$
where $a_0 = D$ and other coefficients belong to $[0, p)$.
3. Compute n different 2-tuples $(1; D_1), \dots, (n; D_n)$ so that
$$D_i = q(i) \bmod p.$$
4. Distribute the 2-tuples among the participating parties.

Figure 3-3 Shamir's simple secret sharing scheme

3.3.4 Retrieving records in a private manner

Data retrieval is a fundamental operation in computing. Therefore, there is little wonder that PIR is one of the most researched privacy-preserving primitives. Initially, PIR protocols were designed with a basic requirement of acquiring an

interesting data record, or just a specific data bit, from a dataholder, the *sender*, in a way that this dataholder is unable to judge which record is of interest to the requestor, the *chooser*. These protocols were not concerned with the secrecy of other records stored in the database, thus in its least optimised state, a PIR could have been achieved by transferring the whole database from the *sender* to the *chooser*, as this would allow the *chooser* to retrieve a record in a private manner. Consequently, the main motivation behind the research in this field is to achieve PIR with a minimal communicational and computational complexity [90, 91].

There is a firm distinction between a single- and a multi-database PIR protocols. It is possible to achieve PIR with information-theoretic privacy by making a number of requests to the database(s) with a distribution that does not allow the dataholder to identify object of the interest. As expected, this operation is more efficient in the multi-database scenario, and requires minimum of $O(n^{1/2k-1})$ communication complexity, where n is the number of records in the set, and k is the number of databases where this set can be obtained from. However, PIR protocols with only polynomial time assurance can achieve much smaller communication complexity, by introducing balancers [90]. These employ trapdoor functions that can be used to change the ratio of computation to communication, and thus can be used to minimise the amount of data that needs to be transmitted for a given run of a PIR protocol. Such a solution was also suggested by Naor and Pinkas [92], who emphasise that selection of the trade-off between computational and communicational complexity depends on the specific problem at hand.

In [93], Shundong discusses a retrieval system that uses symmetric cryptography in order to lower the cost of cryptographic operations, as trapdoor functions are about 1,000 times less efficient than symmetric operations. However, this solution still requires the use of some trapdoor operations. Consequently, the proposed solution is analogical to PKI, where data is encrypted using symmetric algorithms, and the symmetric keys are then hidden using a trapdoor function, such as these provided by RSA encryption.

A stronger notion than PIR is the Oblivious Transfer (OT) primitive introduced by Rabin [94]. In its original form it allows two parties to engage in a secret sharing protocol that ensures that both parties provide their secret entries without cheating.

Normally, if two parties want to exchange their secrets, one of the parties could provide false data in return for the secret from the other party, or it could back away from the protocol once the secret of the other party has been revealed, but before revealing its own secret. Rabin's OT ensures that the parties learn nothing if one backs out from the protocol before it completes, and that, in case one of the parties cheats in the protocol, this can be proven at the later date. However, this primitive is the mostly widely used in a more sophisticated form that can enable the *chooser* to select one out of two values, or records, held by the *sender* in a way that the *sender* cannot learn which record has been retrieved, and the *chooser* cannot learn anything about the other record. This extension to the OT primitive is referred to as 1-out-of-2 OT (1-2 OT) or OT_2^1 , and in [41] Schneier provides its basic protocol to achieve an 1-2 OT (Figure 3-4).

An Oblivious Transfer protocol

Objective: Allow a remote party to retrieve only one record from a set without disclosing the identity of the collected record.

Participants: Chooser and Sender

Chooser's Input: key with a private-key encryption protocol

Sender's Input: two sets of *public/private keys* pairs

Output: The *chooser* is able to decrypt only one record, while the sender learns that a record has been retrieved

1. The *sender* generates two sets of *public/private keys* pairs, and sends all *public keys* to the *chooser*.
2. The *chooser* generates a key with a private-key encryption protocol, such as AES, later called the *AES key*. The *chooser* then uses a *public key* received from the *sender* in Step 1 to encrypt the *AES key* and send it to the *sender*.
3. The *sender* does not know which *public key* has been used to encode the *AES key*, or which record has been selected, thus protecting the privacy of the *suspect*.

- The *sender* can then attempt to decode cipher-text
4. received in Step 2, using all *private keys* generated in Step 1, whilst preserving the order in which they have been decrypted. In this way two potential AES keys are created. But only one is the proper *AES key*; the other output is a random set of bits, which cannot be distinguished from an ordinary AES keys.
 5. The *sender* encrypts the two records using appropriate keys decrypted in Step 3. Thus, the first record is encrypted with an *AES key* decrypted using the first private key generated in Step 1, the second record is encrypted with an *AES key* decrypted using the second private key. Consequently only one record includes data about the *suspect*. The record is encrypted using the *AES key* generated by the *chooser* in Step 2, sent to the *sender* encrypted by the appropriate *public key*, and then decrypted using relevant private key. In this way the selected record will be encrypted using the proper *AES key*, while the other record will be encrypted using by the random string of bits unknown to the *chooser*.
 6. The *chooser* gets the encrypted records, but using the *AES key* the *chooser* is able to decrypt only the selected record. The other record is unreadable to the *chooser* provided that the false keys generated in Step 3, and used to encrypt these records in Step 4, are not broken.

Figure 3-4 Basic Oblivious Transfer Protocol by Schneier [41]

In the work describing MPC, Yao provided a technique for scaling up any 1-2 OT protocol into a 1- n OT protocol. However, the 1- n OT primitive that allows for the retrieval of a randomly selected record from the dataset of n elements held by the *sender*, may not be useful apart when playing mental games [58, 87], as Schneier point out in [41]. 1-2 OT is typically based on modular exponentiation, thus, it is resource intensive. Consequently, even though it is possible to derive 1- n OT from 1-2 OT, in practice 1- n OT, designed as such from the ground up are more efficient [92, 95]. This is backed-up by Goldwasser's proofs that MPC protocols designed for a specific tasks perform better than the general-purpose protocols [88].

Just like in PIR, the fundamental primitive is designed to operate on bits, while, for most proposed uses, OT of strings is more practical. Protocols that can allow for efficient OT of strings raises the possibility of transferring control over any digitally controlled, or contained, asset from one party to another, as access keys and passwords can be retrieved [60, 95]. For example, if the *sender* wants to allow the *chooser* to privately purchase an electronic book, it can openly publish the full content of the electronic library with each book encrypted under a unique symmetric encryption key, and then once the *chooser* makes a payment, the parties engage in an $1-n$ OT allowing the *chooser* privately select the decryption key for a given book [73, 80]. This is a common approach for transferring control over a resource from one party to another, however, in the digital bookshop scenario; the *sender* would need to know which book has been purchased in order to charge the *chooser* the correct amount. Alternatively all books could have equal prices.

Clearly, both approaches are impractical. A solution to this problem is published in [73], where it is suggested that a buyer should make an initial deposit allowing them to obtain a number of goods. The seller would then need to ensure that the balance of the deposit is higher than the value of the goods being purchased. However, the seller should not learn the exact balance of the deposit, but only the result of comparison between the value of the goods being purchased and the deposit balance (a general protocol for making such private comparison is later discussed in Section 3.3.5). Consequently, the balance of the deposit is encrypted by the buyer and then stored by the seller, when a buyer makes a transaction the value of the transaction is sent to the seller encrypted under the same homomorphic encryption scheme allowing the seller to deduct the value from the balance.

The OT protocols that allow the *chooser* to actively select a record to be retrieved, and that have linear or sub-linear complexity, can also be referred to as SPIR protocols, as the primitive protects the records of both parties during the information retrieval process. In addition to the already discussed uses of the OT, and SPIR, primitives can be employed in a variety of systems: electronic watch-lists of suspects [74]; cooperative scientific computation [96, 97]; and on-line auctions [98].

Research presented in [99] implies it is unlikely to achieve an OT without a *trapdoor function*, which is a public key operation. For this reason most OT protocols are

based on the asymmetric encryption employing exponentiation modulo of a prime number. Consequently, most OTs can benefit from a technique for fast exponentiation employed in [95] and defined in Brickell-Gordon-McCurley patent [100]. This technique can greatly improve the performance of most protocols based on modular exponentiation. However, researchers have often chose not to discuss this in detail in the relevant publications discussing the use of PIR and OT primitives, as this is purely a technicality, and it tends to obfuscate the cryptographic solutions being presented.

3.3.5 Private Value Comparison – Locating interesting records

The primitives defined in the previous sections provide techniques to retrieve a record from the *sender*, without the *chooser* revealing anything about the record of interest. However, these primitives require that the record is retrieved using an index. Such approach can be justified for protocols designed for information retrieval from online stores, or databases, where directories providing basic information about each and every record are publicly available. However, in the scenario with no publicly available index of the interesting record, such an approach would fail. Thus, there is a need to provide a method for matching (or comparing) the description of the interesting record with the description record held by the sender, so that an index of the interesting record can be identified. This functionality can be provided by the schemes described next.

An efficient technique for value comparison has been described in [101] where it is used in the context of private bidding. It is also suggested that a protocol that allows for comparing two values privately, where the values are the maximum price a bidder is willing to pay for an item and the minimum price the seller is willing to sell for, can allow for on-line haggling, or *bargaining*, in order to determine a price of an item. A semi-trusted third party is introduced by [101] in order to minimise the communications and computation required by the protocol. This third party is oblivious to the results of the protocol, and it is only trusted not to collude with any of the participants. Thus, an auction house would be a suitable third party for an implementation of the protocol. The protocol compares values bit-by-bit using PIR circuits based on the difficulty of factoring (as per RSA [44]) and higher-residuosity assumption, as discussed in [102].

Privacy-preserving approaches to compare information are essentially different approaches to solve Yao's millionaires' problem. While the millionaires' problem is a good example for an academic discussion, in practice the comparison circuit can be used to facilitate Internet second-price auctions [103], and any other operations where the value comparison must be run on secret inputs. A number of interesting and unconventional approaches to performing data comparison are provided in [104]. The scenario published in [104] requires comparing two secret entries, which in this case is the name(s) of an individual or individuals that made complaints to two managers participating in the protocol, in order to check whether the complaints were made by the same individual. This calls for a special case of value comparison, which is Private Equality Test (PEqT), where in [74] it is referred to as an asymmetric equality test. PEqT is the key primitive in the area of private matching protocols. This primitive allows two parties to compare their secret inputs for equality without revealing these inputs. There are two cryptographic concepts that the PEqT can be based on: commutative cryptography [47, 58, 104]; or homomorphic cryptography [80, 98]. Each of these techniques has its benefits depending on the problem at hand.

The first published solution to the private matching problem is the commutative cryptography scheme used in the protocol for playing Mental Poker over a distance [58] first drafted in 1979, and further analysed by Shamir in [47]. For two parties, Alice and Bob, each holding different commutative cryptography keys the operation of the protocol is summarised in Figure 3-5.

This scheme employs a modification of the Pohlig-Hellman (PH) algorithm described further in Chapter 4. Thus, each encryption/decryption operation requires only a single exponentiation. To date, a number of different PEqT schemes have been proposed, but the complexity of the other schemes is usually higher than this of the commutative encryption solution presented above. Boa and Deng [80] described an efficient method for equality testing based on homomorphic encryption. However, this method requires a series of multiplications, an exponentiation, as well as a round of homomorphic encryption and decryption. The homomorphic encryption used in their scheme is ElGamal, which itself requires two exponentiations modulo a prime during the encryption process, and another for the decryption operation.

Consequently, the computational complexity of their protocol, as well as the protocol described in [98], is higher than this of the PEqT scheme illustrated in Figure 3-5. However, only the slight difference in performance means that the decision of using one or the other method should be based on factors other than efficiency alone.

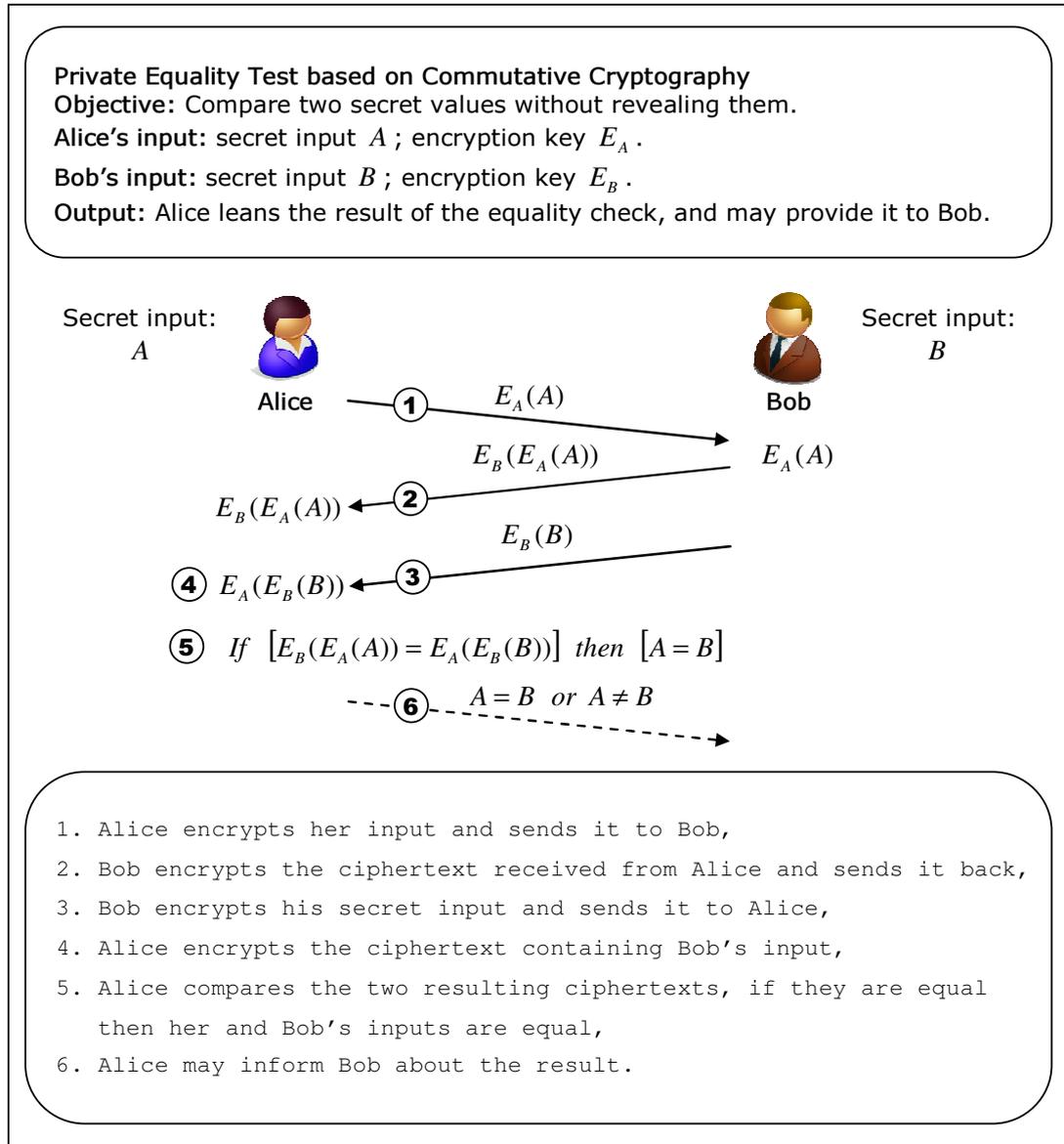


Figure 3-5 Private Equality Test based on Commutative Cryptography

When two parties engage in equality test protocols, often there are a number of inputs to be compared. Thus, a scenario exists where the *chooser*, Alice, wants to compare her value with a number, n , of values held by the *sender*, Bob. In such a scenario, if the homomorphic scheme of Boa and Deng [80] was to be used, then for each record held by Bob, four exponentiations would be required. However, the commutative cryptography-based PEqT shown above can be modified so that the $1-n$

PEqT protocol could be achieved with $O(n + 3)$ exponentiations. The modification of the protocol is in Step 3, where Alice computes a value equal to her input encrypted only by Bob. In this way she can now compare the resulting value with Bob's inputs encrypted by only him. The resulting protocol allows Alice to compare her secret input with n inputs held by Bob.

As far as the openly available literature goes this protocol is likely to be the most efficient *1-to-n* PEqT (or *1-n* PEqT) protocol available. Interesting extensions to the concept of *1-n* PEqT are private intersection and private intersection size protocols. Thus, this protocol can be extended into the intersection size protocol described [60]. A similar approach to computing secure intersection size is also provided in [105]. Whereas, Freedman, Nissim, and Pinkas presented an efficient secure intersection protocol in [106] that is improved in [107], Weis argues that these protocols share a fundamental security flaw, as for a malicious party it is trivial to convince the other party that an intersection exists [55].

3.3.6 Combined approaches to selective information retrieval

The Private Equi-Join (PE) protocol can enable two parties, the *chooser* and the *sender*, to privately compare their sets of unique values V_C and V_S , and allows the chooser to retrieve some extra information $ext(v)$ about records V_S , that match records V_C on a given parameter [60]. The PE protocol involves the steps described in Figure 3-6 and represented graphically in Figure 3-8.

Researchers have shown that using a TC device such as PCI-attached IBM's 4758 SCOP, it is possible to perform efficient hardware-based PIR that allows for the selection of a record based on given match criteria [108]. In such a scenario SCOP can easily match any record based on selection criteria, however, the problem is still in retrieving the record in a way that the host computer cannot identify the record that is sent back to the *chooser*. In an ideal scenario the SCOP would collect a number of records stored on, or accessed through, the host machine so that is impossible to identify which record is being sent to the *chooser*. However, the difficulty lies in the fact that SCOPs often do not have enough memory to store and process many records, as they have a limited amount of RAM. The solutions presented in [108] involve SCOP performing the steps detailed in Figure 3-7.

Private Equi-join protocol

Objective: Chooser obtains data elements linked to identities from its request, and nothing else. Sender learns nothing from this transaction.

Chooser's Inputs: a set of identities V_C , a PH key pair $\langle E_C, D_C \rangle$, and a hashing protocol h . Symmetric encryption and decryption functions K and K^{-1} respectively.

Sender's Inputs: a set of identities V_S with corresponding data elements referred to as $ext(v)$. PH encryption keys E_S, E'_S , and a hashing protocol h . Symmetric encryption and decryption functions K and K^{-1} respectively.

Output: The above objective is met, but the parties also learn the sizes of each-other's sets.

1. Both parties apply hash function h to the elements in their sets, so that $X_C = h(V_C)$ and $X_S = h(V_S)$. Chooser picks a secret PH key E_C at random, and sender picks two PH keys E_S and E'_S , all from the same group Z_p^* .
2. Chooser encrypts entries in the set: $Y_C = E_C(X_C) = E_C(h(V_C))$.
3. Chooser sends to sender set Y_C , reordered lexicographically.
4. Sender encrypts each entry $y \in Y_C$, received from the chooser, with both E_S and E'_S and for each returns 3-tuple $\langle y, E_S(y), E'_S(y) \rangle$.
5. For each $h(v) \in X_S$, sender does the following:
 - (a) Encrypts $h(v)$ with E_S for use in equality test.
 - (b) Encrypts $h(v)$ with E'_S for use as a key to lock the extra information about v , $\kappa(v) = E'_S(h(v))$.
 - (c) Encrypts the extra information $ext(v)$:
$$c(v) = K_{\kappa(v)}(ext(v))$$

Where K is a symmetric encryption function and $\kappa(v)$ is

the key crafted in Stage 5b.

(d) Forms a pair $\langle E_s(h(v)), c(v) \rangle$. These pairs, containing a private match element and the encrypted extra information about record v , are then transferred to *chooser*.

6. *Chooser* removes the encryption E_c from all entries in the 3-tuples received in Step 4 obtaining tuples α , β , and γ such that $\langle \alpha, \beta, \gamma \rangle = \langle h(v), E_s(h(v)), E'_s(h(v)) \rangle$. Thus, α is the hashed value $v \in V_c$, β is the hashed value v encrypted using E_s and γ is the hashed value v encrypted using E'_s .
7. *Chooser* sets aside all pairs received in Step 5, whose first entry is equal to one of the β tuples obtained in Step 6. Then using the γ tuples as symmetric keys it decrypts the extra information contained in the second entry in the pair $\langle E_s(h(v)), c(v) \rangle$.

Figure 3-6 Operation of the Private Equi-join protocol

Hardware PIR:

1. Retrieve records one-by-one.
2. Compare each record to the match criteria.
3. Encrypt each record and store in the host's memory system, keeping a note of the memory location belonging to the record matching the selection criteria.
4. Once all the records have been retrieved by SCOP and stored encrypted in the host's machine, shuffle and re-encrypt all records. In this way the host machine can no longer link records retrieved in Step 1 to their encrypted form.
5. Pick the record matching the selection criteria and send it securely to the requestor.

Figure 3-7 Efficient PIR based on Secure Coprocessor

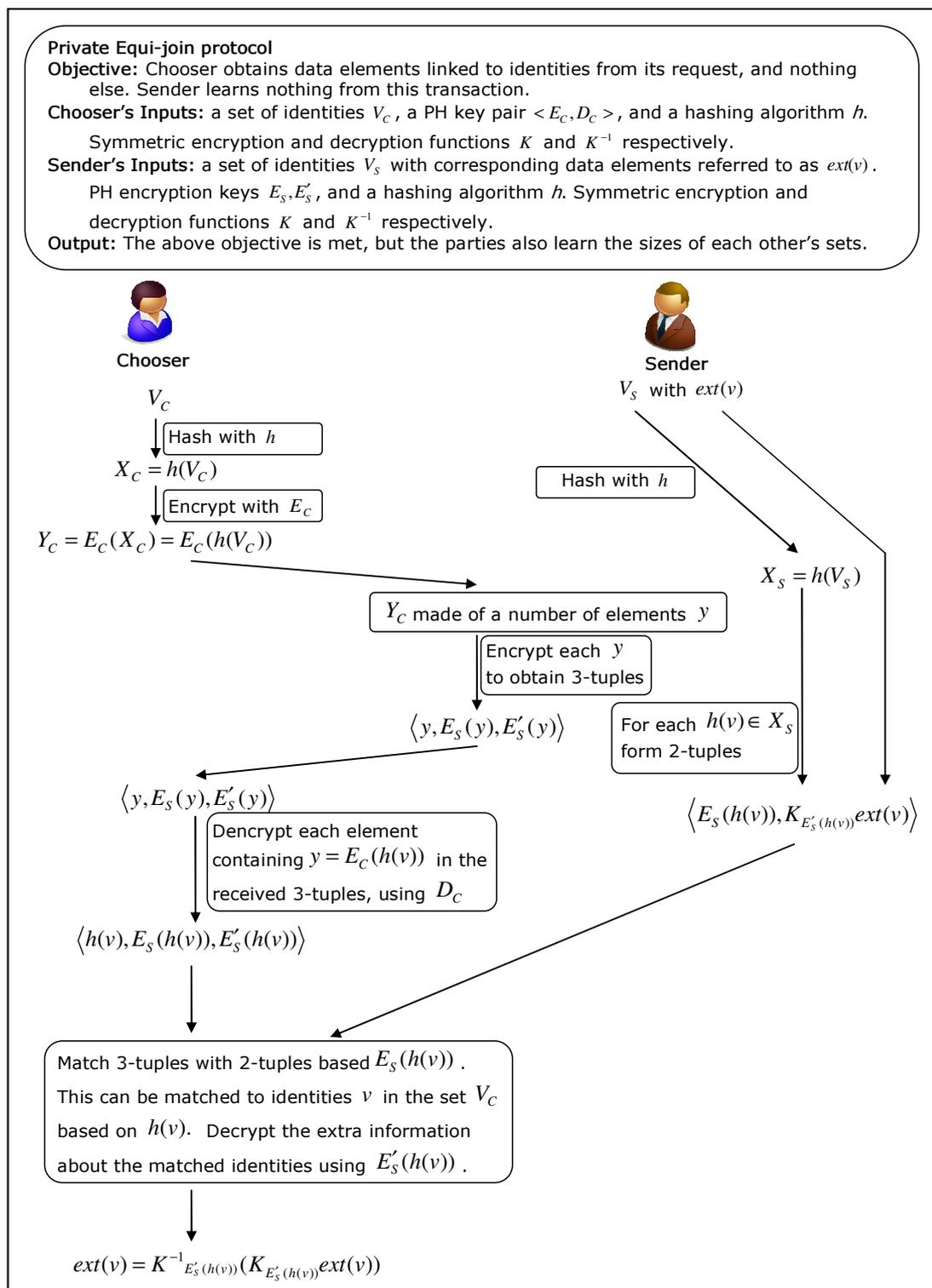


Figure 3-8 Graphical representation of the Private Equi-join protocol.

Introduction of the hardware into the PIR has not greatly lowered the complexity of this primitive. SCOP needs to encrypt each record at least twice (Step 3 and Step 4), and the operations related with loading and unloading the data from the SCOP also add a delay into the protocol. By introduction of the *square-root* algorithm i different PIR requests can be allowed to run following a single shuffle of the records. This

improves greatly the performance of the protocol, but, in order for the *sender* not to realise whether a given record has already been requested or not, the SCOP needs to pick from the host's storage every record that already has been picked in a given shuffle round, as well as one new record. This should be either the record selected by the *chooser*, or a randomly selected one if the *chooser* requested a record that has already been picked in this round. A new shuffle round is run every $i = \sqrt{n}$ records [108]. It is worth noting that this is a PIR and not SPIR protocol, since there is no way for the *sender* to judge how many records have been retrieved.

3.3.7 Security Considerations

Security measures should always be considered in relation to realistic threats to a given system. Thus, Goldwasser discusses four adversaries specific to MPC protocols [88], these are:

- **Passive.** One or more participating party aiming to obtain the secret input of the other participants.
- **Byzantine.** Otherwise referred to as malicious adversary. A party that does not follow the protocol, and provides other parties with specially crafted inputs in order to obtain secrets or compromise other participating parties in another way.
- **Mobile.** A coalition of the passive and byzantine adversaries formed by different set of parties at each round of the MPC protocol.
- **Coercing.** This can force users to provide specific inputs, for example vote in a specific way during electronic elections.

The security of a cryptosystem often depends on correct design and implementation, just as much as on the strength of underlying cryptographic protocol [109]. Most attacks on common systems are possible only due to erroneous design, implementation or maintenance, and not weaknesses of the underlying cryptographic algorithms [110]. Consequently, just because a given cryptographic protocol used strong underlying algorithm, it does not mean that the protocol is secure as shown by the example of a watermarking protocol that fails to sufficiently link the watermarks to the digital goods being signed described in [48]. Because of the nature of cryptanalysis, where any cryptosystem or ciphertext may be attacked in an arbitrary

way by a previously unknown attacker, the cryptography is being compared to *programming Satan's computer* [111]. The key lesson to learn from this approach to cryptography is that it is likely that errors in implementations of certain protocols will occur, and there needs to be a mechanism that will allow for these errors to be fixed (this is also highlighted in [112]).

When RSA was first published its authors encouraged the readers to attempt breaking their algorithm and protocol, as they wanted to make sure that they had not overlooked any potential flaws. RSA was the first to use *trap-door one-way permutation* in a cryptographic algorithm, and thus the exact strength of the protocol was not known [44]. Nowadays, over 30 years later, RSA is still considered as secure if certain conditions are met. RSA has been an inspiration for a number of privacy-preserving solutions. Over the years its security has been addressed by a number of academic and industry studies. These have been summarised in [113], where the *RSA Problem* has been formally defined as the problem of obtaining the plaintext message from the ciphertext and the public key used to produce this ciphertext. It is shown that the RSA Problem is no harder than integer factoring, however, taking into consideration that RSA modulus n is sufficiently high, then RSA Problem is hard to solve. However, the randomness of the plaintext over the range $[0, p-1]$ is also crucial. Some studies of the RSA suggest that using strong primes in the algorithms is necessary in order to safeguard the systems from factoring and *cycling* attacks. However, Rivest and Silverman prove that using strong primes in RSA yields limited benefits to the strength of the protocol, as long as the primes used are reasonably large [114]. This is unlike the DH and the ElGamal, two protocols commonly employed in MPC that need to be based on strong primes [110]. In [115], Sakurai and Shizuya discuss the security of various protocols (including DH, 3Pass, and ElGamal) based on the Discrete Logarithm Problem (DLP). Their research suggests that ElGamal and DH can both be reduced to 3Pass in a polynomial time, and all three protocols should be considered as equally strong. No efficient attacks against DLP are currently known, however, this does not necessarily mean that the schemes based on DLP cannot be broken without breaking the DLP [115].

A number of primitives described in this chapter assume an existence of a secure commutative encryption scheme. However, traditional means of testing the security

of encryption schemes are not capable of evaluating schemes where ciphertexts commute [55, 116]. In [116], a new model for assessing such protocols is presented, and the security of commutative protocols based on RSA is shown to be in NP. Additionally, [55] presents a technique that can transform any semantically secure homomorphic encryption scheme into semantically secure commutative scheme, thus, allowing already existing crypto libraries to perform commutative operations. Finally, in Section 2.5.4, the Massey-Omura algorithm was suggested as one of the possible algorithms that could be employed in systems relying on commutative cryptography. This cryptosystem performs operations in $GF(2^n)$ in order to allow hardware-accelerated implementation of the Pohlig-Hellman cryptosystem. However, [117] shows that discrete logarithms are easier to compute in this field than in $GF(p)$, and therefore using this field for cryptographic operations should be carefully considered.

It is possible to cheat in some protocols and provide the other party with crafted input data that has not been created using an encryption key on the originator, but prepared in order to reveal the secret of the second participant. Such a scenario can be mitigated by the use of ZKP on the inputs from another party, just to prove that the inputs have been generated according to the protocol [73]. In academic discussions and few specific real-life scenarios, it is possible to ignore the threat from the possible exploit by assuming that the participants are honest, but curious. Often PET protocols are presented in *honest-but-curious* form in order to simplify the analysis of the protocols [73, 74, 91, 103]. This assumes that the participating parties follow the protocol (honestly) but will try to compute and imply any information they can with any data obtained during the process (curiously). These protocols can then be transformed into malicious mode with use of ZKPs.

Another security problem in complex systems is the fact that a number of privacy-preserving primitives may need to be used in order to perform a given task. Such composite system would likely reveal more information than required, since apart from the final output the intermediate results would be revealed. As [105] suggests, it is possible to define the intermediate results as a part of the output, in order to evaluate a system under the rules of MPC. This is not an ideal solution, thus, such

composite systems should be avoided if possible. However, this solution ensures controlled disclosure and, in most cases, this is sufficient.

3.4 Conclusion

In theory the DPA provides individuals using UK-based services with a full control of their personal data. An organisation wishing to collect data about an individual must obtain consent (this can be *implicit*) from this individual. This consent can be withdrawn by the individual/data subject at any time, by an opt-out procedure that all organisations storing personal data must provide. However, even if a data subject perceives a given system as intrusive the convenience and economical factors can force this data subject to keep using the system. In such scenarios individuals often decide to use anonymising technologies on the Internet, and tend to use cash in face-to-face transactions. The anonymising technologies, including TOR, are created by a large number of users creating a virtual network (over the Internet) that can hide the identity of an Internet user among 100,000 other TOR users. It uses onion routing based on sequential re-encryption of network packets in order to stop the ISPs and other TOR users from tracing the network packets back to the user that generated them. While TOR is implemented, and is commonly known to Internet users interested in maintaining their privacy, there are other systems proposed that can allow for anonymised purchases, and on-line auctions. These usually rely on concepts from the area of MPC.

MPC protocols can also be used to facilitate privacy-preserving investigations. Literature shows that it is possible to build systems that allow investigators to trace data subjects marked as suspects, without revealing the identity of the suspect or affecting the privacy of other data subjects in a given system. Also, it is possible to create pseudonym-based auditing systems that only reveal the identity of an individual if the actions performed by this individual have reached a threshold of malicious activity.

The solutions employing MPC are often impractical. In early MPC protocols the computational complexity was exponential to the number of bits used to store the private records and to the size of the data records. Fortunately, it is possible to manipulate computational and communicational complexity of different schemes by

using different cryptographic techniques and introducing semi-trusted third parties that could proxy the requests. In similar fashion to PKI, MPC often employs symmetric encryption to lock data that is transferred between the participating parties, while computationally-expensive trapdoor functions, such as public-key cryptography, are used only to conditionally exchange the symmetric encryption keys. It is worth noting that it is unusual for MPC-related research to provide empirical evaluation of protocols. Most research into MPC has focussed on perfecting previously developed schemes, with little attention paid to their practical use [91-93, 95]. A comparison of the different schemes is usually done on the basis of computational and communicational complexity, which, some researchers assert should not be directly compared. In general the efficiency of encryption schemes based on modular exponentiation (used by most trapdoor functions) is approx. 0.1% of the symmetric encryption protocols. Thus, a protocol that takes $O(1000 \times n)$ symmetric operations, would take a similar amount of time to a protocol with $O(n)$ trapdoor function operations, while the computational complexity expressed in terms of the number of operations would suggest otherwise.

While there is a number of PIR and OT primitives that allow private retrieval of records, most require the interesting record to be identified by its index in a given dataset. This approach is optimal in scenarios with an index or a catalogue of the database being available publicly. Such scenarios include purchasing goods or services from an on-line retailer or a service provider, which is the main motivation for a number of PETs. However, in order to retrieve records matching certain selection criteria, it is necessary to run equality tests on the data. It is possible to combine a PEqT primitive with PIR or OT in order to achieve such functionality. But it is also suggested that complex MPC protocols made up of a number of privacy-preserving primitives can release more information than required by a given scenario and the complexity of such protocols is usually less optimal than custom-made protocols. However, the PE protocol based on commutative cryptographic algorithm is a purpose built system for retrieving records that match given selection criteria. Its authors suggest that it should be suitable for use in sharing data between hospitals and other organisations with large databases. It is also possible to achieve similar system if TC hardware device SCOP is deployed to the database server (or a host

attached to the database), but it is much harder to provide guarantees of privacy of other data records stored by the dataholder if SCOP is deployed.

Chapter 4

Improving the Acquisition Process

4.1 Introduction

This chapter documents the initial evaluation of the PET protocols in an investigative context. For this purpose a set of requirements for data acquisition process is drafted and refined. These are based on the available literature such as regulations, guidelines and procedures. Often these requirements are inferred rather than obtained from these sources, and thus expert opinion was obtained as to their validity (see Section 4.7.3 *Feedback from practitioner*).

Gathered requirements are used to select candidate PET primitives. Two different approaches to building a suitable solution are proposed based on related research. These include a solution built from 1- n PEqT and 1- n OT, as well as a solution based solely on the PE primitive. These approaches are empirically evaluated based on tables of the computational complexity produced for each solution and experimentally established timings for applicable cryptographic operations.

Advantages and disadvantages of both solutions are identified based on the requirements.

The outcome of this chapter is a set of requirements for data acquisition process and the suggestion of a protocol capable to satisfy most of these requirements. This contributes towards the design of the data acquisition platform proposed in this thesis and presented in Chapter 5.

4.2 Methodology

Chapter 3 concluded that it is possible to protect the interests of two parties wanting to compute a function without revealing the secret inputs or to conditionally exchange some data between the parties, with the use of PETs. Arguably, it should be feasible to construct a process that employs PETs to retrieve investigative data in a privacy-preserving manner. With the range of PIR, OT and SPIR primitives available, it is likely that a single protocol can perform the required operation, and, if not, then a combination of existing primitives should be able to achieve this. However, before such primitives can be identified, a set of requirements for the data acquisition system needs to be drafted. This can be done based on the literature discussing the data acquisition process, UK legislation and digital forensics research. The gathered requirements can then be used to define evaluation metrics for the platform and to analyse the available PET primitives for the suitability of use in an investigative scenario. Thus, the identified primitives can be evaluated more thoroughly, and compared against each other. Summarising the methodology of the work presented in this chapter is as follows:

- Define the requirements for a data acquisition platform.
- Identify the evaluation criteria.
- Select the types of protocols that can improve privacy in data acquisition process
- Evaluate known protocols in order to select the most suitable.

It should be noted that at this stage only the existing primitives are taken into consideration.

4.3 Initial Requirements

The requirements for a data acquisition platform that can be derived from Chapter 2 and Chapter 3 can form a guideline for design, implementation and evaluation of a data acquisition platform. Following is the list of the requirements derived in this way together with their explanation:

Req. 1 Allow for the gathering of multiple suspect records per enquiry, or have low overhead per each additional query run on the database. (Maximum anticipated number of suspects in one enquiry is 150.)

Description: There is a suggestion that the Internet is now used in organised crime [2] which can mean that some inquiries will require the retrieving of data about a group of suspects, rather than a single individual. Consequently, the protocol chosen for the data acquisition process will need to allow for the retrieval of a number of interesting records at the time, or, if this is not the case, multiple sequential runs of the protocol should bear low overhead. Even that this requirement can be represented as a number it is hard to quantify it. Statistics relating to the average number of the suspect could be gathered from the public authorities using the Freedom of Information Act 2000 [118], but resulting number could be misleading and would not include anti-terrorist enquiries, etc. Possibly a limit of the potential suspects can be set to the Dunbar's number [119], which describes the maximum number of meaningful relations a human can have with others, which would most likely limit the size of criminal networks as well. The Dunbar's number is not strictly defined, but a value of 150 is one of the commonly suggested values and this will be used as the maximum for this requirement.

Req. 2 Keep the data controller in charge of the data. A data record cannot be transferred or made available, to the public authorities, without the data controller's verification of the request.

Description: Investigators need to provide justification for the acquisition requests under the DPA and the dataholder can refuse to provide any data without a valid warrant from the court of law [10]. Whereas, data acquisition notices served under RIPA do not need any form of justification to the dataholder and the dataholder will face a penalty if the relevant data is not provided to the requesting public authority within two weeks. Still, the dataholder may choose to accept the penalty and refuse to provide any data without a subpoena [8, 12]. Consequently,

the platform must leave the dataholder in control of the data, since the data retrieval can only be performed with the dataholder's consent.

Req. 3 Allow for efficient and timely data retrieval. The current maximum time for returning data under RIPA is two weeks, however, it is expected that in urgent enquiries investigators would have access to data in a reasonably short time.

Description: Taking into consideration that a dataholder has two weeks to provide the data under RIPA the computational complexity of the protocol can be reasonably large [12]. However, shortening the time required by the data acquisition process is one of the main reasons provided in [2] as a justification for the proposed modernisation of the data acquisition process, and therefore this number should be revised.

Req. 4 Be cost-effective, as the platform will need to be deployed by a variety of organisations.

Description: Under RIPA the public authorities must make a contribution towards the costs incurred by a CSP during the fulfilling of the data acquisition notice [12]. Thus, the cost of the solution should be low. If the costs were not covered by the authorities, the dataholders would transfer the costs of handling the enquiries to the end-users and such a solution would typically be unacceptable by society.

Req. 5 Retain an audit trail of the processing performed on the potential evidence.

Description: All processes applied to computer-based electronic evidence should be preserved in an audit log so that an independent third party could examine these processes and achieve the same result [17]. Also, any evidence collected may need to be presented in front of court of law, which will require that the electronic evidence must be provided as a true image of the data gathered. So data records should be retrieved from the dataholder on a record-by-record basis, so that if only one of many records is required for the investigation, other records can be discarded. Otherwise the public authorities can end up storing a large amount of unnecessary data, and this can prove costly, taking into consideration the level of security and auditing involved in storing digital evidence [17].

Req. 6 Gain acceptance from the general public.

Description: This thesis is written in response to the worries of the general public about their privacy. Therefore, one of the requirements must be to make the system appeal to the public.

Req. 7 Handle large datasets (such as datasets with more than 15 million records).

Description: To put this in the context, BT has 15 million broadband users [120],

and thus, the system developed in this thesis must handle datasets of similar size.

4.4 Overall design

According to [8] two fundamental parts of data acquisition process are: serving the notice to the dataholder; and the subsequent retrieval of data. The retrieval is often achieved by the dataholder sending back the data to the requestor. The request and the response both need to be performed under the guidelines of the DPA. Thus, some form of secure channel needs to be established between the parties, or the messages need to be encrypted while in transit, with a technique that is either FIPS140 compliant or at least FIPS198 compliant. Currently, such a data acquisition notice would include the specification of the requested records that the dataholder would then use to build a database query. Almost all relational databases support Structured Query Language (SQL) queries [121] and most likely any data acquisition notice would be translated into such a query. Thus, the notice and the SQL query should contain the following parameters:

- 1) Identification of the type of the information that is required. These could be number parameters that contain answers to investigator's questions. (Represented in SQL representation introduced below as H different return parameters, $rp_1 - rp_H$.) For example in an enquiry for the recent use of a given credit card, the return parameters would consist the location where the card was last used, together with the transaction amount and the date of the transactions.
- 2) Specification of any circumstantial request constrains. (illustrated in the following examples as L different input parameters, $ip_1 - ip_L$, with values $ip_val_1 - ip_val_L$.) Using the above scenario this could include the time constrains specifying the time window for interesting transactions, such as "all transaction between 12/08/2010 and 18/09/2010".
- 3) Specification of the relevant data subject, the individual whose data is being retrieved, by providing the ID of the interesting record (such as the mobile phone number of the suspect). (This parameter is later referred to

as the Record of the Interest ri , with value of ri_val). In the above scenario this would be the credit card number.

Then, if we refer to the dataset as the *source*, the request for investigative data could be mapped into the following SQL query:

```
SELECT rp1, rp2, ..., rpn
FROM source
WHERE ri=ri_val AND ip1=ip_val1 AND ... AND ipL = ip_valL
```

Figure 4-1 Typical request for investigative data mapped into SQL

In most cases the names of the return parameters, as well as the names of the input parameters, and values of these input parameters, can be openly communicated. But the value of the interesting record (ri_val) is used to uniquely identify the suspect and therefore in order to provide privacy to the potential suspects, it must be hidden. This can be achieved by running a database query for the return parameters of all the records that satisfy the conditions defined by the input parameters, and then collecting the interesting record from the dataholder using the some privacy-preserving protocol based on the OT primitive. Consequently, the query that is actually run on the dataholder's database can be rewritten as:

```
SELECT ri, rp1, rp2, ..., rpn
FROM source
WHERE ip1=ip_val1 AND ip2=ip_val2 AND ... AND ipL= ip_valL
```

Figure 4-2 Request enabling privacy-preserving queries mapped into SQL

The input parameters ip would need to be selected so that the above query is guaranteed to return a sufficiently large set of results. Chapter 3 concluded that in order to retrieve data in a privacy-preserving manner there needs to be a publicly available index of data records or 1- n PEqT protocol needs to be run in order to obtain such index of (or other form of pointer to) the records of interest. Only then the OT or SPIR primitive can be used to retrieve data from the dataholder. Consequently, there are three distinct operations required in the process of data acquisition. For simplicity in this, and further consideration of the scenarios, a relevant public authority requesting the data is referred to as the *chooser*, while a

dataholder is called the *sender*. Thus, the following are the key operations needed to acquire investigative data in a privacy-preserving manner:

- 1) **Querying.** The *chooser* specifies the type of information that is required for the investigation. This can be achieved using SQL, since it provides standardised format for database querying.
- 2) **Searching.** Allows the *chooser* to find an index of, or a pointer to, the interesting record in the *sender's* database, by the means of private-matching techniques, such as PEqT.
- 3) **Retrieval.** Finally, the interesting record is retrieved from the *sender* using the OT or SPIR primitive.

The above list excludes some elements of the data acquisition process that are derived from RIPA, DPA and the guidelines on data acquisition [8, 10, 12]. Those excluded processes are the steps required to obtain authorisation and the definition of the roles in the data acquisition process (Applicant, Designated Person, SPoC, or Senior Responsible Officer). They have been discussed in Chapter 3, and it has been established that these processes are fit-for-purpose and protect the integrity of the investigation and privacy of the involved parties as much as possible, without the involvement of PET technologies. On the other hand, PET technologies appear to be capable of fitting into this process smoothly, by replacing the current technologies used during for exchanging notices and the data between the SPoCs of the relevant parties.

There are two possible solutions to address searching and retrieval operations. These can be achieved using a combination of PEqT and SPIR, but also a combined approach, such as PE, can be used. There are advantages to both of these approaches. Using a combination of primitives it may be possible to keep more detailed audit logs and provide verification for requests (Req. 5 and Req. 2), as the searching phase is independent and this would potentially allow for running independent checks on the records being requested by the chooser. Also, such a solution could prove to be the least costly (Req. 4) as some of the primitives are built using standard cryptographic protocols that can be found in existing cryptographic libraries, which would cut the development, compliance testing and maintenance costs, and also

make the solution more transparent (contributing towards Req. 6). On the other hand a combination of different privacy preserving primitives can reveal some extra information that is needed to link the two primitives [105], which is not the case if a problem-specific solution is derived straight from cryptographic algorithms. The downfall can be reduced; however, the mitigation can increase the complexity of the protocol, and hence its cost and the time required for queries would also increase. Therefore, since most problem-specific approaches are usually more efficient than a combination of two primitives, it is likely that a protocol such as PE would better fulfil Req. 3.

It is often impossible to compare protocols based on theoretical evaluation criteria such as communicational and computational complexity. These parameters can only be used to compare the efficiency of protocols built in a similar way, thus an improved version of a protocol and an original version of this protocol can be directly compared using these parameters. But protocols built on different concepts can not be directly compared in this way [93]. Consequently, in order to find a suitable protocol for data acquisition purposes, it was necessary to find a good combination of PET primitives that can perform actions similar to those required by the requirements, and also a combined problem-specific primitive that matches most closely the requirements, and then compare these two approaches. Section 4.5 describes the design of the Searching and Retrieval functionality for the data acquisition framework with use of separate primitives for these functions, while Section 4.6 provides details of a design based on a problem-specific primitive. In Section 4.7 these designs are evaluated side-by-side against the requirements.

4.5 Approach 1: Combination of PET primitives

Both Searching and Retrieval phases of the process can be performed by a number of different primitives. These are analysed and suitable candidates for implementation are selected. These candidates are then put together to provide the required functionality as described by the Design and Implementation section. The selected solution has been published in [9].

The Searching phase needs to establish a pointer to the interesting records in the database, or, more precisely, to the *source* table resulting from the Querying phase.

In practice this can be achieved by privately comparing the identity of the interesting record to the records in the *source* for equality. From the protocols discussed in Section 3.3.5 the one that truly stands out as the most likely to be efficient $1-n$ PEqT protocol is derived from $1-2$ PEqT presented by Shamir et. al. in [58] and illustrated in Figure 3-5.

Some homomorphic encryption based protocols, such as the one presented in [80], can compete with the efficiency of the commutative solution in $1-2$ PEqT operations. However, such protocols normally need to be completely rerun for each record being compared, and, thus in $1-n$ PEqT, the complexity would simply be increased by n times. In the case of the protocol derived from [58] this is not the case and while the computational complexity of the $1-2$ PEqT based on commutative cryptography is $O(4)$, the derived $1-n$ PEqT is characterised by computational complexity of $O(n + 3)$. This would satisfy Req. 1 and Req. 3 that deal with efficiency and rapidness of the enquiry. Since, the protocol can be based on a specific case of RSA or Pohlig-Hellman algorithms, the development time necessary would be minimal and the solution should be easier to explain to decision makers that are aware of current encryption standards (which is key to gain acceptance for the protocol). This is the case where the protocols are well known, and thus well researched. Thus, Req. 4 and Req. 6 would most likely be satisfied as well. The *sender* would not have to disclose any information apart of the confirmation whether the interesting record exists in the *source*, and the location (or index) of this record.

The selected $1-n$ PEqT is based on commutative cryptography, and thus a suitable commutative cryptography protocol has to be selected first. Section 2.5.4 discussed different commutative algorithms. From those protocols the one based on ElGamal encryption cannot be employed in the $1-n$ PEqT test, as the plaintexts encrypted under two different keys used in arbitrary order are not equal under ElGamal. Also, the literature suggests that Pohlig-Hellman is a better choice than Massey-Omura cryptosystem, since the discrete logarithm problem is harder in the $GF(p)$ field than it is in $GF(2^n)$ [117] (as discussed in Section 3.3.7). Thus, the protocol selected for the implementation needs to be either the Shamir's commutative protocol (based on Pohlig-Hellman algorithm) or a modification of the RSA scheme, sometimes referred to as Shamir-Rivest-Adleman (SRA) protocol [116]. Since, the OT selected for use

in the Retrieval can be based on RSA, SRA is chosen. This means that the solution could be fully implemented using common cryptographic suites such as Legion of the Bouncy Castle cryptography library [122], with small changes to the way that RSA keys are exchanged, as in SRA both the encryption and decryption keys need to be kept private, and only the modulus, and the primes used to generate it, are shared between the parties.

Most OT and SPIR protocols can be used to perform the Retrieval phase and the difference is mainly in performance. One of the protocols that stands out for its use of common encryption protocols in the design is the OT discussed by Schneier [41] and presented in Figure 3-4. Schneier's example illustrated 1-2 OT protocol; however, it can be extended to perform 1- n OT functionality. It can be built using a combination of virtually any type of public-key and private-key encryption algorithm. Thus, it would be possible to obtain FIPS-140 accreditation for the Retrieval part of the data acquisition process. Also, this OT protocol is relatively easy to comprehend by the professional audience, including the decision-makers, with basic understanding of PKI. Therefore, such a solution could be presented to the relevant decision makers regulating the investigative data acquisition field.

In fact, the operation of most OT protocols requires that the *sender* provide the *chooser* with an encrypted copy of all records in the table resulting from the Querying phase. Then, the control over the interesting record (its decryption key) is retrieved using an OT protocol. The operation of the chosen 1- n OT protocol is described in Figure 4-3. The choice of the asymmetric public key encryption is limited by the fact that the encryption protocol cannot issue errors if the wrong key is used to decrypt a ciphertext. Thus, ElGamal cannot be used, but RSA is a good choice, especially that it ties-in with the SRA used in the Searching phase. On the other hand, this OT protocol can be based on virtually any symmetric encryption primitive, despite the fact that similar restrictions (key verification) are put upon the symmetric operations in the process. This is due to the fact that key verification is seldom implemented in symmetric encryption protocols and if it exists it is usually part of the protocol implementation, and not the actual maths used in the algorithm. Thus, AES has been chosen as per current industry standards and FIPS-140 specification.

Approach 1: Retrieval Phase

Objective: Allow the *chooser* to retrieve a record without from *sender's* database. Sender must not be able to identify the retrieved record, while *chooser* can only obtain one record.

Chooser's inputs: Index of the desired record in *sender's* database and a private encryption key (AES key).

Sender's inputs: Dataset of records related to individuals. As many *public/private* keys pairs as there are individuals described by the data.

Output: The above objective is met however, the number of records in *sender's* database is revealed to *chooser*.

1. The *sender* generates n sets of *public/private* keys pairs, and sends all *public* keys to the *chooser*, preserving the order in which they have been sent.
2. The *chooser* generates a key with a private encryption protocol, such as AES, later called *AES* key. It then uses the i^{th} *public* key received from the *sender* in Step 1 to encrypt the *AES* key and send it to the *sender*.
3. The *sender* does not know which *public* key has been used to encode the *AES* key, or which record has been selected, thus protecting the privacy of the *suspect*. The *sender* can then decode the cipher-text received in Step 2 using all *private* keys generated in Step 1, whilst preserving the order in which they have been decrypted. In this way n potential *AES* keys are created. Only the i^{th} one is the proper *AES* key; the other outputs are random sets of bits, which cannot be distinguished from ordinary *AES* keys.
4. The *sender* encrypts all records using appropriate keys decrypted in Step 3. Thus, the first record in *selected records* is encrypted with an *AES* key decrypted using the first private key generated in Step 1. Consequently the i^{th} record, which includes data about the *suspect*, is

encrypted using the *AES* key generated by the *chooser* in Step 2, sent to the *sender* encrypted by the i^{th} *public* key, and then decrypted using i^{th} private key. In this way the i^{th} record will be encrypted using the proper *AES* key.

5. The *chooser* gets n encrypted records, but using the *AES* key it is able to decrypt only the i^{th} record. Other records are unreadable to the *chooser* provided that the false keys generated in Step 3, and used to encrypt these records in Step 4, are not broken.

Figure 4-3 Retrieval Phase

4.6 Approach 2: Combined PET primitives

Protocols created for the purpose of searching datasets and retrieving objects of interest in private-manner do exist. These are examined in this section and a suitable protocol is selected and put forward for comparison with a solution made from a combination of PEqT and OT primitives (presented in Section 4.5).

Some pseudonym-based systems, such as those proposed by Biskup and Flegel in [75, 76], provide adequate functionality and could, in theory, fulfil the requirements of the data acquisition process. In fact, a solution based on pseudonyms would most likely gain the acceptance of Society, since it is an easy to comprehend approach that can provide information theoretic security for the parties involved. However, it would not meet Req. 7, as it does not scale well, and would be impractical for a system with large amounts of records. The protocol for private on-line transactions presented in [73] provides some functionality of what is required by the data acquisition process. It allows for the retrieval of digital goods based on publicly available index; however, it also provides a private comparison functionality that ensures that the buyer has enough funds to purchase the goods. Therefore, it would be possible to modify this protocol in order to create an adequate data acquisition protocol. On the other hand, the PE protocol provides all the basic functionality required, and it is designed to handle multiple records in the request (which would help to satisfy Req. 1). Consequently, the PE protocol has been chosen as the suitable combined approach.

The operation of the PE protocol is described in Figure 3-6. It uses three different commutative keys to facilitate Searching and Retrieval. Identifiers of the records of interest are hashed by the *chooser*, while the identifiers of the records resulting from the Querying phase are hashed by the *sender*, and then compared using commutative 1-*n* PEqT primitive. The records themselves are encrypted under symmetric encryption with keys crafted from the hashed identifiers encrypted using commutative keys. The *chooser* then retrieves all the records in the dataset, uses hashes of the identifiers commutatively encrypted by the *sender* to locate the records of interest, and decrypts these records using keys obtained in a fashion similar to the 3Pass primitive. Thus, the protocol requires two different forms of encryption: commutative and symmetric encryption. As previously defined, AES is a good choice for the symmetric cryptographic operations. However, in this approach the design is based on commutative properties of the PH protocol and not RSA, as public-key functionality is not required and PH cryptosystem does not provide the additional avenues for attack, namely the large-number factoring problem, that RSA is based on.

4.7 Evaluation

4.7.1 Experiment Design and Implementation

For the purpose of performing an initial evaluation on the different ways to implement the PET solution for the data acquisition process, testing needed to be performed in order to establish which approach fulfils the requirements. The requirements that could be considered without experimentation have been discussed already in Section 4.5 and Section 4.6. The key purpose of this empirical evaluation is to allow the comparison of protocol performance. Thus, the two key performance factors that need to be evaluated are: the time required for computations; and the amount of communication taking place. Typically, such protocols are evaluated using the notions of computational and communicational complexity. However, as discussed earlier, the notion of computational complexity cannot be used where the time for a single operation is different between the protocols. This is clearly the case in the protocols shortlisted, as the first makes use of public-key cryptography, and the other does not. Consequently, some measurements are required to establish the

average times for the operations in the process as well as the total time for the protocol run. These are evaluated on a single machine acting as both the *chooser* and the *sender*, with a single database to store the records. This choice is natural for empirical evaluation of OT and SPIR protocols since balancers can be used to distribute the load between the participating parties [90, 92] and semi-trusted parties can be used to take-up some of the computational burden [101]. Consequently, performance does not have to be measured on a per-party basis, and only the total time required to reach the result is needed to evaluate the performance for given set of input parameters (as shown in [78]). During the performance testing, the processes were organised into a series as not to affect each other. This means that some optimisation can be added to speed-up the operation of the protocols, but the measurements illustrate the worst-case scenario, and thus allows for direct comparison of the two approaches.

The main variables in the experiment were the size of the dataset being queried n , and the number of interesting records m . In the experiments the output from the Querying on the *sender's* dataset is simulated by a data table containing 128-bit MS SQL Globally Unique Identifiers (GUIDs) acting as record identifiers (ri) and randomly assigned text of 1kB in size that acted as the data content – the information about the record (as illustrated in Figure 4-4). The test script has randomly selected m different GUIDs prior to the simulation, in order to act as the identifiers of the interesting records requested by the *chooser*.

Between the two approaches there are four encryption protocols and one hashing protocol required to build the experiments. Since, it is not advisable to create a new implementation of encryption protocols [66], where tried and tested crypto suites exist, Bouncy Castle API is employed as the basis for the test implementations. The proof-of-concept protocols themselves are implemented using the C# .NET programming language in order to speed up the development. Since, both approaches are developed using C# .NET the fact that there are other languages, such as C++, that can produce applications performing faster is irrelevant to this demonstration.

GUID	DATA
0E138AC0-BD34-40DC-A1FB-0000238D746B	Cras nec tellus elit. In hac habitasse platea dictumst. Proin lectus elit, molestie sit amet iaculis

	quis, consectetur in metus ...
19E5B1CF-F6FC-41DB-9779-0000562A56A7	Phasellus pulvinar consectetur metus, vel auctor magna malesuada auctor. Suspendisse potenti. Donec eu leo non diam ultricies eleifend. Cras sed lorem elementum erat auctor egestas in at nulla ...
4629E748-7A74-42D4-9D5C-00006633D3EC	Donec et neque dui, at volutpat urna. Praesent ipsum sapien, laoreet quis tincidunt at, semper at ante ...
7FA90D8F-40E7-44F2-BBF4-0000BDAADC75	Aliquam interdum lectus sagittis mauris sodales sodales. In id aliquet elit ...
...	...

Figure 4-4 Test Dataset

SRA Implementation

An advantage of SRA over other commutative encryption protocols is that it can be implemented using common cryptographic suites, with only small changes necessary. In SRA and RSA, the encryption (Eqn. 4-1) and the decryption (Eqn. 4-2) operations are identical. These are performed modulo n , which is a product of two large primes p and q . From these primes $\varphi(n)$ is produced (Eqn. 4-3), which is used to generate the encryption keys. The encryption exponent e is generated randomly from the range shown in Eqn. 4-4, and so that e is co-prime with $\varphi(n)$. Then the decryption d is calculated as the multiplicative inverse of e modulo $\varphi(n)$ (Eqn. 4-5). All this is identical for SRA and RSA with the exception that in RSA the parties share their public keys that consist of the encryption exponent e and modulus n , but keep $\varphi(n)$ and the decryption exponent d secret, whilst in SRA both primes p and q are shared, or even public, and the both exponents need to be kept private.

$$C = M^e \bmod n \quad \text{Eqn. 4-1}$$

$$M = C^d \bmod n \quad \text{Eqn. 4-2}$$

$$\varphi(n) = (p - 1)(q - 1) \quad \text{Eqn. 4-3}$$

$$1 < e < \varphi(n) \quad \text{Eqn. 4-4}$$

$$de \equiv 1 \bmod(\varphi(n)) \Leftrightarrow d = e^{-1} \bmod(\varphi(n)) \quad \text{Eqn. 4-5}$$

In order for the keys to commute they need to be generated using the same primes, and therefore the crypto suite needs to be modified to accept the primes as inputs to the key generation process, which in most crypto suites is performed by an atomic procedure. In the case of the Bouncy Castle library, the `RSAPublicKeyPairGenerator` class

class from the `Org.BouncyCastle.Crypto.Generators` had to be modified to achieve this.

Commutative PH

RSA is based on the PH protocol, and thus there is some deal of similarity between them. However, PH algorithm, does not support public-key operations, as the decryption key can be easily calculated from the encryption key. Nor is it a symmetric algorithm, as two different keys are used for encryption and decryption. Therefore, PH can be considered as an asymmetric private-key encryption algorithm, and this can explain why PH cannot be found in any openly available cryptographic suite. However, thanks to its common elements with RSA, only small modifications are required to the cryptographic suites. Eqn. 4-6 and Eqn. 4-7 show PH encryption and decryption functions respectively.

$$C = M^e \text{ mod } p \quad \text{Eqn. 4-6}$$

$$M = C^d \text{ mod } p \quad \text{Eqn. 4-7}$$

Both operations are performed modulo of a large prime p , and different keys, exponents, are used for encryption (exponent e) and decryption (exponent d) in this algorithm. Thus, the main difference between RSA and PH is that RSA uses modulus made of a product of two primes, while PH uses only a single prime p for the modulus. Consequently, the RSA engine can be used to perform the operations, if the prime p is provided as an input instead of n . In addition, the encryption exponent e is randomly chosen in a way analogous to the RSA exponent, with the difference being that the upper limit of the range for e is different, and that e needs to be co-prime with $(p - 1)$:

$$1 < e < (p - 1) \quad \text{Eqn. 4-8}$$

Then, exponent d is calculated as:

$$de \equiv 1 \text{ mod } (p-1) \Leftrightarrow d = e^{-1} \text{ mod } (p-1) \quad \text{Eqn. 4-9}$$

Unlike RSA, it is easy to calculate the decryption key from the encryption key, thus, e and d must remain secret. Again, the modifications necessary to implement PH protocol with use of a crypto suite are limited to the generation of the keys that in the case of the Bouncy Castle library need `RSAKeyPairGenerator.cs` to be modified.

4.7.2 Empirical Evaluation

The Bouncy Castle crypto library has been used to produce implementation of RSA, SRA, PH and AES encryption schemes, as well as the SHA-256 hashing protocol. These implementations are used to gather performance data relating to the generation of cryptographic keys, as well as the encryption and decryption operations. The results of the measurements are based on an average time for the execution of 1 million operations for each protocol. GUIDs acted as input to hashing protocols, while the produced hashes are used as an input to the asymmetric algorithms (as in the OT and PE protocols). The AES128 protocol is tested using a 1kB input (that is approx. 150 words of ASCII text) that is expected to be larger than necessary to simulate records returned by the dataholder (similar amounts of data are used in [108] and [78]).

The test is conducted on a test machine running Microsoft Windows XP Professional with an AMD Turion 64 X2 Mobile 1.58GHz CPU, and 3GB of RAM. The results are provided in Table 4-1. From this comparison table it can be gathered that operations such as hashing and AES key generations are performed almost at wire-speeds and can be safely considered as negligible in this consideration. The large difference between key generation times for different asymmetric protocols can be explained by the fact that RSA generates a new pair of primes p and q each time, while in PH and SRA, the primes are common between the parties, and therefore are generated only once per protocol run, or are part of the system. RSA has the smallest encryption time when compared to other asymmetric protocols examined, since most RSA keys use a default encryption exponent e that is reasonably small (such as 0x10001) in order to speed up the encryption and signature verification processes, but not too small as not to expose the protocols to attacks described in [113]. On the other hand, the SHA and PH implementations cannot use such common choices for the encryption exponent, since these protocols expose the primes (used to calculate the decryption exponent) to the other parties, and thus the decryption exponent can

be easily derived from the encryption exponent. PH and SRA use larger randomly generated numbers as the encryption exponents, and this results in longer encryption times.

	Conditions (bits)		Results (ms)		
	strength	Data Size	key generation	encryption	decryption
RSA	1024	128	723,344	766	47,266
SRA			7,128	21,125	51,297
PH			7,250	22,389	50,594
SHA	256		-	31	-
AES	128	1k	31	551	564

Table 4-1 Cryptographic operation performance measurements in nanoseconds (*ns*)

Table 4-2 shows the complexity of OT-based Approach 1 in an arithmetical format, while the complexity of the PE-based Approach 2 is shown in Table 4-3. The simplicity of operation of the OT-based approach contributes greatly to its high computational complexity. Since RSA is used as a trapdoor function of the OT, each record that needs hiding requires a separate RSA key. Key generation times for RSA are large in comparison to other asymmetric protocols discussed in this chapter, since new primes p and q are generated for each new key. This is compatible with the typical use of the RSA protocol since each user needs only two different sets of keys (a separate set for signing and encryption). On the other hand, the PE-based approach only uses three asymmetric encryptions (commutative PH) keys through the protocol, so the preparation phase for the protocol run is almost negligible. The computational complexity tables have been used to plot graphs illustrating the performance differences between the two approaches. Figure 4-5 depicts the total running time for both OT- and PE-based approaches including the preparation time, which is the time used to perform operations that are independent from the enquiry and can be performed prior to the protocol execution. The logarithmic graph shows that the total running time for the OT-based solution is more than an order of magnitude higher than this characteristic for the PE based protocol, the exact values calculated are available in Appendix C. Figure 4-6 illustrates that when the preparation time is eliminated the performance of both protocols for $m = 1$ is of the same magnitude, with PE-based solution taking on average two thirds of the time used by the OT-base solution. It is worth noting that for both approaches the run-time is almost linear for the varying size of the dataset.

PE uses commutative encryption which, employed in a 3Pass protocol, adds similar benefits to m - n SPIR protocols as public-key encryption had on securely exchanging information with multiple parties. Prior to emergence of public-key cryptosystems a party wanting to communicate securely with n other parties would need n different symmetric keys. Likewise, in order to achieve m - n SPIR using encryption other than commutative m different runs of 1- n SPIR would often be necessary, while with commutative encryption and systems like PE as little as a single additional encryption operations is necessary to retrieve one more record. This is evident looking at the above complexity tables. In the OT-based approach in Step 3 of the OT phase the *sender* needs to decrypt the ciphertext received from the *chooser* $O(m \times n)$ times and then encrypt all n records m -times in Step 4 resulting in $O(m \times n)$ complexity of symmetric encryption operations. The equivalent operations in the PE-based solution include the generation of the symmetric keys by encrypting (using the PH cipher) the hashed record ID and using the hashed result to encrypt the records using AES. Consequently, the equivalent operations in the PE require only $O(n)$ operations, each. Figure 4-7 presents this difference between the two approaches. Still, both protocols need to be praised for the use of symmetric encryption in hiding the data and using asymmetric ciphers to selectively transfer the symmetric keys between the parties. This is an optimal technique inspired by the PKI, and praised by Shundong et. al [93].

As mentioned already the OT protocol, used in Approach 1, is a simplistic protocol that is useful for illustrating the process of data acquisition. Its extensive preparation step requires generating n different RSA keys, which makes it suboptimal for the requirements. However, the characteristic depicted in Figure 4-7 is similar to other OT protocols. Namely, OT protocols are usually optimised for handling a single request from a large dataset per round of the protocol, hence the 1- n OT is the most common type of OT in use.

		Symmetric Crypto. operation	SRA Cryptography			RSA Cryptography		
			key gen.	encrypt.	decrypt.	key gen.	encrypt.	decrypt.
PEqT	Preparation	-	$O(2)$	-	-	-	-	-
	Step 1	-	-	$O(m)$	-	-	-	-
	Step 2	-	-	$O(m)$	-	-	-	-
	Step 3	-	-		$O(m)$	-	-	-

	Step 4	-	-	$O(n)$	-	-	-	-
OT	Step 1	-	-	-	-	$O(n)$	-	-
	Step 2	-	-	-	-	-	$O(m)$	-
	Step 3	-	-	-	-	-	-	$O(m \times n)$
	Step 4	$O(m \times n)$	-	-	-	-	-	-
	Step 5	$O(m)$	-	-	-	-	-	-
Total Complexity		$O(m(n+1))$	$O(2)$	$O(2m+n)$	$O(m)$	$O(n)$	$O(m)$	$O(m \times n)$

Table 4-2 Computational complexity of the OT-based approach.

This table illustrated the complexity of each step in the OT-based approach to data acquisition being evaluated in this chapter. Since, each operation require different amount of computation each column represent different operation. The total complexity is the sum of the complexity for the given operation.

		Symmetric Cryptography		Asymmetric Cryptography		
		encryption	decryption	key gen.	encryption	decryption
PE	Step 1	-	-	$O(3)$	-	-
	Step 2	-	-	-	$O(m)$	-
	Step 4	-	-	-	$O(2m)$	-
	Step 5	$O(n)$	-	-	$O(2n)$	-
	Step 6	-	-	-	-	$O(2m)$
	Step 7	-	$O(m)$	-	-	-
Total Complexity		$O(n)$	$O(m)$	$O(3)$	$O(3m+2n)$	$O(2m)$

Table 4-3 Computational complexity of the PE-based approach.

This table illustrated the complexity of each step in the PE-based approach to data acquisition being evaluated in this chapter. Since, each operation require different amount of computation each column represent different operation. The total complexity is the sum of the complexity for the given operation.

The literature review has identified the PE protocol as the only protocol that is optimised for retrieval (in a single round) of m records from a dataset. Since the data acquisition process calls for a solution that allows for retrieving of multiple records per enquiry (Req. 1) the PE protocol is the most likely choice to satisfy this condition. The load is almost linear to the number of records in the database. Since each record can be processed independently, there are no technical limitations to processing any number of records that can be stored in database. Therefore, PE would most likely meet Req. 7. The fact that it is possible does not necessarily mean that it is feasible. It would take at least eight days to process 15 million records on a standalone PC with a specification similar to the one used to generate the test data. Consequently, Req. 3 would not be met, as if more than one enquiry would be run on such PC, the system would take more than two weeks to provide the response.

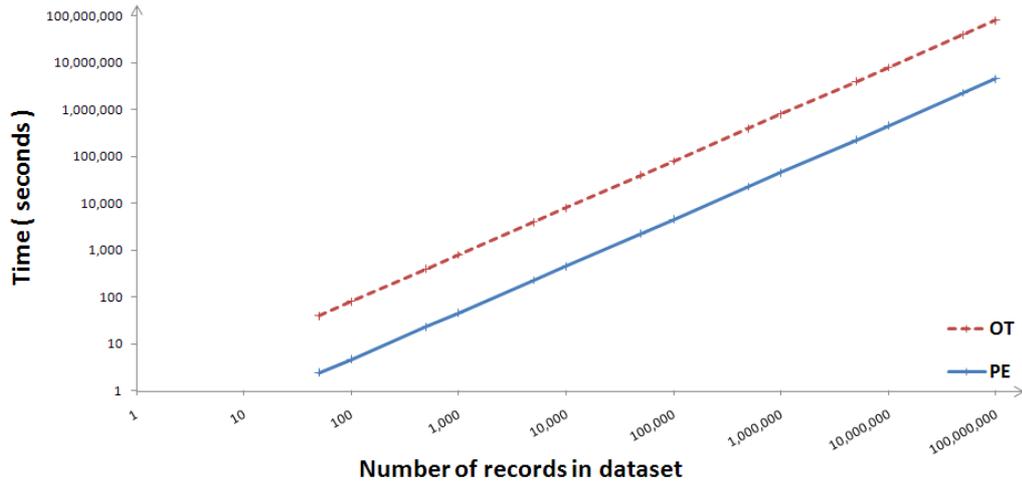


Figure 4-5 Total running time for both approaches including preparation time. Plotted for n varying from 50 to 100 million, and constant m equal to unity.

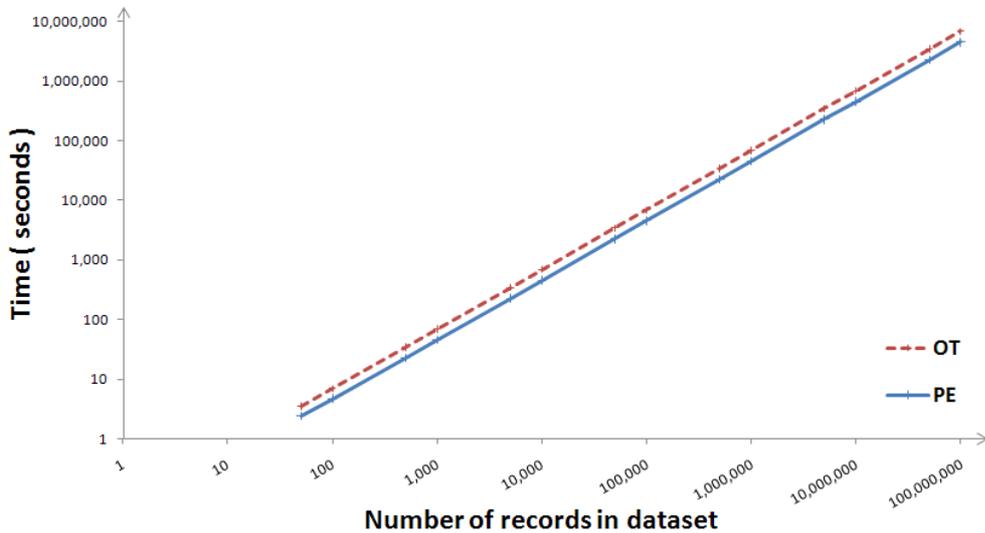


Figure 4-6 Data Acquisition processing time excluding preparation time. Plotted for n varying from 50 to 100 million, and constant m equal to unity.

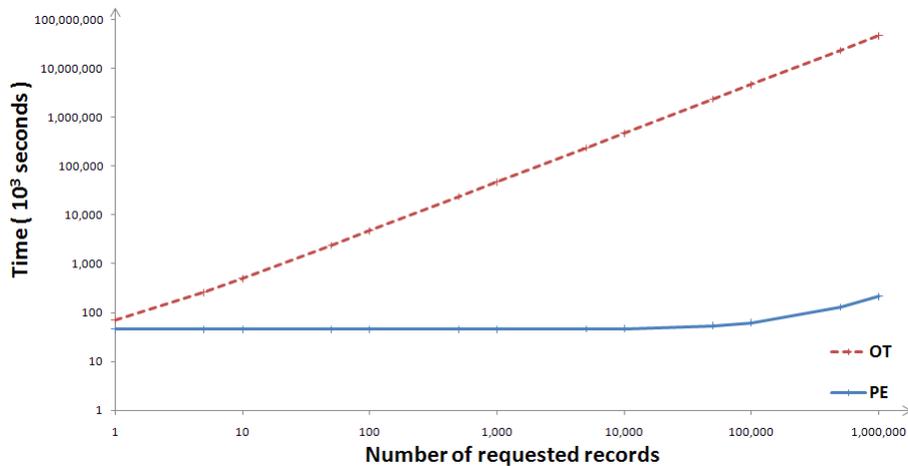


Figure 4-7 Performance of both approaches for varied number of records ($1 < m < n$) and constant size of the database ($n = 1$ million) [excludes prep. time for OT].

If the results achieved are compared to other, similar systems for which data is available, PE still looks to out-perform the competition. The hardware-based TC solution presented in [108] initially requires the shuffling of all records before any request can take place. Shuffling 15 million records of 1.5kB in size would take more than a year using the SCOP. While the results discussed above show eight days as the total time required for the data acquisition of one record on a dataset of the same size, they concern records that are 1kB in size. However, increasing the record size to 1.5kB in the PE-based solution, increases the processing time by less than two hours. The shuffle needs to be performed only every \sqrt{n} requests, while the PE protocol requires full run each time. Nevertheless, under PE a single request could contain \sqrt{n} records.

Research presented in [78] describes protocols similar in functionality to the suggested approaches to the data acquisition process. The empirical results presented there suggest that the therein-discussed protocols are more efficient than PE, however, this empirical evaluation does not specify the size of the data, nor does it include the preparation phase in the considerations. The table illustrating the complexity of the protocols described by Cristofaro et. al. shows that it is similar to this of PE for $m = 1$, while the shorter total processing times are due to the use of the C as programming language, and a more powerful test PC. However, just like the OT-based solution, the total processing time increases in line with the increase of the number of interesting records, whilst it is almost constant for PE-based solution.

Both approaches as presented, without the use of any balancers, require only a few rounds of communication. Consequently, if the time given for an enquiry was the statutory 14 days, the data could be exchanged using physical media rather than over the Internet, thus eliminating any limitations for the record and dataset size.

4.7.3 Feedback from practitioner

Following the initial evaluation, a member of ACPO and ACPO in Scotland (ACPOS) has been approached in order to obtain qualitative feedback for the results and clarify the investigative process. Also, a research poster has been presented during the second SIPR Annual Conference. The Detective Superintendent (DS) that cannot be named as the interview was based on Chatham House Rule is a member of

the UK Data Communications Group – which comprises members of ACPO (and ACPOS), HMRC, and representatives of different CSPs. The DS was very interested in the research and stated the following in respect to the drafted requirements and assumptions:

- 1) Currently, an enquiry for communication data needed in a case where life of an individual is in danger takes minimum of 30 minutes.
- 2) Police have direct access to subscriber data, such as name and address, for all major CSPs.
- 3) Collateral damage to a data subject should be minimal as all enquiries are inspected by a designated person before being sent to dataholders, and are scrutinised in court of Law (if charges are pressed against the data-subject, or data acquired is used as evidence).
- 4) On some occasions, investigators need to postpone their enquiry, until they have enough background for the check (in order to protect data subjects against the collateral damage) or the subject is in custody (if there is suspicion that the data subjects may be informed about the enquiry taking place).
- 5) In face-to-face enquires, law enforcement officers can ask general question allowing the person being interviewed to choose to amount of detail provided to the investigators. The technique is sometimes referred to as *dilution*. This is impossible in the current state-of-art in the digitalised enquiries, as it is often considered as *fishing for evidence*.
- 6) In some occasions, location data or other leads may need to be used to identify possible suspects and witnesses of an incident.

An example of a typical investigation is the use of evidence gathered from third parties given by the DS is the case of the Soham Murders (with accused being Ian Huntley), Maxine Carr provided an alibi for Huntley by stating he was with her at a specific location at the time of the murder. However, her phone location records obtained from her CSP showed she was 100 miles away. Therefore, the investigators could prove that she was lying. Newspaper articles confirm that the communication

data was extensively used during the case, where the call timings are used to place individuals at different locations in a timeline of the events [123].

4.8 Requirements review

In light of the feedback from the practitioner and the empirical results, Req. 3 needs to be altered to reflect law enforcement expectations. Since, the minimum time for the complete enquiry is 30 minutes the protocols must be able to provide results even from large databases in less this time. However, preparation time that takes place before the enquiry can be permitted. The update is:

Req. 3 Allow for efficient and timely data retrieval. (The protocol run excluding preparation should take less time than 30 minutes that it currently takes investigators to obtain investigative data in emergencies.)

Additionally, since different clues may need to be used to identify a potential suspect or a witness the data acquisition process should allow for a complex private matching criteria, that would allow selection of the records based on more than one column of data, and possibly allow for fuzzy matching. Thus, the following requirement should be added:

Req. 8 Provide a technique for multiple selection criteria of interesting records, and allow for fuzzy matching on the selection criteria different than record ID.

4.9 Conclusions

This chapter has identified and refined requirements for the data acquisition process. While some of these requirements are technical, such as the expected performance, other relate to the legal and social aspects of the process. These requirements were used to select a suitable PET primitive needed to facilitate privacy-preserving investigative data acquisition platform, and are used later in this thesis to evaluate the platform itself.

The PET primitives identified in the Literature Review have been scrutinised, and PE primitive allowing for private retrieval of records forming an intersection between two sets (in this case the set of potential suspects, and the set of all data subjects in a

database) was chosen as a suitable protocol for the task-at-hand. In comparison to other protocols PE, run-time is almost independent from a number of suspects in an enquiry, while other protocols show almost linear increase in the total run-time with an increase in m .

Despite the PE-based approach being the best performing, it does not meet all the requirements (neither does any other considered protocol). Whilst it is capable of processing datasets of any size, as the records are processed one at the time independently from each other, it would take eight days to perform a retrieval on a database of 15 million using PE-based solution implemented in a managed C# .NET code and run on a computer similar in specification to the test setup. This could be improved with the use of different programming languages for the implementation and fast exponentiation, but still an enquiry would not complete in 30 minute as necessary (30 minutes is the current minimum time taken for an enquiry).

It is interesting whether such a solution could gain the acceptance of the general public. By choosing PE, a protocol based on commutative encryption that is relatively easy to explain, it would be likely to gain approval from the decision makers, as they already understand basic encryption terminology. However, the system sends all the records in the database to the *chooser*, using encryption to hide the unselected records from the authorities, and members of the general public suspect public authorities of having computational power to break cryptosystems. For this reason, some additional measures should be built into the platform, in order to ensure that the *chooser* can prove that the data irrelevant to investigation has not been decrypted.

Another functionality that PE does not seem to provide is handling multiple selection criteria to identify interesting records. For example, if law enforcement officers are looking for a white female in her twenties, they cannot make such an enquiry privately against a corporate HR databases. It would be ideal if such cases could also be catered for, enhancing the authorities ability to identify potential suspects.

Chapter 5

Novel Data Acquisition Platform

5.1 Introduction

Investigative Data Acquisition Platform (IDAP) is formed by improving on the shortcomings that the PE primitive has in an investigative scenario. These shortcomings have been derived from the initial evaluation presented in Chapter 4 and include:

- Long processing times.
- Lack of capability to retrieve records matched on multiple selection criteria.
- Potentially low acceptance of the SPIR-based techniques.

The improvements that aim to address these shortcomings introduce a *dilution factor*, which is a numeric value that specifies the level of anonymity required for a given investigation. With this factor the data subject behind the interesting record should

feel assured that a constant level of privacy is provided to all individuals independently from the number of interesting records in an investigation.

The chapter presents a technique for forming complex privacy-preserving queries, without affecting the complexity of the protocol. This relies on joint hashing of the different selection criteria together, and using these as an input to the PEqT protocols. In order to gain approval of the general public a semi-trusted proxy is added as a novelty, in order to ensure information theoretic privacy of data-subjects whose records are not defined as interesting. Quantitative and qualitative evaluation of the complete approach to investigative data acquisition is planned and test implementation of IDAP is implemented.

5.2 Methodology

Chapter 4 has defined the requirements for data acquisition process and has shown that an information retrieval system based on PE primitive would be capable of meeting most of those requirements. PE is possibly the only information retrieval PET protocol that has almost constant processing time for enquiries, with a varying number of interesting records m . This suggests that it is likely to be the most efficient m - n SPIR primitive. Still, processing of 15,000,000 records would take 8 days on the test bed used in experiments presented in Chapter 4. This could be shortened to less than 14 hours if the program is written in C (or C++) and run on a host similar to the one used in producing results for [78]. Thus, there is a clear need to improve the performance if the data acquisition process is going to employ the PE primitive. Other drawbacks of using SPIR-based PETs in obtaining investigative data are the lack of explicit functionality to retrieve records based of multiple selection criteria, and possible low levels of public acceptance (or understanding) of the SPIR concept, as it requires transferring data unrelated to an enquiry.

In order to design and implement IDAP, the shortcomings of the PE-based solution for data acquisition identified in Chapter 4 need to be mitigated. The modifications proposed in this thesis are based on the results of the initial evaluation and inspired by controls used in other PET primitives. The complete IDAP system is then defined and evaluated. The performance and the security of IDAP are discussed against the PE protocols and the requirements. This provides the pure quantitative evaluation of

the platform, while, in order to gain understanding of the public attitudes towards the protocol, a survey was carried out among IT security and privacy experts. Results from both experiments and survey are presented and discussed in Chapter 6.

5.3 IDAP Design

The Investigative Data Acquisition Platform (IDAP) is formed on the basis of the PE primitive extended to fulfil the requirements outlined in Chapter 4. There are three modifications to this primitive that are required in order to facilitate the requests for investigative data. The resulting IDAP is a novel privacy-preserving approach to the data acquisition process.

5.3.1 Lowering Processing Time

There is a clear need to minimise the processing time required for each run of the protocol in large databases, such as those belonging to ISPs and mobile telephony providers. Theoretically, in order to maintain the privacy of the suspects, the *sender* needs to process all the records in the database per enquiry. Only in this way no information about an interesting record is revealed and the correctness of the PE scheme can be proven under the rules of MPC [124]. Thus, if the data acquisition platform would use the PE primitive without modifications, the system would not be capable of processing any urgent requests due to the run-time required per enquiry, and this would be a major drawback. A possible mitigation against this could be to limit the numbers of records that are processed and sent by the *sender* per enquiry. This would also lower the communicational complexity that has not been taken into direct consideration in this thesis.

Privacy of the alleged suspect should be protected, but if the probability of the *sender* guessing the ID of the interesting record is for example 1:100,000 and not 1: n (for n being the size of the population or a large dataset), and the dataholder has no other information that could help infer the identity of the suspect, this research argues that the privacy of the suspect and the investigation is maintained. On occasion during traditional, i.e. face-to-face, information gathering exercises, Police Officers would use a concept of *dilution* – hiding the suspect's identity by asking open-ended questions about a larger group of individuals rather than about a single person. This

is a widely accepted technique, however, in a digital environment it is impossible to build a system that would maintain privacy, while providing answers to such general questions. Consequently, any attempts of investigators to *cast their net wide* during electronic investigations are prohibited and treated as *fishing-for-evidence*. Taking in consideration that using the PET-based system the investigators will not get more data than required for their inquiry, limiting the set of records that are processed per enquiry should be acceptable. In comparison to methods used by the Internet users to protect their identities, there are an estimated 100,000 active TOR clients any point in time and 1,500 TOR relays [70]. Thus, at best TOR users can expect to have only 1:100,000 privacy ratio.

The problem is to decide on the technique of narrowing down the scope in a way that ensures interesting records are among the results returned. If the list of the record identifiers is public, such as the list of the IP addresses or telephone numbers served by a given network operator, the *chooser* could simply select a number of random records from such directory in order to hide the true target of the investigation. Possibly, the *chooser* would first need to obtain a list of unused addresses from the provider, or at least know the percentage of unused addresses, in order to ensure that the number of unused addresses accidentally included in the request does not reduce the level of privacy. However, in the case where a list of IDs is not publicly available, it would be possible to split the PE protocol back into separate parts: PEqT; and SPIR. In this way, the PEqT can be used during an initial preparation phase run against the whole dataset, and that the information retrieval would be performed against a smaller set of records.

It was previously mentioned that PE has almost constant processing time for enquiries with increasing number of interesting records (for low m as shown in Figure 4-7). However, if the number of records retrieved per enquiry is lower than the size of the dataset it would be ideal if there is a constant level of privacy provided to each potential suspect. In the data mining field, there are already k -anonymity models that ensure that any privacy-protected statistical data record links to at least k different identities [85]. Consequently, providing a controlled level of privacy to the data-subjects. Relating to the concept of *dilution* used by the Police, a number of records requested per each interesting record can be defined as the *dilution factor* –

o. This factor could be changed before each protocol run in order to allow investigators to dynamically choose the appropriate level of protection for the given investigation, the data subject, and the data controller.

The proposed improved PE protocol operates by creating a single encrypted table of identities and allowing the investigators to privately match (using PEqT primitive) the identities of their suspects against this table. As the outcome of the private match operation the *chooser* would find out encrypted IDs of the interesting records. Then to perform an investigation the *chooser* would select $(o - 1)$ records at random per each interesting record from the encrypted table of IDs. The double encrypted IDs of the selected records would be communicated to the *sender* and remaining operations of the PE protocol would be run only on the selected records. Thus, the total number of requested records would be a product of the number of interesting records and the dilution factor, $(m \times o)$.

The described technique would introduce the potential for few different data controllers to collaborate and possibly identify the records of interest by checking for overlaps (intersection) of the requests made by the investigators to the collaborating data controllers. However, in the cases when the data is being retrieved from large databases that require use of the dilution technique during data retrieval process, the interesting records would usually be identified by a mobile phone number, or an IP address. Phone numbers and IP addresses are thus unique to the operators and their assignment can be obtained from call and network routing tables, respectively. Consequently, in most cases, the investigators would only need to ask a single operator for information about a given identity, and there would be no intersections of the requests. This fact makes most investigations equivalent to a single database SPIR allowing for dilution to be applied, with no adverse affect on the privacy of the data-subjects. The description of the improved protocol is as Figure 5-2, Figure 5-3 and Figure 5-4.

In this improved protocol the initial processing depends on the size of the dataset – n , but it needs to be performed only once in a given period of time. However, the remaining operations are run on limited dataset. Figure 5-1 illustrates the processes taking place in this improved version of PE protocol.

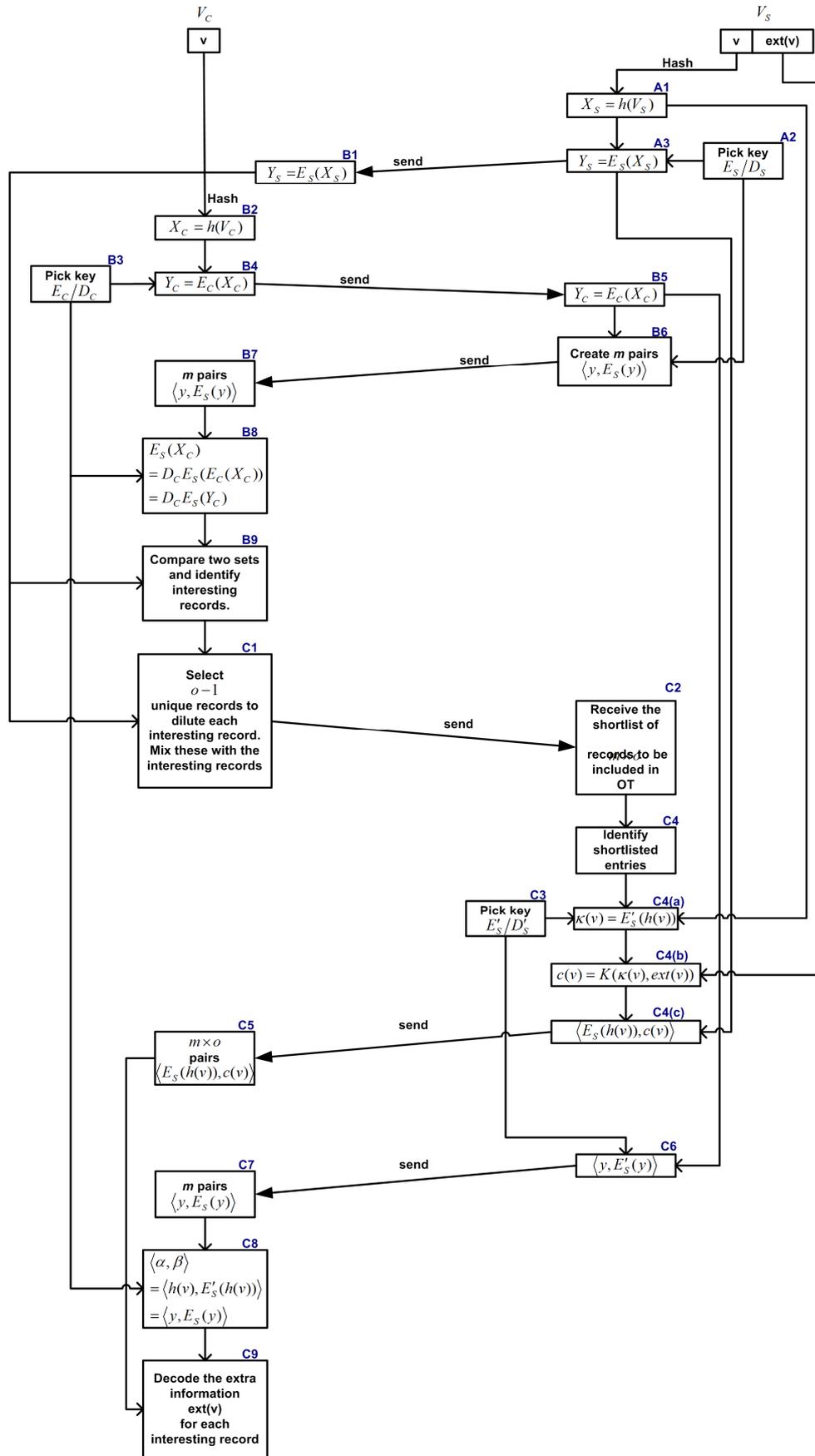


Figure 5-1 Process flow of the protocol incorporating the *dilution factor*

Phase A – Preparation:

Objective: Encrypt the identities in the *sender's* database, in order to facilitate private matching of the records.

Chooser's inputs: *chooser* does not take part in this step.

Sender's inputs: dataset V_S containing identities related to the records held in the *sender's* database. Selected group Z_p^* .

Outputs: A list of PH encrypted identities for use in private equijoin.

1. *Sender* applies hash function h to the elements in the input set V_S , so that $X_S = h(V_S)$.
2. *Sender* picks an encryption PH key E_S at random from a group Z_p^* , where p is a strong prime.
3. *Sender* encrypts each $h(v) \in X_S$ with the key E_S , the result is a list of encrypted identities $Y_S = E_S(X_S) = E_S(h(V_S))$

If more record needs to be added to the set these can be processed using steps 1 and 3, and then added to the list.

Figure 5-2 Lowering Processing Time Phase A – Preparation

Phase B – Searching:

Objective: Allow *chooser* to match identities of the interesting records to the privacy-protected identifiers of the records in the *sender's* database.

Chooser's inputs: A list of the interesting records V_C . Knowledge of the group Z_p^* being used for PH key generation.

Sender's inputs: A list of PH encrypted identities for use in private equijoin. Knowledge of the group Z_p^* being used for PH key generation.

Outputs: *Chooser* obtains a list of identifiers encrypted by the *sender* that can be used to request specific records.

1. Following a request for data, *sender* provides *chooser* with a complete list of encrypted identities prepared during Phase A, reordered lexicographically.
2. *Chooser* applies hash function h to the elements in set containing the identities of the interesting records, so that $X_C = h(V_C)$.
3. *Chooser* picks a commutative cryptography key pair, encryption key E_C and decryption key D_C , at random from the same group Z_p^* that was used by *sender* in the Phase A.
4. *Chooser* encrypts entries in the set X_C , so that $Y_C = E_C(X_C) = E_C(h(V_C))$.
5. *Chooser* sends to *sender* set Y_C reordered lexicographically.
6. *Sender* encrypts with key E_S each entry $y \in Y_C$ received from *chooser*.
7. *Sender* returns set of pairs $\langle y; E_S(y) \rangle$ to *chooser*.
8. *Chooser* decrypts each entry in $E_S(Y_C)$ obtaining $E_S(X_C) = DCESEC(X_C) = DCESYC$.
9. *Chooser* compares each entry in $E_S(X_C)$ to the entries of Y_S received in the Step B1 (Step 1 of Phase B). This way the interesting records can be identified.

Figure 5-3 Lowering Processing Time Phase B – Searching

Phase C – Retrieval:

Objective: This phase is used to retrieve the interesting records from the *sender*. It is in this phase that the enquiry can be narrowed down to the subset of records in the *sender's* database. The identifiers of all the records in this database are known from Phase B, as well as the list of the identifiers that belong to the interesting records. Based on this knowledge the *chooser* can request the interesting records, plus some additional records to minimise the data processing,

while maintaining privacy of the data-subjects.

Chooser's inputs: A list of identifiers encrypted by the *sender* that can be used to request specific records. Symmetric encryption scheme K .

Sender's inputs: A list of PH encrypted identities for use in private equijoin. Knowledge of the group Z_p^* being used for PH key generation. Record data $ext(v)$. Symmetric encryption scheme K .

Outputs: *Chooser* obtains the data $ext(v)$ for the interesting records.

1. After identifying the interesting records in Y_S the *chooser* selects at random $o-1$ other unique records from Y_S for each interesting record in V_C . These are the diluting records, that together with the records of interest form a shortlist for the enquiry. If the number of interesting records multiplied by o is greater than n , the size of the dataset V_S , then the complete Y_S is shortlisted.
2. Send the shortlist to *sender*.
3. *Sender* picks an encryption PH key E'_S at random from the group Z_p^* .
4. *Sender* identifies entries $h(v)$ from X_S that have been shortlisted and processes each shortlisted record in the following way:
 - (a) Encrypts $h(v)$ with E'_S to form the key used to lock the extra information about v , i.e. $ext(v)$, $\kappa(v) = E'_S(h(v))$.
 - (b) Encrypts the extra information using a symmetric encryption function K and the key $\kappa(v)$ crafted in the previous step: $c(v) = K_{\kappa(v)}(ext(v))$
 - (c) Forms a pair $\langle E'_S(h(v)), c(v) \rangle$.
5. The pairs formed in C4(c), containing a private match element and the encrypted extra information about record v , are then transferred to *chooser*.

6. *Sender* encrypts each entry $y \in Y_C$, received from *chooser* in Step B5, with key E'_S to form set of pairs $\langle y; E'_S(y) \rangle$.
7. Pairs $\langle y; E'_S(y) \rangle$ are then transferred to *chooser*.
8. *Chooser* removes the encryption E_C from all entries in the 2-tuples received in Step C7 obtaining tuples α, β such that $\langle \alpha; \beta \rangle = \langle h(v); E'_S(h(v)) \rangle$. Thus, α is the hashed value $v \in V_C$, and β is the hashed value v encrypted using E'_S .
9. *Chooser* sets aside all pairs received in Step C5, whose first entry is equal to one of the first entry of any two-tuples obtained in Step B9. Then uses the appropriate β tuple associated with a given interesting record as a symmetric key to decrypt the extra information contained in the second entry in the pair received in C5. This is performed for all the matching entries.

Figure 5-4 Lowering Processing Time Phase C – Retrieval

5.3.2 Allow multiple selection criteria

The PE protocol can be used to privately retrieve data if the data is identified by a single parameter, such as ID number, credit card number, IP address, and so on. However, this is not always the case. If the data acquisition process is used to find a suspect based on circumstantial knowledge, or a suspect's profile, the PE protocol would need to be modified. The query shown in Figure 5-5 shows the way the request from Figure 4-2 would be modified for such enquiry, here $sip_{1..j}$ stand for j secret input parameters (sip):

```
SELECT sip1, sip2, ..., sipj, rp1, rp2, ..., rpn
FROM source
WHERE ip1=ip_val1 AND ... AND ipi = ip_vali
```

Figure 5-5 mapped into SQL

A computationally expensive solution to this problem can be achieved by using symmetric encryption to lock the return parameters and then hiding the symmetric keys used with the commutative encryption keys unique to each value of the secret

input parameter. The *chooser* would then perform a separate PE-based retrieval of the asymmetric key for each interesting value of the secret parameters (such as age equal to 25 years). Since these asymmetric keys are commutative, the *chooser* would be able to decrypt the ciphertext containing the symmetric key that was used to lock records matching the selection criteria. Despite being computationally-expensive this solution has a unique benefit of allowing semi-fuzzy matching of the results if the underlying commutative protocol is ElGamal-based. This is the case as ElGamal (and its commutative form suggested by Weis in [55]) uses checksums that allow for verifying whether a given ciphertext can be decrypted with a given key. Thus, it would be possible to establish how many records match each secret input parameter. This solution has been published in [125], however, it is not suitable for large databases due to its high computational complexity.

In this thesis a simplified approach is proposed. Since, the query from Figure 5-5 replaces the *ri* parameter with J different *sip* parameters then the list of these J parameters could be used as a complex *ri* for use with PE-based data acquisition protocol. Thus, in steps B2 and A1 of the protocol presented in Section 5.3.1 a list of all values of given *sip* parameters would be hashed together to form records in sets V_C and V_S . In this way neither the security, nor the complexity of the protocol is affected by this improvement (if processing time required to produce hashed values is considered to be negligible).

5.3.3 Reassuring the Public

The initial design of IDAP proposed in the form of two approaches to data acquisition process investigated in Chapter 4 would shift the balance of the privacy protection from innocent individuals towards the suspect and the secrecy of investigation. Currently, the data acquisition process employed by the public authorities does not affect privacy of the data-subjects whose records are not of interest to the investigators, as there is no need to process these records. IDAP changes this as per each enquiry there are a number of records unrelated to the investigation returned to the *chooser*. The fact that the *chooser* is unable to decrypt these records does not change the fact that the records are being *processed* (according to the DPA definition of processing). As the anonymity in the PE protocol is based on hiding the interesting records among other records, some records

unrelated to the investigation will always be retrieved by the *chooser*. Thus, there is a need to ensure that the *chooser* does not abuse the system. As Juvnal put it:

Sad quis custodiet ipsos custodias? (Juvenal, Satires VI, 347)

Which translates to: *But who will watch the watchers?*

It is likely that providing government agencies with records of innocent individuals unrelated to any investigation would worry the general public. This is despite the data being encrypted in the way that renders these records unusable to the authorities i.e. secure against attacks in polynomial time. However, the public may worry that government organisations have enough computing power to break the encryption used in IDAP. There are few actions that may reassure the public that the data is safe. First, if the technique for minimising the processing time presented in Section 5.3.1 is employed, the chances that investigators will retrieve encrypted records of a particular individual that is not a suspect is small in large datasets. Thus, for a dataset with n records, during investigation with m interesting records and the dilution factor o , the probability of this event A can be defined as:

$$P(A) = \frac{(o-1) \times m}{n-m} \quad \text{Eqn. 5-1}$$

Consequently, for investigation with five interesting records, with dilution factor of a 1,000 and dataset consisting of a 15 million records, the probability of this event occurring during a single run of the protocol would around 3%. This also means that the investigators would need to first break the encryption key used by the *sender* to hide identities (Phase A), before they could attempt to obtain the data about a specific individual that is not a suspect, otherwise the probability of the encrypted record for this individual being provided to them would be small. Thus, if the identity of a data subject were never encrypted under the same key as the data records then investigators would need to successfully brute force two separate keys in order to link any record not declared as interesting to an identifiable individual. Otherwise, the information would be unintelligible, or random.

The merits of the above discussion could certainly improve the perception of the system. Still, currently there is grater trust in security processes than encryption, as

inferred from [79, 126] . The solution proposed in this thesis in order to reassure the public is to introduce a semi-trusted party into the protocol. Such a party is often used in PET protocols in order to balance the computational and communicational complexity between the participating parties or to off-load processing from the participants [101]. This party would become a *proxy* between the investigators and the dataholder, however, unlike in other PETs, the purpose of this proxy is to ensure that investigators get only the interesting records and all other records sent to them during the PE protocol are discarded, thus the solution to gain the public's trust is proposed in Figure 5-6.

The semi-trusted party should thus have no interest in finding out the object of the investigation (the interesting record ID) or the content of the data records returned by the dataholder. For this reason, it is suggested that the role of this party should be conducted by an independent body trusted by the public. In the UK for example, the Information Commissioner's Office (ICO) is such an independent body that ensures DPA and RIPA are adhered to. Thus, the ICO would be an ideal organisation to become the *proxy* and could help restore the natural order, where the rights of the innocent are put ahead of the secrecy of the investigation.

1. All communication between *chooser* and *sender* goes through *proxy*.
2. Chooser provides *proxy* with the identifiers of the interesting records encrypted by *sender*, $E_S(h(v))$. This is done over a secure channel or with use of a 3Pass protocol once the parties are authenticated.
3. At the stage where data is transferred from *sender* in Step C4, *proxy* filters the response and discards the records that were not specified by *chooser's* request, this is the records other than the ones identified in Step 2.

Figure 5-6 Reassuring the public by introducing semi-trusted third parties

It must be noted that the party that is chosen to become the *proxy* must not cooperate with the *sender* or the protocol will be broken, since simple matching

exercise would reveal the identities of the interesting records (but not the data). A key concept is that the *proxy* has no incentives to find out the detail of the investigation so it is going to cooperate with the *sender* in order to establish the identity of the suspect. On the other hand, if the need arises to verify the *chooser's* requests in front of a court of law, the *proxy* and the *sender* could work together to establish the identities of the records requested by the *chooser*. This is analogical to the two-man rule commonly used in security to ensure that a single person is unable to abuse the system.

In law if a party refuses to provide the evidence needed by another party a commissioner can be appointed to gather evidence listed on the Specification of Documents prepared by the requesting party. This commissioner then verifies whether the requesting party needs any given piece of potential evidence, and provides this party with only the documents relevant to the investigation. The process is referred to as *commission and diligence* [28, 29], and this is the legal justification for the introduction of the semi-trusted *proxy* into the data acquisition process.

In a scenario involving two parties, the *chooser* and the *sender*, the *chooser* would be capable of keeping the retrieved records that are irrelevant to the investigation for a period needed to decrypt them using a brut-force approach. This period should be long enough for the data to become outdated. However, there is a residual risk and worry that some public authorities would attempt such an attack on the system. In order to limit this vulnerability of the system, and provide a way of auditing the process by an independent organisation the functionality of the proxy is introduced into the platform design. Nevertheless, it should be noted that in order to improve the security of the platform a correct organisation needs to be chosen to provide this functionality. If the *proxy* colludes with the *sender* then collectively they would be able to reveal the identities of the interesting records. This would be the worst case scenario for IDAP, but the current procedures reveal this information up-front. Thus, both parties would still be better off using IDAP with a proxy than the current system, even if the *proxy* colludes with one of the parties. The additional benefit of this approach is the limited amount of storage that the public authorities would require to keep the original evidence gathered during the electronic enquiry, as without the proxy the *chooser* would likely be required to keep all the retrieved records until the

any possible court case is resolved, and with the proxy only the interesting records would form the set of the gathered evidence.

The *proxy* would be a single party, thus, a single point of failure in the process. However, the *proxy* would be formed by a distributed system possible with its elements being capable of operating independently in order to provide fault tolerance. Also, just in case the organisation in charge of the *proxy* would become malicious or stop providing the services, then the platform would fall back to direct connections between the *chooser* and the *sender*.

5.4 Implementation

For evaluation purposes only, IDAP was implemented using C# .NET programming language allowing it to run on Microsoft Windows NT 5.0 and higher platforms. However, the design is not specific to a programming language. In fact, as it is possible to implement the symmetric encryption and hashing with no modifications to known cryptographic programming suites and only small modifications are required to implement PH encryption, the system can be deployed using virtually any popular programming language onto any operating system and hardware that supports RSA encryption.

The implementation is performed in order to facilitate the discussion of the platforms performance. Although, it is possible to compare the platform to other PET protocols using their theoretical complexity, this is insufficient to evaluate the platform against the requirements. The running time can be estimated using the average cryptographic operation times provided in Table 4-1, but there is no certainty that the model used to produce these estimates is not too simplistic. There is possibility that the operations considered as negligible, in terms of the required processing, by the model in the implemented solution do affect the performance. Therefore, IDAP is implemented in order to confirm the results achieved from the model. In this way, if the modelled complexity gives a good illustration of the results then it will suffice as a source of data for other experiments.

Section 5.3 has provided the overall design of IDAP, which has highlighted that IDAP is only a tool that should be used within the well-vetted data acquisition

processes already used by the public authorities, and other public authorities. Figure 5-7 depicts the way that IDAP fits into the data acquisition process. In this figure blue arrows signify requests, which between the *chooser* and the *sender* include randomly generated ($m \times o$) identities, while the response from the *sender* is filtered by the *proxy* to provide the *chooser* with m interesting records, only.

Implementation of the encryption protocols has already been discussed in Section 4.7.1. It is possible to implement these with only small changes to common open source cryptographic suites. Most OT and PIR primitives are evaluated on a single machine, since the communicational complexity could be easily modelled if necessary and there is no real need to run processes in parallel during testing. However, in order to establish whether this test methodology is correct, the implementation of IDAP makes it a distributed application. For this reason, data must be transferred between the parties in a way that is most optimal for the large number of records being transferred. In the Microsoft .NET framework, data can generally be transferred in the raw binary format or using XML. The first method is usually preferred for large files, since the raw format is thought to contain less overhead and does not need to be specially encoded for transport. However, the first is true only for certain binary data items such as pictures, videos and sound clips, since even binary files contain metadata describing the file type and structure. The more complex the structure, the more metadata is required. On the other hand, in ordinary XML files, binary data (such as the encrypted records) needs to be encoded, which can result in an unnecessary overhead. As an example, the encoding-overhead ratio is 4:3 if binary data is Base-64 encoded for transport. Finally, Simple Object Access Protocol (SOAP) technology called Message Transmission Optimization Mechanism (MTOM) can be used to embed raw binary data within XML. This combines the benefits of the raw format and XML structure for objects in cases where the size of records is higher than 1kB [127]. Consequently, MTOM is employed in data exchanges between the parties.

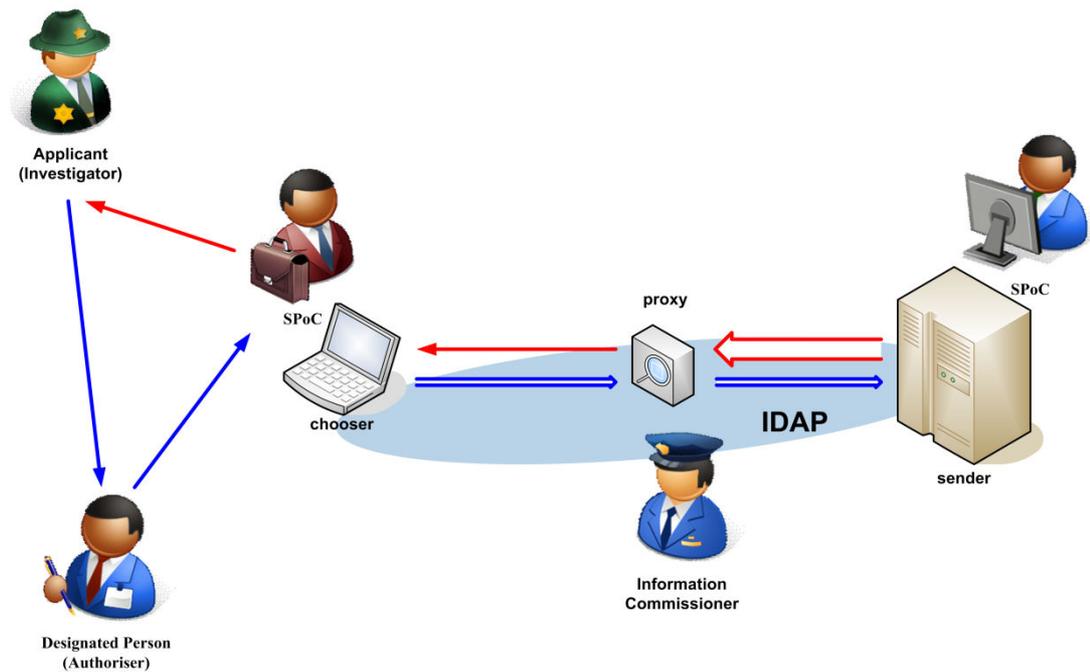


Figure 5-7 IDAP

Common programming practice is to receive the whole piece of data into a buffer, when transferring it over a network. This technique, although perfectly valid, is not suitable for larger pieces of data and streaming should be used instead for large chunks of data. However, in case of IDAP, the individual pieces of data are small in size but large in numbers so streaming is not necessary, and the data should be sent across the network in relatively small messages, as to enable buffering of these messages [127].

The control messaging between the three different players in the system (*chooser*, *proxy* and *sender*) is handled by the Windows Communication Foundation where the *chooser* is the client of services, provided by the *proxy*, and the *proxy* is the client of the *sender*.

5.5 Proposed Quantitative Evaluation

The quantitative evaluation of IDAP assesses the validity of the customisations to the PE primitive discussed in Section 5.3. Results from the evaluation are provided and discussed in Chapter 6.

5.5.1 Overall design of experimental environment

IDAP consists of three applications: chooser (client); proxy (client and server); and sender (server). The performance of each of these applications needs to be evaluated, however, it should be noted that due to the way the platform has been designed the proxy does not need to process records, as it simply relays messages between the choosers and the senders, and filters the results returned from data acquisition queries. The major metrics that can be used to evaluate the protocols are

- Processing time per operation.
- Bandwidth used.

However, in order to establish the strain that IDAP puts on hardware a number of secondary metrics needs to be collected during the experiments:

- CPU usage.
- Memory usage.

The method used to collect these data cannot be resource intensive, as not to affect the results. For this reason, the processing time is measured by taking a timestamp before and after an operation, and the time frames for each operation are calculated and stored into results database at the end of the protocol run. An example, shown in Figure 5-8, demonstrates measuring processing time for symmetric encryption performance test.

```
byte[] messageBytes, outputBytes;
Hashtable encrypted = new Hashtable();
start = DateTime.Now;
for (int i = 0; i < count; i++)
{
    messageBytes = (byte[])Encoding.Unicode.GetBytes(input);
    outputBytes = symHlpr.performEncrypt((byte[])keys[i], messageBytes);
    encrypted.Add(i, outputBytes);
}
end = DateTime.Now;
step_2 = Convert.ToInt32(((TimeSpan)(end-start)).TotalMilliseconds);
```

Figure 5-8 Measuring processing time

```

string networkCard = "";
static PerformanceCounter dataSent;
static PerformanceCounter dataReceived;
ArrayList bandwidthSamples = new ArrayList();

public static KeyValuePair<float, float> GetNetworkUtilizationNow()
{
    return new KeyValuePair<float, float>(dataSent.NextValue(),
        dataReceived.NextValue());
}

private void initialiseNetCounters()
{
    dataSent = new PerformanceCounter("Network Interface", "Bytes Sent/sec",
        networkCard);
    dataSent.NextValue();
    dataReceived = new PerformanceCounter("Network Interface", "Bytes Received/sec",
        networkCard);
    dataReceived.NextValue();
}

public void client_StartFullEnquiryCompleted(object sender,
    StartFullEnquiryCompletedEventArgs e)
{
    ...
    initialiseNetCounters();
    bandwidthSamples.Add(GetNetworkUtilizationNow());
    ...
}

private void decryptData()
{
    ...
    KeyValuePair<float, float> full_stop_band = GetNetworkUtilizationNow();
    bandwidthSamples.Add(full_stop_band);
    KeyValuePair<float, float> full_start_band = (KeyValuePair<float, float>)
        bandwidthSamples[0];
    float totalSent = (full_stop_band.Key - full_start_band.Key);
    float totalReceived = (full_stop_band.Value - full_start_band.Value);
    ...
}

```

Figure 5-9 Measuring bandwidth used during a protocol run

Bandwidth used, memory and CPU load are measured in a similar manner by an external application. Prior to the system entering a given stage in the program the

performance measurement subroutine is run that reads the number of the bytes transmitted and received on a given networking adapter (Figure 5-9), then the probing of the memory and CPU usage takes place. Once the given stage of the IDAP program is over, the subroutine is terminated, statistics for the network adapter read again, the used bandwidth, as well as average memory and CPU usage calculated.

5.5.2 Experiments

The experiments forming the quantitative evaluation of IDAP mainly consist of performance measurements in different scenarios. However, the initial experiment assesses whether modelled complexity of IDAP is a suitable source of data for evaluation. This is achieved by comparing the results obtained from a model of IDAP with measurements taken from scenarios run on the implementation of IDAP.

Evaluation of the complexity model

The main objective of this experiment is to check how closely modelled complexity matches performance of the implementation, which should reflect on the programming language and the run-time environment used for the implementation. Thus, for simplicity, the experiment is conducted using PE implementation and not IDAP. Both are made with the same components, so results achieved will apply to IDAP, just as well.

Varying the number of interesting records

Some proposed modification to PE set out to improve significantly the processing time required for each protocol run. This simulation run on PE and IDAP is used to compare the performance of these protocols for different queries with varied number of interesting records requested. It has already been shown that processing time for PE is almost constant for increasing number of interesting records. On the other hand, in the graph outlining the performance of the PE in similar simulation (Figure 4-7) it can be seen that for high number of interesting records m the processing time is no longer linear, and grows almost exponentially. This simulation will be run a number of times for different values of n in order to evaluate the maximum limit of interesting records per enquiry.

Varying number of enquiries

In IDAP there can be multiple parties acting as *choosers* and *senders*, but logically there can be only one *proxy*. Physically, the *proxy* can consist of a number of devices, but only one party can be put in the position of trust to oversee the data acquisition requests. Otherwise, it would be possible for a *chooser* to cheat and make a number of similar requests to the same *sender* via different proxies, potentially allowing the *chooser* to *fish for evidence*. Taking into consideration that there are many different public authorities that may require investigative data from third parties, and the major CSPs are the most likely targets of the DPA and RIPA notices the load on the *proxy* and the *sender* need to be evaluated for growing number of simultaneous enquires.

Evaluation of IDAP with access to a directory of identities

Many efficient SPIR solutions exist for scenarios where directory of records is public. For example, such SPIR techniques are proposed for privacy-preserving purchases [72]. While on the first sight it is unlikely for the investigators to have a list of the identities in the dataset, this is not always the case. The Police has already direct access to subscriber data of large CSPs, and HMRC knows the names (and other details) of employees in an organisation. In this simulation performance of the IDAP using a dictionary of the dataset is going to be measured against performance of ordinary IDAP that runs the PEqT protocol against the complete dataset. This is going to be measured for varying number of enquiries and varying number of interesting records.

5.5.3 Proposed Qualitative Evaluation

It is possible to quantitatively evaluate the performance of IDAP; however, this thesis sets out to provide a complete solution to the problem of privacy in data acquisition. A problem that is not purely technical. An integral part of privacy evaluation is an assessment of perception of a given system. This is because privacy is different for each and every individual, and the legislations in UK and other European countries enforce this, by giving all individuals certain level of control over their data. For this reason, qualitative evaluation is necessary for IDAP. This major evaluation is carried out in a form of survey targeted at security and privacy experts. This evaluation is complemented by a discussion of IDAP's security.

Survey

A website [128] has been set-up to host the survey. On the survey website the area of the research is explained and brief explanation of IDAP provided. After collecting some details relating to the participants' interest in the subject, the introductory questions ask the participants about their (and their organisation's) security practices, such as the use of secure communication and storage, and their attitude towards digitalisation of the data acquisition process. The remaining questions introduce the privacy problems in this process and propose solution, requesting quantitative and qualitative feedback from the participants. Thus, each single- and multiple-choice question is followed by a text box allowing the participants to express their answers, and opinions, in less rigid manner. Please see Appendix C for to see the survey questions.

Correctness and Security

Security of IDAP is verified against the requirements of MPC [78] and in contrast to the cryptanalytic attacks outlined in Section 2.5.5. The PE primitive has already a good proof of security, and this will be presented here based on [60], while the discussion will focus on the effect that changes to PE introduced by IDAP influence security.

5.6 Conclusion

Building on the results presented in Chapter 4, where PE has been chosen as most likely primitive to facilitate privacy-preserving data acquisition process, necessary improvements to this primitive are proposed. It has been identified that despite PE being most likely the efficient protocol that can privately match and retrieve investigative data. However, retrieving records from large databases with (15 million and more records) is not viable. Consequently, a technique for narrowing down the scope of data retrieval is designed and implemented for evaluation.

Generally, when using privacy-preserving information retrieval techniques, it is only possible to match records on a single selection criterion, such as a record ID. In this thesis it is proposed that a number of different selection criteria can be combined together by linking the sought-after values in one string and hashing this string for use in the same way as a record ID. In this way, it is possible to privately search

databases with minimal complexity increase. However, this technique may place considerable load on the database, thus an experiment is proposed to evaluate this load.

The third, and final improvement potentially needed by the PE primitive to become viable for data acquisition process is addition of semi-trusted third party. This thesis, and the whole field of MPC, is based on the real-life assumption that trusted third party does not exist in most cases. Looking at the currently deployed privacy-preserving technologies procedures are more trusted than encryption / technical measures. Thus, a semi-trusted model where a party proxies the communication between the investigators and the dataholder, and ensures that the investigators do not receive any other data records than those specified as interesting, can benefit public's trust. Such party, called *proxy*, would have to be mutually trusted by the public authorities, the dataholders and the data subjects. In UK ICO could potentially become the *proxy* for all data acquisition requests,

The improved PE primitive forms IDAP and then suitable evaluation techniques are discussed. This includes quantitative evaluation aiming to establish whether the performance of the IDAP is satisfactory, as well as qualitative evaluation seeking feedback from security and privacy specialists, and grading IDAP's security based on available literature.

Chapter 6

Evaluation

6.1 Introduction

The PE protocol was selected as the most suitable basis for an investigative data acquisition protocol. However, Chapter 4 identified a number of shortcomings of this protocol in an investigative scenario that have been analysed in Chapter 5. This resulted in definition of a novel approach to data acquisition referred to as IDAP. This approach needs to be evaluated in terms of performance and correctness.

IDAP's performance is, therefore, evaluated against the PE protocol. But first, the methodology of using simulations for assessing the performance of this type of protocol is put to test in an experiment where empirically-gathered performance data is compared to a simulation. This confirms the methodology of modelling complexity of a protocol in order to evaluate it, commonly used in this field, is correct. Finally, the conducted simulations clearly outline the benefits of using IDAP over PE.

Survey results are discussed and suggest that hiding the identities of alleged suspects is the correct solution to the problems of privacy and secrecy in the investigative data acquisition process. This shows that the main benefit of IDAP over the current processes is in-line with the expectations of privacy and security practitioners.

6.2 Presentation of performance impact

Performance of IDAP is put to test using experimentation and simulation. It is shown that IDAP performs better than PE under most circumstances.

6.2.1 Evaluation of the complexity model

The main purpose of this experiment is to establish whether it is possible to evaluate IDAP based on modelled complexity, rather than experimental results taken from the implementation of the complete platform. The literature suggests using the notion of complexity to compare different protocols, however, most often only the most costly operation is included in the consideration of computational complexity. Thus processes such as symmetric encryption and hashing are often ignored, and a number of asymmetric encryption operations is used to express computational complexity [78]. Often in a similar fashion, communicational complexity is often expressed as a function dependant mainly on number of records processed. In this thesis an analogous, but more precise technique, is employed to quantify computational complexity. Chapter 4 provided a complexity table for the PE protocol outlining the number of cryptographic operations required for the protocol run. These outlined operations included symmetric and asymmetric (commutative) cryptographic functions. Also, the cost of each operation evaluated empirically is expressed in milliseconds per operation. Such tables can potentially be used to plot performance graphs for the different uses of the protocols. However, some empirical testing needs to be performed in order to ensure that this simulation technique is fit-for-purpose. Thus, in this experiment the modelled complexity of PE is compared to results obtained from the C# .NET implementation of this protocol (outlined in Appendix A). Figure 6-1 illustrates this comparison in a simple scenario with 1,000 records in the *sender's* database and a varying number of interesting records. The three lines in this graph represent:

- (1) The predicted total time required for a protocol run, derived from a summary computation complexity for the *chooser* and the *server*.
- (2) The total time needed by the implementation to provide results. Measured from the callback *client_StartFullEnquiryCompleted* till the end of processing the records by the *decryptData* method (Appendix A).
- (3) The summary of the time required by the *chooser* and the *server* to compute and to exchange the results. This is the sum of the run times for sections C1, C3 and S1 (Appendix A)



Figure 6-1 Complexity table reading vs. actual measurements

There is no direct link between the total time taken per enquiry and the total time achieved using the complexity tables. The run times are of the same magnitude, but the line illustrating the simulated run-time curves up with a lower count on the number of interesting records. This can be explained by the process flow of the implementation. By default the *chooser* performs calculations on the IDs of the interesting records, while the server prepares the dataset, by creating packages for each record according to the query received from the *chooser*. Thus, some operations may run in parallel. However, this fact is not taken into account in the complexity table presented in Table 4-3 and used for this experiment. So, in fact, the calculated total run-time should be compared to the sum of time required by the *chooser* and the *sender*. In the graph, the curve of the line illustrating the measured run-time matches closely the curve plotted based on the simulated results. Therefore, the complexity tables do reflect closely on the performance of the implemented protocol. However,

the fact that some operations are conducted in parallel by two separate parties needs to be taken into account when discussing IDAP's performance.

Taking into consideration these results, the complexity table for IDAP should allow for different ways to assess the processing times. Consequently, these should include separate definition of complexity for the *chooser* and for the *sender*. Table 6-1 provides a detailed breakdown of IDAP's computational complexity. It includes the parameter k that expresses the number of enquiries subsequent to the periodical encryption of all IDs in the system performed by the *sender* (Phase A).

		Symmetric Cryptography		Asymmetric Cryptography		
		encryption	decryption	key gen.	encryption	decryption
Phase A (run periodically)	Step 2	-	-	$O(1)$	-	-
	Step 3	-	-	-	$O(n)$	-
Phase B (run per enquiry)	Step 3	-	-	$O(1)$	-	-
	Step 4	-	-	-	$O(m)$	-
	Step 6	-	-	-	$O(m)$	-
	Step 8	-	-	-	-	$O(m)$
Phase C (run per enquiry)	Step 3	-	-	$O(1)$	-	-
	Step 4(a)	-	-	-	$O(m \times o)$	-
	Step 4(b)	$O(m \times o)$	-	-	-	-
	Step 6	-	-	-	$O(m)$	-
	Step 8	-	-	-	-	$O(m)$
	Step 9	-	$O(m)$	-	-	-
Total Complexity		$O(k \times m \times o)$	$O(k \times m)$	$O(2k + 1)$	$O(km(o + 3) + n)$	$O(2 \times k \times m)$

Table 6-1 Initial definition of IDAP's complexity.

This complexity table can be further refined. Symmetric encryption and decryption times are almost identical, 551 μ s and 564 μ s respectively, so the complexity of the operations can be summarised and the cost rounded-up to 0.6ms/operation. In IDAP there is no need for a large number of asymmetric cryptographic keys, as for a single run of protocol there are only three keys required, which means that asymmetric key generation should not affect the protocol run times. On the other hand, asymmetric encryption contributes the most towards the computational complexity. This can be observed in the Figure 6-2, as the processing time for asymmetric encryption is few magnitudes higher than any other component of the complexity table.

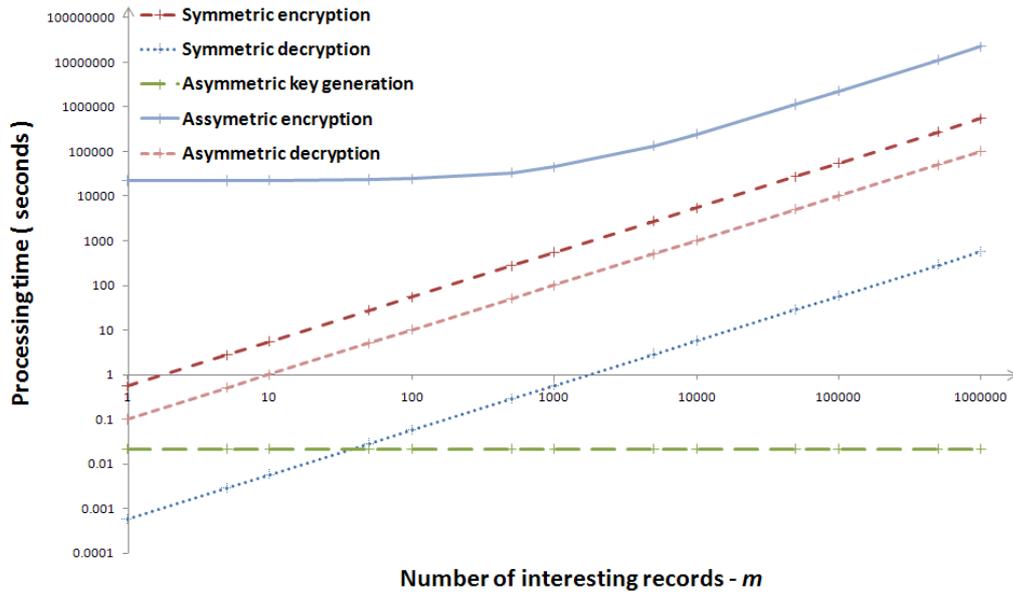


Figure 6-2 IDAP complexity in detail for varied m ; $k=1$; $o=1,000$ and $n=1,000,000$

The presented complexity table, and the resulting graph, are imprecise, as they do not reflect all conditions. First, and foremost, if the product of the number of interesting records, and the dilution factor, is higher than the number of records in the dataset (this means for the cases where condition from Eqn. 6-1 is not met) the pure version of the PE protocol, as presented in [60], needs to be employed:

$$m \times o \leq n \quad \text{Eqn. 6-1}$$

But such scenarios should be avoided in practice, as they cannot guarantee k -anonymity of the alleged suspects. Nevertheless, in order to evaluate IDAP's processing, the complexity of PE is used for such cases. This is reflected in Table 6-2.

	Condition	Symmetric Crypto. operation	Asymmetric Crypto		
			Key gen.	Encryption	Decryption
Total complexity	$O(m \times o \leq n)$	$O(k \times m \times (o + 1))$	$O(2k + 1)$	$O(km(o + 3) + n)$	$O(2 \times k \times m)$
	$O(m \times o > n)$	$O(k \times (n + m))$	$O(2k + 1)$	$O(k(3m + n) + n)$	$O(2 \times k \times m)$
Chooser's complexity	-	$O(k \times m)$	$O(k)$	$O(km)$	$O(2km)$
Sender's complexity	$O(m \times o \leq n)$	$O(k \times m \times o)$	$O(k + 1)$	$O(km(o + 2) + n)$	-
	$O(m \times o > n)$	$O(k \times n)$	$O(k + 1)$	$O(k(m + n) + n)$	-

Table 6-2 IDAP's complexity.

The second omission from the complexity show in Table 6-1 and Table 6-2 is the scenario where the *chooser* can obtain a directory of the entries in the dataset. For

example the Police have direct access to the subscriber data of large CSPs in UK, thus, they can use this access to narrow down the scope of the enquiry so that the complete dataset of a large CSP does not need to be encrypted. Also, the CSPs make public the ranges of the identifiers that they manage, such as IP addresses, or telephone numbers. Thus, with a certain probability, it is possible to generate an enquiry that ensures there are at least o different active identities (addresses) that can be linked to every interesting record. The computational complexity for this variation of the protocol is shown in Table 6-3.

	Condition	Symmetric Crypto. operation	Asymmetric Crypto		
			Key gen.	Encryption	Decryption
Total complexity	$O(m \times o \leq n)$	$O(k \times m \times (o+1))$	$O(3k)$	$O(km(2o+3))$	$O(2 \times k \times m)$
	$O(m \times o > n)$	$O(k \times (n+m))$	$O(3k)$	$O(k(3m+2n))$	$O(2 \times k \times m)$
Chooser's complexity	-	$O(k \times m)$	$O(k)$	$O(km)$	$O(2km)$
Sender's complexity	$O(m \times o \leq n)$	$O(k \times m \times o)$	$O(2 \times k)$	$O(km(2o+2))$	-
	$O(m \times o > n)$	$O(k \times n)$	$O(2 \times k)$	$O(k(m+n)+n)$	-

Table 6-3 IDAP's complexity for datasets with publically available dictionaries.

6.2.2 Varying the number of interesting records

The purpose of this test is to evaluate IDAP's performance in contrast to PE for varying number of interesting records. By customising PE to a specific application of investigative data acquisition IDAP has introduced a number of potential improvements to the efficiency of this information retrieval protocol.

Figure 6-3 illustrates processing time for IDAP and PE for varying number of interesting records (Eqn. 6-2). It can be seen that initial processing time for IDAP is slightly lower than the PE processing time, and, as designed, both are the same in the range where the product of m and the dilution factor o is higher than the size of the dataset n (Eqn. 6-3).

$$1 < m < n \quad \text{Eqn. 6-2}$$

$$m \times o > n \quad \text{Eqn. 6-3}$$

The graph shown in Figure 6-4 focuses on the range of values that do not meet condition in Eqn. 6-3. IDAP performs significantly better than PE for low values of m , or, to be more precise, when ratio of m to n is small. This experiment needs to be

repeated for different values of n and m in order to establish the optimal value of the ratio of $m:n$, and the effect of the dilution factor on the optimal value of this ratio. The graphs that help establish this connection are presented in Figure 6-5 and Figure 6-6. Figure 6-5 shows that IDAP's processing time is almost constant for low values of m if other parameters (n and o) are constant. The point where the line showing the processing time curves up represents the optimal operation of IDAP. However, Figure 6-6 shows that the processing time is also dependant in a similar, but inverse, way on o the dilution factor. If the proportion of the three key parameters (m , n , and o) is referred to γ it can be expressed as in Eqn. 6-4 (this follows from Eqn. 6-3). Thus, since the processing time is constant for $m \leq 100$, $o = 1000$, $n = 1000000$, and so on, the operation of IDAP is optimal for $\gamma = 0.1$ (Eqn. 6-5), if there is no overlap between the diluting records.

$$\gamma = \frac{m \times o}{n} \tag{Eqn. 6-4}$$

$$\gamma = \frac{100 \times 1000}{1000000} = \frac{1}{10} \tag{Eqn. 6-5}$$

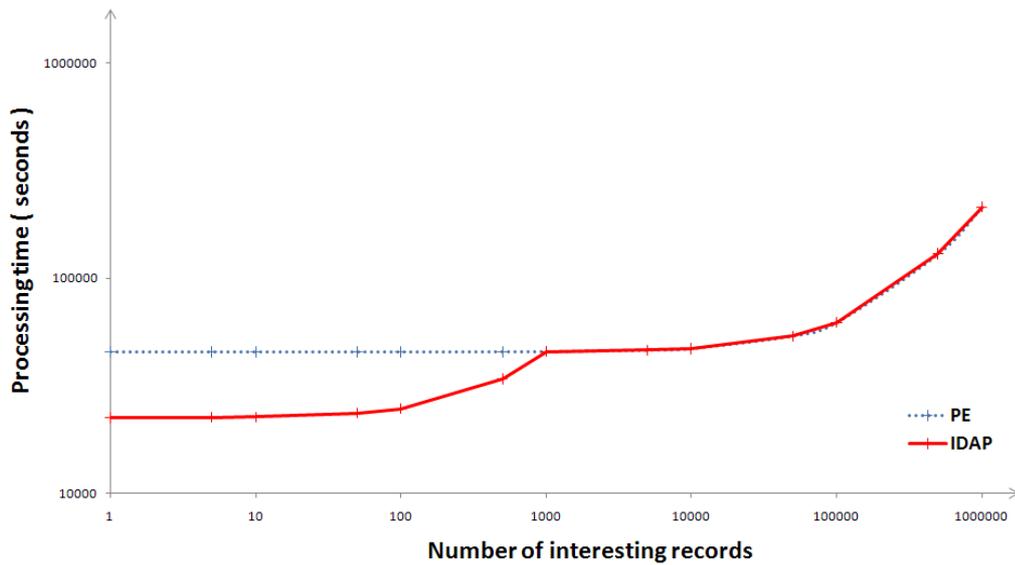


Figure 6-3 Computational complexity of IDAP and PE for increasing m (logarithmic scale)

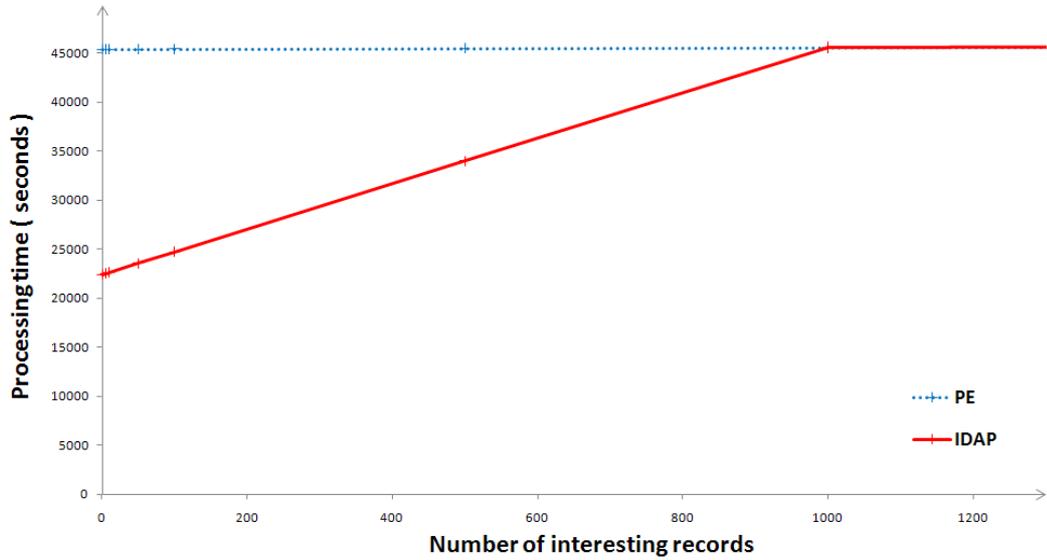


Figure 6-4 Computational complexity of IDAP and PE for increasing m

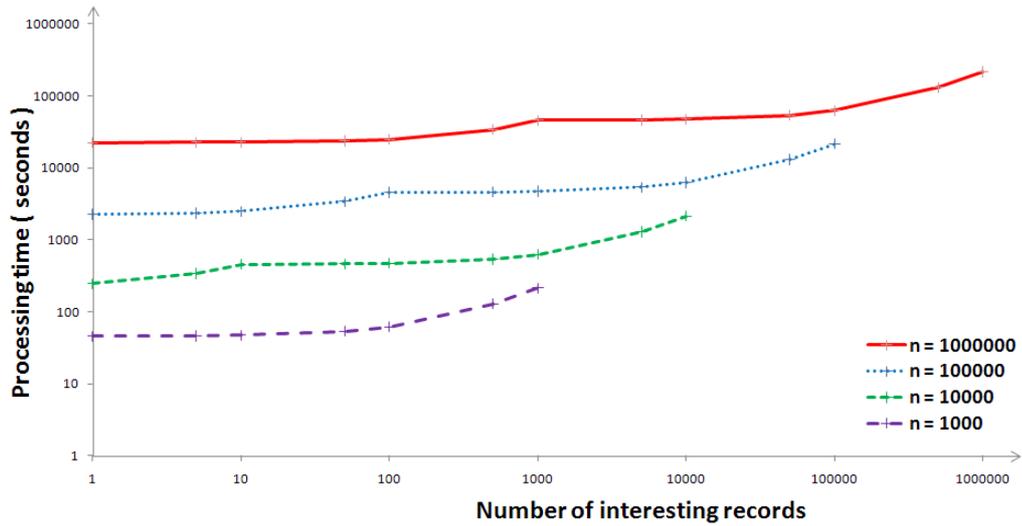


Figure 6-5 IDAP's processing time for varying m and different values of n

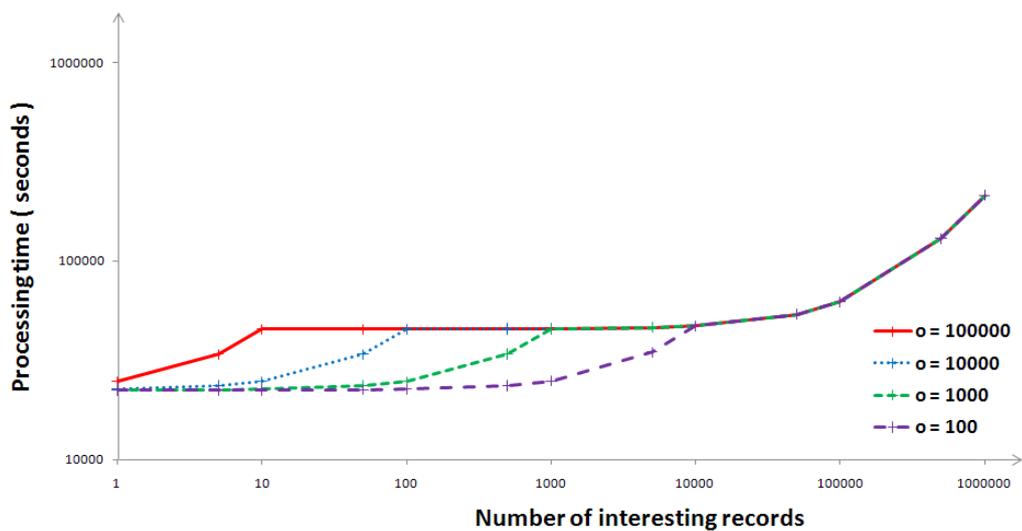
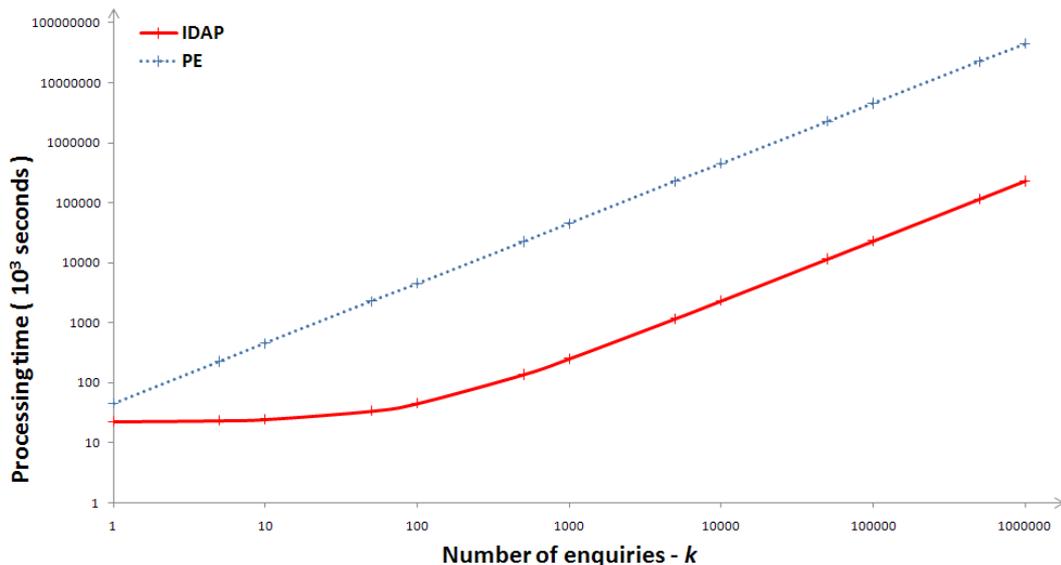


Figure 6-6 IDAP's processing time for varying m and different values of o

6.2.3 Varying number of enquiries

This simulation focuses on the evaluation of IDAP, in contrast to PE, for scenarios with multiple enquiries against the same dataset. In IDAP k different enquiries can be made against a single encrypted list of identities within a dataset. This would be typical scenario for organisations that often provide data to the public authorities, such as CSPs. Figure 6-7 and Figure 6-8 depicts the results of this experiment. The key point to take from these graphs is that under the test conditions ($n = 1\text{mln}$; $m = 10$; $o = 1000$) the processing time for PE is two orders of magnitude higher than IDAP's processing time. In addition, IDAP's processing time is almost constant for $k \leq 10$.

The value of γ affects the results and therefore the effect that the ratio defined in Eqn. 6-4 has on these results is further evaluated in Figure 6-9. The maximum value for γ is one in environment where diluting records do not overlap, but even for $\gamma = 1$ IDAP performs significantly better than PE in a data acquisition scenario.



**Figure 6-7 Comparison of processing time for IDAP and PE.
For varying k , $n=1\text{million}$, $m=10$, and $o=1000$**

As expected the lower the value of γ , the lower the processing time is for the same parameters n and m . Consequently, looking at the performance alone, there is no limit for the number of enquiries that can be run following a single encryption of the identities in the system by the *sender* (Phase A of IDAP). But there are possibly security limitations to this procedure and the maximum number of subsequent

enquiries k that can be run out of a single encryption of a directory. These will be discussed later in this chapter (Section 6.3.1).

The graph shown in Figure 6-9, also illustrates that when condition from Eqn. 6-6 are met, the processing time is almost constant. This knowledge can be used to maximise the return-on-investment, thus fine-tuning the platform to use the maximum computational capacity of the protocol. Consequently, to provide maximum number of records with minimum effort:

$$\gamma \times k \leq 0.1$$

Eqn. 6-6

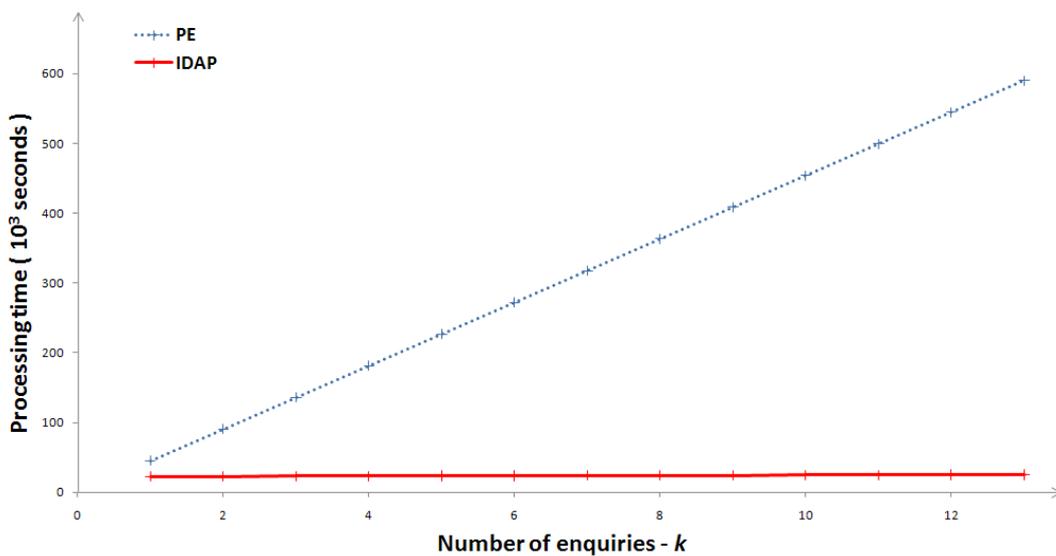


Figure 6-8 Detailed comparison for IDAP and PE with. For varying k , $n=1$ million, and $m=10$

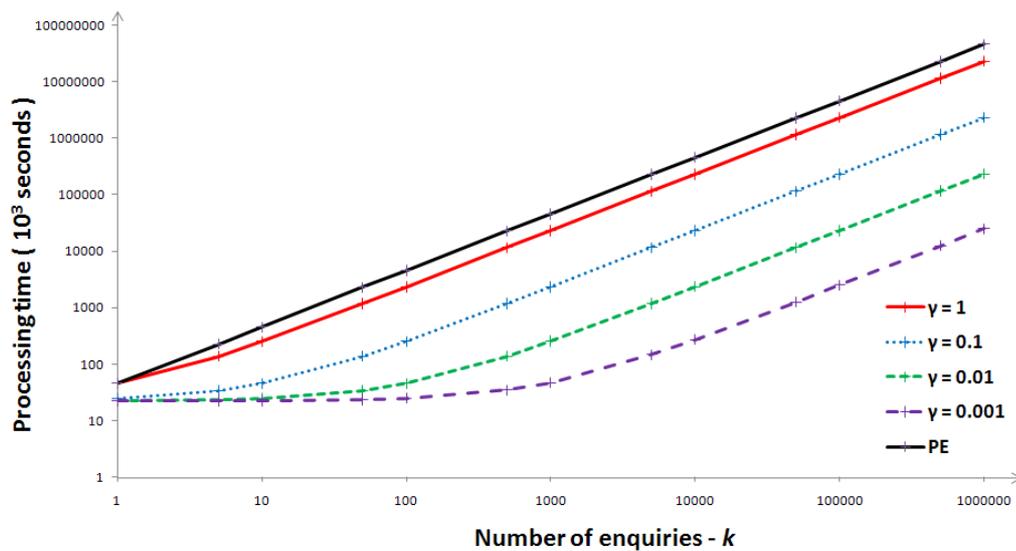


Figure 6-9 IDAP's performance for different values of γ , as compared to PE. For varied o , $n = 1$ million, and $m = 10$

6.2.4 Evaluation of IDAP with directory of identities

In theory it should be possible to improve IDAP's performance in scenarios where directory of identities in the dataset is available. In such cases there is no need to encrypt the whole dataset with a commutative encryption scheme. Figure 6-2 shows that commutative encryption is the major factor in the total processing time, thus reducing the number of records that need to be encrypted is likely to lower the processing time.

The results of this experiment are shown in Figure 6-10 and Figure 6-11. As expected, with access to a directory of identities, IDAP performs significantly better for $\gamma < 1$ (Figure 6-10). This is for cases where the product of the number of interesting records, and dilution factor, is smaller than the number of records in the dataset. However, Figure 6-11 depicts that IDAP using a directory is only more efficient than ordinary IDAP for $k < 10$ under the test conditions ($n = 1 \text{ million}; m = 100; \text{ and } o=1000$). To be more precise, this occurs if condition from Eqn. 6-7 is met:

$$\gamma \times k < 1$$

Eqn. 6-7

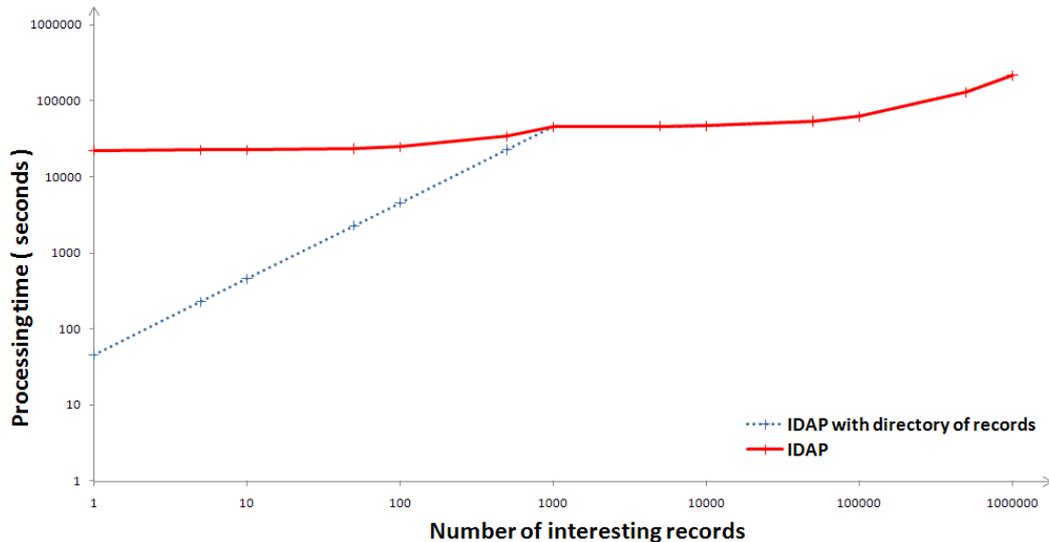


Figure 6-10 Performance gain for IDAP run on dataset with a directory.
Plotted for varying m , $n=1\text{million}$, $o=1000$, and $k=1$

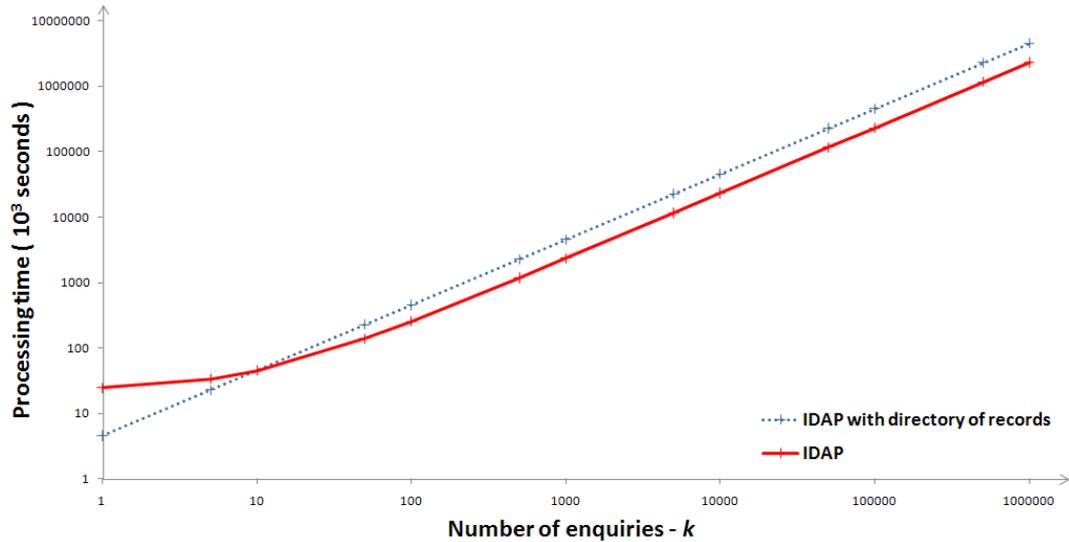


Figure 6-11 Performance gain for IDAP run on dataset with a directory.
Plotted for varying k , n =million, m =100, and o =1000

6.2.5 Use of dilution factor with different protocols

The technique of dilution the enquiries can also be applied to other SPIR protocols, as long as it operates in a single database scenario. As an example the combined approach formed from PEqT and OT protocol presented in Chapter 4 can be trivially modified to perform such enquiries, as presented in [9]. Figure 6-12 depicts a comparison of how IDAP compares to PE- and OT-based approaches described in Chapter 4, as well as the OT-based solution that benefit from the notion of dilution. (Large preparation time of the OT-based approach is omitted here, and for this reason, the results are comparable to those of other SPIR protocols.) Looking at the $\gamma \leq 1$ range it can be seen that the introduction of the dilution factor to the PE protocol in IDAP cut the processing times nearly by half. However, the notion of dilution factor provides even larger benefits to non-commutative solutions, such as this based on the OT protocol.

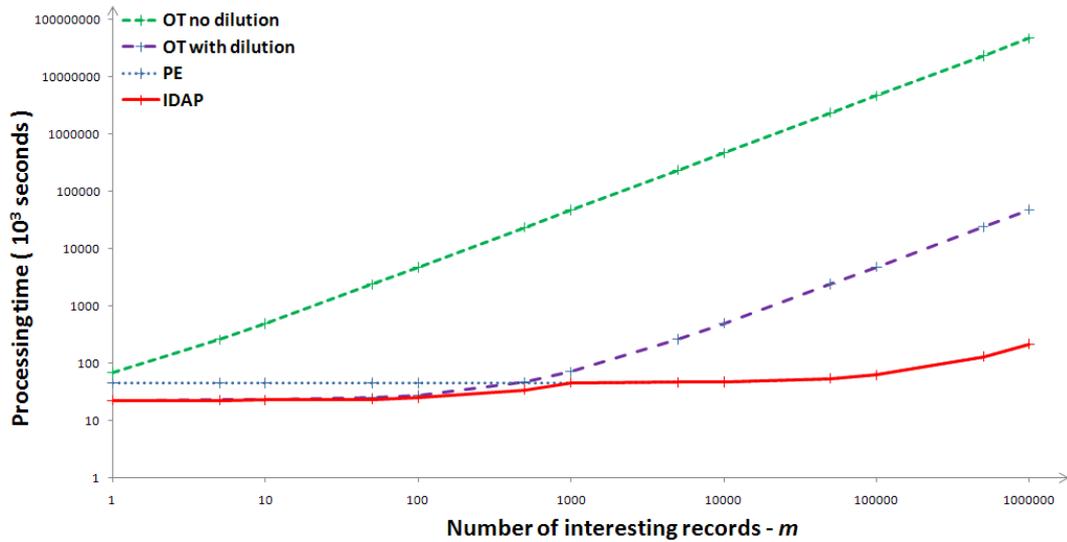


Figure 6-12 Comparison of IDAP to a modification of OT-based approach, that allows for dilution of the enquiry. ($n = 1,000,000$; $o = 1,000$; $k = 1$) [excludes OT preparation times]

6.2.6 Controlling the balance between privacy and feasibility

One of the key objectives of this thesis is to provide a level of control for the balance between privacy and feasibility. The discussion provided in Chapter 4 makes it apparent that even the modern SPIR protocols are unable to retrieve records from large datasets, as their performance is highly dependant of the size of the dataset. For this reason, the notion of the dilution factor has been proposed in this thesis. However, while adding a dilution factor to a protocol such as the one based on simple $1-n$ OT presented in Section 4.5 has a significant effect on performance as shown in Figure 6-12, its effects on the PE protocol are limited. Still, the dilution factor used in PE can half the run times for a single enquiry. However, it is clear from Figure 6-6 that the dilution factor has a direct effect on the performance. The higher the dilution factor, the lower the performance. Simply put with a higher value of dilution factor there is more records to process per enquiry, while the processing time is almost linear to the number of records processed (as shown in Figure 6-2). At the same time we know that higher values of the dilution factor o carry greater protection of privacy, since the identifier interesting of the interesting records are hidden in a greater pool of identities. Therefore, it is possible to control the balance between the privacy and feasibility using the dilution factor. This becomes even more apparent in Figure 6-13 depicting the run times of IDAP with directory against the size of the dilution factor.

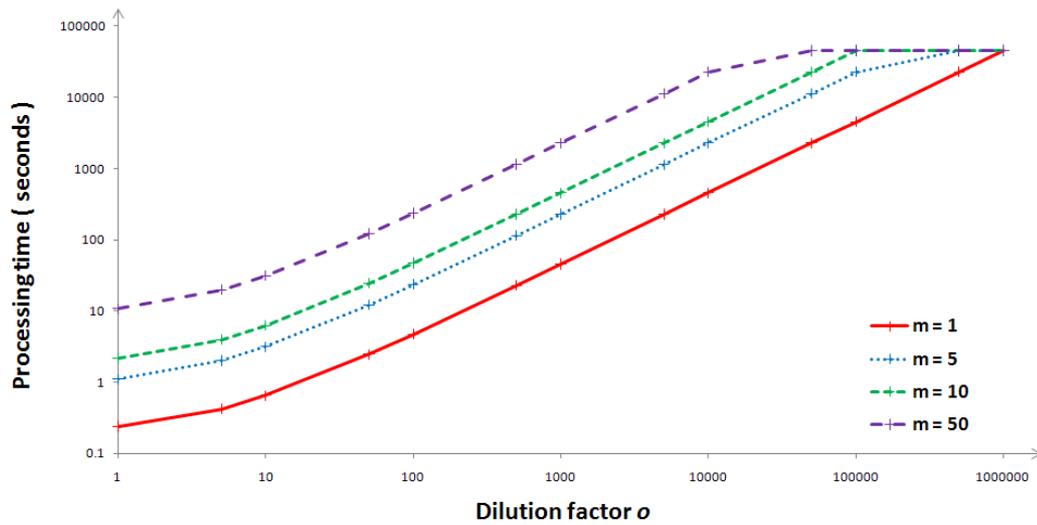


Figure 6-13 Processing time against the dilution factor σ , for IDAP with directory of records. ($n = 1,000,000; k = 1$)

6.3 Presentation of qualitative evaluation

The previous section evaluated the performance gain in IDAP as compared to other PET protocols, such as PE. However, this would be irrelevant if IDAP would not enhance the investigative data acquisition process. In order to establish the benefits that IDAP brings into this process, a survey among of privacy and security professionals has been conducted, and security evaluation is carried out.

6.3.1 Correctness and Security

IDAP is a modification of the PE protocol that has its correctness and security proofs provided in [60]. The goals and logic of IDAP and PE are similar, however, IDAP is streamlined to provide better performance than PE in the specific use scenario of investigative data acquisition. There is an assumption that there is a method of authenticating other parties and securing the channel for communication. In order to evaluate the correctness and security of IDAP the inputs and outputs need to be clearly stated [78], thus these are:

Chooser's input: set V_C containing IDs of interesting records.

Sender's input: set V_S containing IDs of the records in the dataset, together with extra information about these records – $ext(v)$.

Output: *chooser* learns $|V_S|$ (the size of the set V_S), $V_C \cap V_S$, and $ext(v)$ for $v \in V_C \cap V_S$, while *sender* learns $|V_C|$. *Proxy* learns only the sizes of the sets.

Normally both parties learn the sizes $|V_S|$ and $|V_C|$, as by default all the encrypted identities in V_S are sent to the *chooser*, while the *chooser* in order to find the interesting records among these encrypted identities and in order to decrypt the $ext(v)$ for these records provides the *sender* with encrypted elements of the set V_C . There is no requirement by the public authorities to know the size of the dataset, but since there is now a way to run IDAP and avoid providing the authorities with the dataset size, this needs to be accepted as an outcome of the protocol. The fact that the *sender* learns the number of interesting records is beneficial in the data acquisition scenario, as the *sender* can then verify that the *chooser* follows the data acquisition notice that would previously outline the IDs of the interesting records, and under IDAP would specify the number of the interesting records.

IDAP is based on *Shamir's commutative protocols*, a variant of PH protocol where the prime p is public and common between the communicating parties. An adversary with the knowledge of the ciphertext C and the prime p would need to solve the following hard problem to break the commutative PH protocol [41]:

$$e = \log_p C \text{ mod } p \tag{Eqn. 6-8}$$

Just like RSA, the ciphertext created using the PH algorithm may leak some information about the input plaintext message. Therefore, this algorithm is suitable for uses where the input is formed from random data. This is the case in the PE and IDAP, as the commutative PH is used to encrypt hashed IDs of the records. While it is normally recommended to use padding schemes in any implementation of RSA[113], and thus PH implementation as well, the PE and IDAP mitigate this requirement by using fixed size hashes as the input.

The proofs of the correctness and security of PE can be found in [60], while IDAP has modified this protocol by introducing the following improvements:

- Lowering processing time, by narrowing the scope of the enquiry.
- Allowing for multiple selection criteria.
- Restoring the balance between the privacy of the innocent and the suspects.

In order to narrow down the scope of the enquiry IDAP splits the PE protocol into three parts. However, the only way the operations of the protocol are affected is the fact that under IDAP the *chooser* request extra information for only $(m \times o)$ records, rather than for the whole dataset n . The main consequence of this approach in respect of security of the protocol is that the *sender* knows that there are m interesting suspects in the set of identities the size of $(m \times o)$. This could become an issue if the same request is run against a number of parties and the parties collude, but this thesis has shown that the investigative data acquisition process can be treated as a single database scenario, if requests are made against CSPs. So colluding is not possible. On the other hand, for small organisations with less than 100,000 IDs, there is no need to narrow down the results. Consequently, in IDAP, the privacy of the suspect is affected by the dilution factor o , and the *sender's* probability of guessing the interesting records IDs is $1:o$ and not $1:n$. As long as o is reasonably large, and the *sender* has no other sources of information about the suspects, the privacy of the suspects should not be unaffected.

IDAP allows for the multiple selection criteria by hashing together different selection criteria and using it within the PE protocol as an ID of a record. This does not affect the security of the PE protocol. On the other hand adding the semi-trusted third party – the *proxy* – in order to restore the balance between the privacy of the innocents and the suspects somewhat modifies the security of the protocol. The *proxy* filters out the records not classified as interesting from the *sender's* response. Assuming that the semi-trusted party behaves as expected, the security of the $ext(v)$, the data records contained in the *sender's* database is information theoretic from the *chooser's* perspective. On the other hand, if the *proxy* and the *sender* cooperate, they can easily work out the identities of the interesting records. The main aim of IDAP is to hide those identities from the *sender*, however, currently, the identities of the suspects are provided in every data acquisition notice. Consequently, if the semi-trusted party cooperates with the *sender*, the only result that would reveal the information is

currently openly communicated to the dataholders. Still, it is important that the semi-trusted party is chosen so it does have an incentive in maintaining the privacy of the investigations and the data subjects.

6.3.2 Survey Results

A survey has been carried out for the purpose of qualitatively evaluating IDAP. Due to the nature of the subject, the survey was aimed at specific security and privacy professionals, as well as the law enforcement professionals. This means that the responses are from the practitioners that would likely be involved during possible roll-out of IDAP. The graphs illustrating the results can be found in Appendix C.

According to the results, the participants are aware of the encryption technologies used for storing data and for communications, but are sceptical about the positive impact that introduction of digital technologies could bring to security or privacy during the investigative process. Some respondents suggested that while security techniques can increase the security and privacy in an investigative process, the availability of data can balance-out these benefits, as it will open new avenues to abuse the access.

Most respondents agree that currently a data acquisition request can breach suspect's human and natural rights, by affecting the relationship between the suspect and the *sender*. Also, similar views were shared in respect to the effect that a data acquisition request may have on an investigation. On the other hand, the plans to provide the public authorities with direct access to CSP data were met with a mix of responses. Respondents agree that these plans can allow the maintaining of secrecy of investigations and can provide faster access times to urgently needed data. However, the respondents were unsure what effects the proposed changes will have on the privacy of the individuals-under-investigation. The more verbose responses show worries that the extended availability is likely to cause excessive use the communications data. On the other hand, most of the respondents agree that hiding the identity of the suspect from the provider of the data can protect integrity and security of investigations, as well as the rights of the suspect.

Most respondents would accept IDAP as a solution to the privacy concerns in investigative data acquisition, and would expect their organisations to accept it, as

well. However, the more verbose responses suggest that there are worries about the correctness of the non-trivial implementation of any mathematically sound protocol. These worries are also reflected in the respondent's attitudes towards introduction of the *proxy*, as the number of individuals that would accept this solution matches the number of individuals that are against it.

6.4 Conclusion

The experimental results obtained from the implementation of the PE protocol show the same trends as the results achieved via simulation of the PE protocol using computational complexity tables. Since, IDAP uses the same cryptographic mechanisms as PE, these results allow for the simulations to be used as the main way of evaluating IDAP's performance in contrast to PE. Other research has used similar, but less precise, approach of evaluating the protocols based on the number of asymmetric operations required [73, 80, 95], as this is the most costly operation in most OT and SPIR protocols. Thus, using simulations to compare these kinds of protocols is a generally accepted practice.

Results achieved from the simulations of IDAP under different conditions show that the performance of the protocol is highly dependant on the three key parameters:

- m – number of interesting records.
- o – dilution factor.
- n – size of the dataset.

To be more precise the performance depends on the ratio of $(m \times o)$ to n . This proportion is defined as γ in this thesis. It is shown that the operation of IDAP is optimal for $\gamma = 0.1$, as, at this point, the processing capacity of the protocol is used to maximum effect without affecting the total time required to process a query. If $\gamma > 1$ then PE should be used for the operation, rather than IDAP, as under such conditions the IDAP protocol would call for more than n records to be encrypted on the *sender's* end of the process. However, such scenarios should generally be avoided, as they cannot guarantee k -anonymity of the interesting records.

According to the simulation results, IDAP is uniquely placed for multiple enquiries that run against the same set of identities. This includes enquiries where IP addresses or telephone numbers are used to match the records. The number of enquiries run against a single set of encrypted IDs is referred to k in this thesis. For low values of k and γ , this is for $k \times \gamma < 0.1$, the processing time is constant. As expected, the lower the γ , the higher the benefits of using IDAP over PE. Additionally, for $k \times \gamma < 1$ IDAP has a smaller processing time than a single enquiry using PE. Thus, as long as γ is small, a number of enquires can be run using IDAP without incurring any considerable costs, where under PE the computation costs would grow in-line with the number of enquiries.

The dilution factor can be used to control the balance between privacy and performance in the system. Generally, from the privacy point of view, this factor should be as large as possible, as it specifies how many records dilute the identity of each interesting record. However, the processing times are directly affected by the size of the dilution factor, and for higher values, the enquiries are often not feasible. On the other hand, it is interesting to see that the dilution factor can be successfully applied to other types of single database SPIR protocols. In fact, the use of the dilution factor benefits non-commutative protocols more, and can enable other SPIR protocols to perform on par with IDAP and PE in m -to- n enquiries.

The survey conducted among of information privacy and security practitioners suggests that hiding the identities of the suspects during the investigative data acquisition process could enhance the privacy of the investigations and protect the human rights of the suspects. Also, most respondents would accept IDAP as a solution to the privacy concerns related to the acquisition process.

Chapter 7

Conclusions and Future Work

7.1 Introduction

This chapter concludes the thesis. It summarises the steps taken to achieve the objectives of this thesis, including the brief digests of the findings of the Literature Review, initial experiments conducted to further explore the field of SPIR protocols and their potential use in investigative scenarios, as well as the design, implementation and evaluation of the novel IDAP. This final chapter also provides solutions to the motivating scenarios defined in Chapter 1.

This thesis has made a number of contributions to knowledge. It defines a platform that can assist making the investigative data acquisition process more ethical. In order to achieve this it introduces SPIR protocols into this process, and shows that thanks to employing the notion of k -anonymity in such protocols it is possible to achieve a satisfactory level of performance when conducting requests from large databases. This approach is analogical to the technique of *dilution* used, on occasion,

by investigators in order to hide object of their enquiry by combining a number of requests for data together. Thus, this thesis introduces a *dilution factor* used to express the level of privacy and secrecy that a given investigation requires by specifying how many records are used to hide a single interesting record. Furthermore, this parameter can be used to control the balance between performance and privacy in an enquiry. In order to ensure fairness in the protocol and to reassure the public that the records used to dilute the enquiry cannot be decrypted by the authorities this thesis proposes the use of an independent party to monitor the enquiries, and filter-out the responses from the dataholders. The thesis also proposes a technique for making complex queries on the datasets with minimal impact on computational complexity. These aspects of the thesis will be now summarised.

Finally, this chapter discussed possible areas of the further work that can develop and evaluate the concepts of SPIR-assisted investigative enquiries, as well as use the findings of this thesis in different areas of knowledge.

7.2 Achievement of objectives

This thesis set out to meet five objectives defined in Chapter 1. These have been achieved and discussed at various points of this manuscript.

Construct a literature review within the PET sphere

The literature review of the PET sphere is provided in Chapter 3. This review focuses on the measures that individuals can use to protect their privacy in information systems, as well as those that organisations can, and possibly should, use to improve the privacy of the data-subjects in operations on personally identifiable data. PETs are often misperceived as computationally expensive methods for obfuscating sensitive data. However, there are various technologies that can be classified as PETs [79]. For example, security measures such as access-control are also vital for protecting privacy of data subjects, and therefore the distinction between ordinary security measures and PETs is disputable. The literature review has identified that it is possible to use a PET protocols referred to as SPIR protocols to retrieve data from almost any relational database, in a private manner [74, 96-98].

A drawback to the use of SPIR protocols is that most require an index, such as the index of a row in a database table, of the interesting record as an input. In many cases a directory of the records in a dataset exist, for example in an on-line purchase of digital goods scenario, the seller can publish a list of products together with their descriptions, and the buyer can use the index of the item of interest to purchase it and retrieve it using SPIR [73, 80]. However, for the datasets that do not have a public directory often a separate PEqT protocol needs to be run in order to find out the index of the interesting record [60, 105]. Another drawback of most SPIR systems is that they are often designed to perform retrieval of a single record from a dataset of n records. So, in the cases where more than one record would need to be retrieved from the same dataset, the run-time of the operation would be linear to the number of interesting records. One of the protocols that mitigate both of these drawbacks is the PE protocol based on commutative cryptography that was originally designed to share extra information about common records between two or more databases [60].

It is worth noting that SPIR protocols are generally based on a trapdoor function, thus, they are often characterised by a high computational complexity. In addition, such protocols often have high communicational requirements. The balance between the communicational and computational complexity, though, can be altered by the use of balancers [90, 92], as well as with the advances in the area of TC and the use of SCOP devices [108].

Define set of requirements that data acquisition process must meet.

In order to provide a PET-based solution to the privacy concerns related to the data acquisition process, it was necessary to define the requirements that such a solution would have to meet. The process of selecting and refining these requirements, mainly based on available literature, is described in Chapter 4:

- Req. 1** Allow for the gathering of multiple suspect records per enquiry, or have low overhead per each additional query run on the database. (Maximum anticipated number of suspects in one enquiry is 150.)
- Req. 2** Keep the data controller in charge of the data. A data record cannot be transferred or made available, to the public authorities, without the data controller's verification of the request.
- Req. 3** Allow for efficient and timely data retrieval. (The protocol run excluding

preparation should take less time than 30 minutes that it currently takes investigators to obtain investigative data in emergencies.)

- Req. 4** Be cost-effective, as the platform will need to be deployed by a variety of organisations.
- Req. 5** Retain an audit trail of the processing performed on the potential evidence.
- Req. 6** Gain acceptance from the general public.
- Req. 7** Handle large datasets (such as datasets with more than 15 million records).
- Req. 8** Provide a technique for multiple selection criteria of interesting records, and allow for fuzzy matching on the selection criteria different than record ID.

Construct a novel methodology for privacy-preserving investigative data acquisition.

In Chapter 5 a platform for gathering investigative data from third parties was proposed. The platform is called IDAP (Investigative Data Acquisition Platform) and it is well matched to meet requirements listed above. The design of IDAP is influenced by the discussions and simulations presented in Chapter 4. Based on the gathered requirements, SCOP-based solutions were considered to be more expensive and harder to deploy than software based equivalents. The literature review identified that most SPIR protocols can be combined with a PEqT protocol to form a system capable of searching the datasets belonging to third parties and retrieving records in a private manner. Such a combined approach was tested against the PE protocol that natively provides this functionality. Results of the simulations described in Chapter 4 have shown that it is unlikely for any 1- n SPIR protocol to perform better than the PE protocol in scenarios with multiple interesting records.

PE is the basis for IDAP, the novel methodology for privacy-preserving investigative data acquisition. In order to address all the requirements for the data acquisition platform, IDAP modifies the way that PE handles requests by adding the following, domain specific, improvements:

- Introduction of the dilution factor that can limit the scope of a request in order to improve performance.
- Introduction of a method for performing low-overhead dataset searches with multiple selection criteria.

- Innovative use of a semi-trusted third party in order to restore the balance between the personal privacy and the secrecy of investigations.

These allow IDAP to become a scalable platform for privacy-preserving retrieval of investigative data from third parties. If IDAP is treated as a tool exchanging data between the relevant SPoCs, it fits well within already established, and well defined, processes for data acquisition, such as the code of practice presented in [8].

Propose an evaluation framework suitable to assess performance of novel cryptographic enhancements to retrieval of investigatory data.

The literature review has found that most researchers use the notion of complexity to evaluate OT- and SPIR-based privacy-preserving protocols [75, 76, 78, 80, 84, 96, 124, 129]. However, often only the number of computationally expensive encryption operations is included in the function describing complexity of a given protocol. In this thesis, a more precise approach is proposed. IDAP is evaluated based on the simulations plotted according to the complexity tables that depict cryptographic operations used per step of the protocol. With this approach it is possible to compile graphs showing IDAPs performance under various conditions.

Investigate parameters that could be used to assess the balance between the privacy and feasibility.

PET technologies are often computationally expensive, and the higher levels of privacy are usually linked to the higher complexity of protocols. Typically, it is difficult to control the levels of privacy offered by SPIR protocols. Typically, 1- n and m - n SPIR protocols aim to provide total anonymity, as the *sender* is unable to distinguish the interesting records from any other records in the dataset. However, the performance of the protocol depends mainly on the size of the dataset, which, in case of data acquisition enquiries, can be large. This thesis has identified that most enquiries require only a single data provider. For example, if the last known location of a suspect needs to be established based on the data from a mobile telephony provider, then only one provider needs to be queried for this data, as it is based on the telephone number, it is possible to identify the right provider of services. Consequently, it is possible to limit the number of records retrieved from the provider to the interesting records, plus a number of records to obfuscate the

identities of these interesting records. Thus, it is possible to provide a constant ratio of interesting records m , to the number of records that is retrieved from the provider o . In this thesis this is referred to as the *dilution factor* and it provides a form of k -anonymity in the investigative scenario. Finally, it is possible to control the balance between the privacy and feasibility using the dilution factor, which is one of the main contributions to knowledge of this thesis. Currently, the dataholders know with 1:1 probability the identity of the interesting records, with 1- n SPIR protocols this would be a 1: n probability and would depend solely on the size of the database, while the dilution factor allows for custom 1: o hiding of the interesting records. Needless to say, the dilution factor can be dynamically set depending on the characteristics of the enquiry. If an enquiry is urgent, then a low dilution factor can be set in order to speed up the processing.

7.3 Motivating scenarios with solutions

Chapter 1 has provided two motivating scenarios that helped to illustrate the problems in the current data acquisition process. The solutions to these scenarios are provided below. For these, take into considerations that the processing times discussed below are taken from the complexity tables, and the trial runs of the encryption protocols used to build IDAP. It is possible, though, to improve on these if the methods for fast exponentiation are employed and the implementation is done in a compiled C programming language (as shown in [78]).

Scenario 1 – Request for ISP subscriber data:

When using IDAP to retrieve data from an ISP the provider would not be given a list of interesting IP addresses, which are the addresses of the potential suspects in an investigation. Thus, the rights of the *individuals-under-investigation* should be preserved. Since, on the other hand, the data acquisition notice would be served under RIPA, it would not have to carry any justification to the dataholder, and with no identities linked to the notice it would be unlikely to compromise an investigation. Performance-wise, the IP address assignment will tell investigators the name(s) of ISPs that provide these IP addresses. This also means that a directory of the *sender's* database, namely the list of IP addresses served by a given ISP is public, or can be inferred with certain probability (although the percentage of unused addresses for a

given ISP would have to be known). The scenario specifies that there are 14 interesting IP addresses. Let assume that the dilution factor can be as little as 1,000, since such enquiries are common and the data holder does not know the reason for the request. If all IPs are provided by the same ISP, then IDAP is run against $14 \times 1,000$ records. Such an enquiry would only take 10 minutes to complete under the given parameters. Thus, the records would be returned within the 30-minute time window that under currently is achievable only in life threatening situations, as discussed in Chapter 4.

Scenario 2 – Banking transaction details:

In this scenario, if there is yet no official warrant for the enquiry, the request would need to be made under DPA. In such a case, the data controller can use the voluntary disclosure mechanism of the Act to provide investigators with the relevant data. However, the investigators would need to inform the data controller about the nature of the investigation, in order to persuade them to disclose the data records. But, there is no technical reason for the data controller to know the exact identity of the interesting record. The request should, in fact, be considered based on the circumstances described by the request, and not the identity of the suspect. Thus, by making IDAP enquiry for the records related to a given credit card number, the relations between the data-subject and the bank would not be affected, however, the investigation could be compromised as its details need to be disclosed. Thus, if the secrecy of the investigation is important, an IDAP enquiry could be made under a court warrant, still hiding the identity of the interesting record, with the warrant in place the data controller could not question the enquiry.

As to performance, at the time of writing there were 58 million credit cards in UK [130]. The initial few digits of the credit card can easily be used to narrow down the bank of the potential suspect. With at least 10 different banks offering such cards in the UK it is more than likely that the initial digits of the credit card (that specify its type and the issuing bank) the search could be narrowed down to less than 10 million records. Therefore, it would take the bank approximately two and a half days to build an encrypted dictionary (using an ordinary computer similar in the specification to the one used to provide experimental results). If the directory of the card numbers is already encrypted then the inquiry itself would take 38 minutes for a dilution factor

of 100,000, or less than a minute if the dilution factor is 1,000 (note that with the use of streamlined implementation of IDAP the processing time would be at least a magnitude shorter).

In both cases the dilution factor introduced in this thesis would enable larger enquiries, with more interesting records but less than 150 as set out by Req. 1, to run in virtually the same time as an enquiry for a single record (see Figure 6-12). This is a very desirable property.

7.4 Contribution to Knowledge

This study presents a methodology for retrieving investigative data in a private manner. The following contributions to knowledge were made by this thesis:

- 1) This thesis demonstrated the manner in which SPIR techniques can be used to assist public authorities in privacy-preserving retrieval of investigative data from third parties. It has been established that the current process of data acquisition attempts to minimise the collateral damage that can be caused by investigations. However, they are unable to hide the identity of the suspects, and thus stop short of protecting privacy. The novel approach to performing investigative data acquisition presented here can add this ability to the already tried and tested data acquisition framework. Therefore, it can be used as a tool that can provide an enhanced level of privacy during the acquisitions, without the need to redesign the already established and well-defined process.
- 2) The problem of investigative data acquisition has been reduced to a single-database SPIR problem. The simulations presented in Chapter 4 have shown that a SPIR-based system for the data acquisition is not feasible if all the records in a large database must be processed. However, it was found that, in most investigative scenarios, such as those presented in [131], the investigators would need to make a single acquisition request per given identifier. For example, if phone billing information is required for a given phone number, the investigators could identify a relevant dataholder from the number allocation that is publicly available for call

routing purposes. Consequently, with certain exceptions and precautions, the data acquisition process can be treated as a single database SPIR, and therefore k -anonymised queries that use the dilution factor can reduce the complexity of the protocols without leaking any data about the suspects, or to be more precise hiding each suspect in a group of records.

- 3) Chapter 6 evaluated the performance of IDAP, while its legal aspects were discussed in Chapter 5. This thesis has shown that it is possible to deploy IDAP into the real-life scenarios in order to benefit the privacy of the data-subjects and the secrecy of the investigations.
- 4) Section 5.3.1 of this thesis introduces a novel approach of employing k -anonymity principles in SPIR protocols. This is done in order to improve the performance of single-database SPIR systems. The main benefits of using such an approach are presented in Section 6.2.5 that illustrates performance of IDAP against PE protocol, and that of $1-n$ OT-based solution against its modification incorporating the dilution factor.
- 5) The definition of a technique, for building complex privacy-preserving enquiries has been provided in Section 5.3.2. It is based on hashing multiple selection criteria together to form one value that can be compared using m -to- n PEqT protocol. For this reason it carries almost no additional computational (or communicational) complexity, however, it may put a larger strain on databases.
- 6) The key concept in enhancing privacy is to also enhance the perception of the final system by the individuals it aims to protect. This thesis proposes that encrypted requests should be relied on by an independent semi-trusted third party (Section 5.3.3). This party would ensure that the investigators can only get the data specified as interesting at the start of the investigation, thus, it would thwart any potential cryptographic attacks by the investigators on the data-records unrelated to the investigation, as these would be removed from the communication stream between the dataholder and the investigators. It is hoped that this technique should gain support of the general public.

7.5 Critical analysis

This thesis has focused on showing that it is feasible for the public authorities to use privacy-preserving techniques while retrieving personal data from third parties. There are various motivations for the public authorities to request third-party data. It may be required by ambulance crews to pinpoint the location of a casualty using data from a mobile communications provider, or to find out the last transactions on a credit card that is being used by a suspect. The thesis focuses on the investigative scenarios as they carry a greater, and more evident, risk of privacy and human rights violations. However, the resultant IDAP can, and should, be applied to most requests for personally identifiable data made under RIPA and the voluntary disclosure mechanism of DPA.

This thesis employs the notion of k -anonymity to form a dilution factor for investigative data acquisition enquiries. With this factor, the balance between the privacy and performance can be dynamically controlled by the requesting party. For example, the Designated Person that scrutinises the requests for data [8] can decide the level of privacy and secrecy that given request should involve by manipulating the dilution factor. Such approaches, using k -anonymity, are now becoming popular in the PET domain. The inspiration for the dilution factor was the use of the k -anonymity models in the statistical data mining [85]. However, it is important to consider advances in the use of such models. These include the application of k -anonymity to location based services, hiding historical, current or trajectory location of the mobile users [132]. In this domain it is apparent that the algorithms for choosing the right records to obfuscate the request are vital to the privacy considerations. In this thesis it is assumed that a random selection of $(o - 1)$ records that are distributed in the dataset uniformly with the interesting record is a sufficient approach. This should allow for making a single query about a record, however, if the query is for the same interesting record and was to be repeated, the dataholder could possibly infer the identity of the interesting record. A solution to this problem could be similar to the one used in the SCOP-based PIR described in [108]. Namely, every new request can include a number of the identifiers that have previously being requested. In this way a constant number of records would overlap between the requests, and the dataholder would be unable to infer any extra information from

overlapping requests. However, this thesis does not provide enough detail about the algorithm for selecting the records that obfuscate the true object of the enquiry, and this is a potential for further work.

In Figure 6-13 it is clear that whether IDAP is feasible or not depends on the dilution factor o . The higher the value of this factor, the higher the run time for the protocol is, and the less practical the use of IDAP. However, it should be noted that Chapter 6 provides results from simulated runs of IDAP, based on the complexity tables and empirically-obtained measures of the time required to perform different cryptographic operations. It should also be noted that the time measurements for different cryptographic operation have been obtained from Microsoft .NET C# implementation of the encryption techniques. It is certain that with use of machine compiled programming language such as C or C++, the performance of IDAP would be at the least a magnitude better. This is confirmed by the results achieved by other researchers [78]. Additionally, under IDAP, both parties can perform some of the processing in parallel, and the cryptographic operations themselves can also be parallelised. Thus, the performance of IDAP will depend on the hardware used to implement it. Consequently, the level of privacy that the framework can offer will depend on the monetary investment into IDAP. This confirms that *privacy is not free* [77], however, the monetary cost of computational power is decreasing.

The advances of cloud computing allow companies and individuals to purchase computing resources on at ad-hoc and pay-per-hour usage [133]. Running of IDAP in a cloud is not recommended if the data is stored in-house, but if this data is already in the cloud there is no reason why the computational power of the cloud should not be harnessed to perform IDAP, where it is possible to purchase *High-CPU On-Demand Instances* of virtual machines for less than 50 pence per hour [133]. Thus, an operation that would usually take hours to complete can be completed in the space of minutes, if 60 instances of such virtual computing machines are used at a cost of £30 per hour. However, the availability of cheap computing resources is also a threat to the system, as the perpetrators can also harness these resources. This is already evident on an example of cloud-based services offering cracking of WPA keys under the cover of penetration testing services [134].

The cost of communications should also be considered when discussing IDAP. This cost depends on the dilution factor, just as the cost of processing does. Thus, for low values of o , such as 1,000, the cost of communications should be reasonable. However, where higher degree of privacy and secrecy is required, the costs of on-line communications could prove to be prohibitive. In such cases, it would be possible to exchange encrypted data via the post or couriers, as there are a small number of communication rounds between the parties.

This thesis suggests that commutative encryption protocols can provide for efficient privacy-preserving m -to- n equality tests, and also privacy-preserving information retrieval. While 1-to-1 PEqT protocols based on homomorphic encryption, such as [98], perform as well as 1-to-1 equality protocols based on commutative encryption, in scenarios with multiple records commutative encryption performs significantly better (Section 3.3.5 provides more detail). The same applies to the SPIR or OT protocols, while 1- n protocols based on commutative encryption perform on par with other 1- n SPIR methods. In this, the commutative protocols have a large advantage in m - n protocols, as only a single encryption of the records in the dataset is required for multiple interesting records, while with virtually any other encryption technology these records would need to be encrypted m times. For this reason, IDAP is based on commutative encryption. However, even though this type of encryption has been around for a few decades, as it was first scrutinised by Shamir in [47], it has not been explored and evaluated to the level that would allow a commutative encryption algorithm to become an acceptable encryption standard. Thus, before IDAP could become recognised as a suitable solution to mitigating collateral damages in investigative data acquisition scenarios, the commutative protocols employed and the logic of the protocol would need to be scrutinised by cryptanalysts. So far IDAP has been peer reviewed in [9, 125, 135], but in the area of information security this should be treated as a basic *sanity* check. It is likely that flaws will be found, in the protocol or its implementation [66, 109]. For example, in [55], Weis points out that most commutative encryption schemes, including the one presented in this thesis, are not semantically secure. This is caused by the very nature of the commutative encryption expressed by Eqn. 2-3, and cannot be avoided if this property is required. It can be assumed that the reason the commutative encryption has not been evaluated in the fashion that other encryption protocols are, is the fact that under commutative

encryption certain links of the ciphertext to the plaintext are desirable. However, the main contribution to the knowledge of this thesis, are the three improvements proposed to make PE meet the requirements for IDAP. These improvements can be applied in most SPIR protocols and, for this reason, even if commutative cryptography is rejected by the information security community as a basis for IDAP, the findings of this thesis can still be used to build IDAP based on another SPIR protocol. This should be feasible as Figure 6-12 shows the use of the dilution factor allows an approach based on 1- n OT to perform just as well PE does for a certain range of γ .

The commutative encryption PEqT protocol employed in IDAP compares hashes of the identities from the set of suspects to the hashes of the identities in the *sender's* dataset. Therefore, there is no room for error. The identities need to be an exact match, so matching on values such as telephone numbers and IP addresses is preferred. There is no easy way to make IDAP match a different spelling of a name, such as John Smith and John Smyth. However, during the investigation of the subject area, a computationally-expensive method for achieving fuzzy matches has been developed and presented in [125].

7.6 Main findings

The current processes for data acquisition are designed to minimise the probability of collateral damage to the suspects and the investigations [8]. Still, the identities of the data records need to be provided to the dataholders, simply to identify the interesting data records. IDAP fits within this current methodology as it can provide the means to request data without revealing the identities of the interesting records. Thus, it is a tool that SPoCs can use to communicate, that also allows an independent watchdog organisation to monitor the exchange of acquisition notices and data. Consequently, the internal processes associated with data acquisition do not need to change beyond the fact that, under IDAP, the dataholder provides a large number of records as an input to the protocol, with the records locked by the encryption process in a way that renders them unusable to anybody who does not know the relevant encryption keys. The requesting party can then unlock only the records defined as interesting in the privacy-preserving data acquisition notice. In order to make the process more future-

proof, and to gain acceptance of IDAP in the society, the watchdog organisation filters out the diluting records from the response, itself knowing only the encrypted form of the identifiers for the interesting records.

The operations of IDAP must not affect the validity of the gathered information as potential evidence for the use in a court of Law [17]. Since IDAP uses a number of cryptographic techniques to retrieve and decrypt the interesting records, some lawyers could potentially question the validity of the evidence, and lower the status of the data gathered to second-hand, or hearsay, evidence [131]. One solution could be to get the dataholder to retain the original records provided under each enquiry for a set period of time. However, since the dataholder does not know the identities of the interesting records, then all the records provided as an input to IDAP would have to be stored. This would be excessive, and a costly exercise. For this reason, under IDAP, the proxy is responsible for retaining the records identified by the public authorities as interesting, while the dataholder retains only the commutative encryption key used to lock the data records (there is only one per enquiry). Then if a verification of the evidence is required, these two parties could be made responsible for working together to validate the data.

It is interesting that from specifying who needs given piece of data and why, some knowledge may be inferred. For example if a law enforcement officer requests to know the last location of an individual it is likely that the individual is a suspect, while if the emergency services request this information most people would assume that the individual is a casualty. With IDAP using a watchdog organisation as a central hub for all enquiries (Figure 7-1), the requests from investigators will be diluted by requests from other public authorities, including the health authorities and this will contribute to the privacy levels provided by the system.

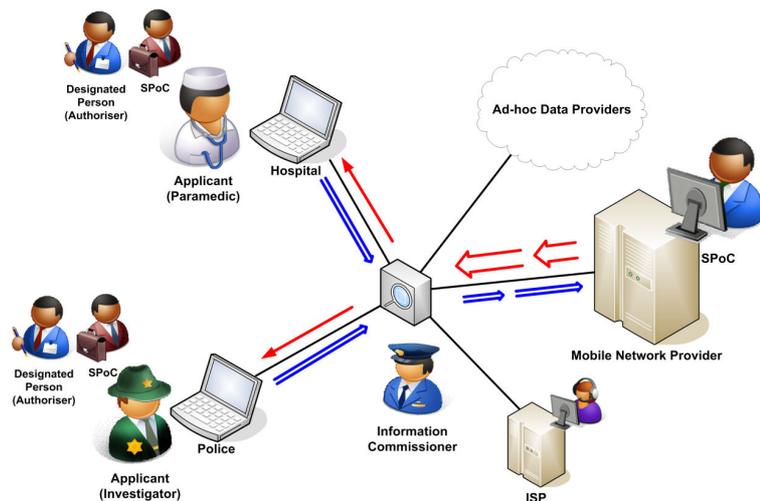


Figure 7-1 Hiding of request originators can improve privacy

This thesis is a response to the perceived trend of limiting privacy of individuals in order to protect the security of a nation. An integral part of the problem of privacy is the public's perception. Therefore, the privacy, just like beauty, is in the eye of the beholder, and it is vital for any privacy measures to be acceptable by the individuals they try to protect, otherwise they are meaningless. This is the reason why achieving privacy is a complex matter beyond any technological solution. For this reason, the following paraphrased statement about security could also be used in relation to privacy:

If you think technology can solve your security problems, then you don't understand the problems and you don't understand the technology!

Ferguson and Schneier, [136], pp. XXII

This thesis suggests (in Chapter 4), that security risk management can be employed to predict and manage risks to the privacy. However, this can only be done to a certain extent, as privacy, just like security, is a *people problem* [137]. Because of the human involvement of the both sides of the parameter, there is no way to predict all the vulnerabilities and potential threats to individual's privacy, and it is far harder to predict individual's perception of the privacy measures applied. Therefore, IDAP has been qualitatively evaluated with the help of a number of carefully selected individuals, and the results show the use of IDAP to perform data acquisition will be welcomed.

It is crucial that the merits of IDAP are communicated well in order to get as high a level of acceptance as possible. Professor Burkhard Schafer, Professor of Computational Legal Theory at Edinburgh University, when asked to give feedback on IDAP suggested an analogy to the “I am Spartacus” defence [138], for the way this platform protects the privacy of individuals under investigation by making a large group of individuals appear as suspect. If the size of the group is large enough then no repercussions can be applied and therefore the third party is unable to act on the information provided. This kind of defence was recently used by Internet users around the world in a protest against the judgement on a case of a Twitter airport bomb threat joke [139], showing that this type of defence is potentially acceptable by the society.

IDAP, as defined in this thesis, is a feasible approach to acquiring investigative data from third parties, and such privacy-preserving solutions are generally welcomed by the governing organisation, the ICO [126]. However, the platform’s computational complexity may limit its use in the most urgent scenarios. Therefore, some alternative arrangements may be necessary for enquiries made to organisations that are required to provide investigative information on daily basis, to the level that justifies employing full-time personnel just to handle requests. In such a case, strong privacy-preserving procedures may be put in place, with certification system for organisations and vetting programme for the individuals handling the requests for investigative data.

7.7 Future Work

There are a number of important areas of further investigation, and different applications for the protocols discussed. The areas of future work include:

- Expand on the discussions on IDAP’s security and correctness by providing a complete security analysis of the platform.
- Evaluate in detail the communicational complexity of IDAP.
- Evaluate database load. Searching databases on columns containing record IDs is usually a fast operation, since such columns are usually indexed by the database system. This thesis proposed a technique for privately matching

records on multiple selection criteria, however, these criteria would potentially include columns that are not indexed. Use of this technique may cause considerable load on the database and therefore needs to be evaluated.

- Investigate of a technique for selecting an appropriate dilution factor and diluting records depending on the circumstances.
- Include warrant signing and verification system in IDAP. Similar system already exist as described in [74].
- Evaluate an add-on to office packages, such as Microsoft Office and Open Office that could perform IDAP on spreadsheets and thus it could enable small organisations to utilise IDAP.
- Investigate using the components of IDAP in privacy-preserving proximity- and location-based dating and social networking. There is a high degree of risk to privacy of individuals actively looking to find a partner or expand the network of people they know. With some components of IDAP, it would be possible to actively advertise ones own preferences using a mobile device, allowing for the individuals matching the profile and advertising similar values and interests to recognise each other. f

Chapter 8

References

- [1] P. Swire and L. Steinfeld, "Security and privacy after September 11: the health care example," Proceedings of the 12th annual conference on Computers, freedom and privacy, San Francisco, California, 2002, pp. 1-13.
- [2] Home Office, "Protecting the Public in a Changing Communications Environment," 2009. [Online]. Available: <http://www.homeoffice.gov.uk/documents/cons-2009-communications-data?view=Binary> [Accessed: 20 Apr 2009].
- [3] Rasmussen Reports, "51% Say Security More Important than Privacy," 2008.
- [4] D. Balz and C. Deane, "Differing Views on Terrorism," The Washington Post, 2006. [Online]. Available: <http://www.washingtonpost.com/wp-dyn/content/article/2006/01/10/AR2006011001192.html> [Accessed: 11 Nov 2006].
- [5] Ponemon Institute, "2010 Privacy Trust Study of the United States Government," Ponemon Institute, Traverse City, MI 06302010, 20 December 2010.
- [6] P. Kennedy, Sir, "Report of the Interception of Communications Commissioner for 2008," Information Commissioner's Office, London SG/2009/138, 21 July 2009.

- [7] European Parliament, "European Data Protection Directive 95/46/EC," *Official Journal of the European Union*, vol. L, 1995, pp. 31-50.
- [8] Home Office, "Acquisition and Disclosure of Communications Data - Code of Practice." London: Home Office, 2007.
- [9] Z. Kwecka, W. Buchanan, D. Spiers, and L. Saliou, "Validation of 1-N OT Algorithms in Privacy-Preserving Investigations," Proceedings of the 7th European Conference on Information Warfare and Security, University of Plymouth, 2008, pp. 119-128.
- [10] Crown, "Data Protection Act 1998," TSO, 1998.
- [11] Crown, "Human Rights Act 1998," TSO, 1998.
- [12] "Regulation of Investigatory Powers Act 2000," Crown - The Stationery Office Limited, 2000.
- [13] G. Palmer, "A road map for digital forensic research," DFRWS, Utica, NY, Technical Report, 7 Aug 2001.
- [14] BCS, "Submission to the Criminal Courts Review," 2000. [Online]. Available: <http://www.computerevidence.co.uk/Papers/LJAuld/BCSComputerEvidenceSubmission.htm> [Accessed: 20 Dec 2006].
- [15] M. Carney and M. Rogers, "The Trojan Made Me Do It: A First Step in Statistical Based Computer Forensics Event Reconstruction " *International Journal of Digital Evidence*, vol. 2, 2004, pp. 1-11.
- [16] W. J. Buchanan, J. Graves, L. Saliou, H. A. Sebea, and N. Migas, "Agent-based Forensic Investigations with an Integrated Framework," Proceedings of the 4th European Conference of Information Warfare and Security, University of Glamorgan, UK, 2005, pp. 47-52.
- [17] Association of Chief Police Officers, "Good Practice Guide for Computer based Electronic Evidence (version 3)," 2003. [Online]. Available: http://www.acpo.police.uk/asp/policies/Data/gpg_computer_based_evidence_v3.pdf [Accessed: 19 Dec 2006].
- [18] DNS, "An Introduction to Computer Forensics," 2006. [Online]. Available: www.dns.co.uk/NR/rdonlyres/5ED1542B-6AB5-4CCE-838D-D5F3A4494F46/0/ComputerForensics.pdf [Accessed: 15 May 2007].
- [19] "Standardizing digital evidence storage," *Commun. ACM, The Common Digital Evidence Storage Format Working Group*, vol. 49, 2006, pp. 67-68.
- [20] Y. Tang and T. E. Daniels, "A Simple Framework for Distributed Forensics," Proceedings of the Second International Workshop on Security in Distributed Computing Systems (SDCS) (ICDCSW'05) - Volume 02, 2005, pp. 163-169.
- [21] J. Tozer, "Retired police chief arrested over 'misuse of data' in hunt for female colleague's missing dog," Daily Mail, 2009. [Online]. Available: <http://www.dailymail.co.uk/news/article-1124625/Retired-police-chief-arrested-misuse-data-hunt-female-colleagues-missing-dog.html> [Accessed: 20 Jul 2010].
- [22] "Addressing the data problem," *Information Security Technical Report*, vol. 8, 2003, pp. 18-31.
- [23] The Crown, *Regulation of Investigatory Powers Act 2000: explanatory notes Chapter 23 2nd impression*. London: The Stationery Office, 2000.

- [24] Home Office, "About RIPA," 2000. [Online]. Available: <http://security.homeoffice.gov.uk/ripa/about-ripa/> [Accessed: 20 Jun 2008].
- [25] D. Spiers, *Intellectual Property Law Essentials*. Dundee: Dundee University Press, 2009.
- [26] Home Office, "Retention of Communications data under art 11: Anti-terrorism, Crime & Security ACT 2001. Voluntary Code of Practice," 2001. [Online]. Available: <http://www.opsi.gov.uk/si/si2003/draft/5b.pdf> [Accessed: 10 Jun 2008].
- [27] E. Kosta and P. Valcke, "Retaining the data retention directive," *Computer Law & Security Report*, vol. 22, 2006, pp. 370-380.
- [28] Scottish Courts, "Recovery of Evidence," in *Rules of the Court of Session*. Edinburgh: Scottish Courts,, 1994.
- [29] A. MacSporran and A. Young, *Commission and Diligence*: W. Green, 1995.
- [30] "IP Litigation Discovery Escrow," Iron Mountain, 2010. [Online]. Available: <http://www.ironmountain.ca/en/whatwedo/technology-escrow/ip-litigation-discovery-escrow.asp> [Accessed: 10 Dec 2010].
- [31] J. DeCew, "Privacy," Stanford Encyclopedia of Philosophy, 2006. [Online]. Available: <http://plato.stanford.edu/entries/privacy/> [Accessed: 10 Aug 2007].
- [32] J. Waldo, H. S. Lin, and L. I. Millett, *Engaging Privacy and Information Technology in a Digital Age*. Washington, D.C.: The National Academies Press, 2007.
- [33] Privacy Rights Clearinghouse, "A Review of the Fair Information Principles: The Foundation of Privacy Public Policy," 1997. [Online]. Available: <http://www.privacyrights.org/ar/fairinfo.htm> [Accessed: 20 Jul 2010].
- [34] G. T. Marx, "An Ethics For The New Surveillance," *The Information Society*, vol. 14, 1998.
- [35] N. McKAY, "Lawmakers Raise Questions About International Spy Network," in *The New York Times*. New York, 1999.
- [36] RAPID, "Telecoms: Commission launches case against UK over privacy and personal data protection " Europa Press Releases RAPID, 2009. [Online]. Available: <http://europa.eu/rapid/pressReleasesAction.do?reference=IP/09/570&format=HTML&aged=0&language=EN&guiLanguage=en> [Accessed: 15 Apr 2009].
- [37] K. Ball, D. Lyon, D. M. Wood, C. Norris, and C. Raab, "A Report on the Surveillance Society," Surveillance Studies Network.
- [38] D. Campbell and R. Evans, "Surveillance on drivers may be increased," *Guardian*, 2006. [Online]. Available: <http://www.guardian.co.uk/frontpage/story/0,,1725360,00.html> [Accessed: 07 Mar 2006].
- [39] M. Gill and A. Spriggs, "Assessing the impact of CCTV," Home Office Research, Development and Statistics Directorate.
- [40] B. C. Young, E. C. Kathleen, S. K. Joshua, and M. S. Meredith, "Challenges Associated with Privacy in Health Care Industry: Implementation of HIPAA and the Security Rules," *J. Med. Syst.*, vol. 30, 2006, pp. 57-64.

- [41] B. Schneier, *Applied Cryptography: Protocols, Algorithms, and Source Code in C*: John Wiley & Sons, Inc., 1995.
- [42] W. Diffie and M. E. Hellman, "New Directions in Cryptography," *IEEE Transactions on Information Theory*, vol. IT-22, 1976, pp. 644-654.
- [43] "Public-Key Encryption," in *Handbook of Applied Cryptography*, A. J. Menezes, P. C. V. Oorschot, and S. A. Vanstone, Eds.: CRC Press, 1997, pp. 283-319.
- [44] R. L. Rivest, A. Shamir, and L. Adleman, "A method for obtaining digital signatures and public-key cryptosystems," *Commun. ACM*, vol. 21, 1978, pp. 120-126.
- [45] N. Ferguson and B. Schneier, *Practical Cryptography*. Indianapolis, Indiana: Wiley Publishing Inc., 2003.
- [46] C. Shannon, "Communication theory of secrecy systems," *Bell System Technical Journal*, vol. 28, 1949, pp. 656-715.
- [47] A. Shamir, "On the Power of Commutativity in Cryptography," Proceedings of the 7th Colloquium on Automata, Languages and Programming, 1980, pp. 582-595.
- [48] R. C.-W. Phan and B.-M. Goi, "(In)Security of an Efficient Fingerprinting Scheme with Symmetric and Commutative Encryption of IWDW 2005," *IWDW*, vol. 5041, 2007, pp. 34-44.
- [49] S. Cheung, H. Leung, and C. Wang, "A commutative encrypted protocol for the privacy protection of watermarks in digital contents," Proceedings of the 37th Hawaii International Conference on System Sciences, Hawaii, vol. 4, 2004, pp. 1-10.
- [50] C. Lei, P. Yu, P. Tsai, and M. Chan, "An efficient and anonymous buyer-seller watermarking protocol," *IEEE Transactions on Image Processing*, vol. 13, 2004, pp. 1618-1626.
- [51] M. Stevens, A. Sotirov, J. Appelbaum, A. Lenstra, D. Molnar, D. A. Osvik, and B. Weger, "Short Chosen-Prefix Collisions for MD5 and the Creation of a Rogue CA Certificate," Proceedings of the CRYPTO '09: 29th Annual International Cryptology Conference on Advances in Cryptology, Santa Barbara, CA, 2009, pp. 55-69.
- [52] X. Wang, Y. L. Yin, and H. Yu, "Finding Collisions in the Full SHA-1," Proceedings of the CRYPTO '05: 25th Annual International Cryptology Conference on Advances in Cryptology, Santa Barbara, CA, 2005, pp. 17-36.
- [53] B. Schneier, "Cryptanalysis of SHA-1," Schneier on Security, 2005. [Online]. Available: http://www.schneier.com/blog/archives/2005/02/cryptanalysis_o.html [Accessed: 20 Jun 2008].
- [54] T. ElGamal, "A public key cryptosystem and a signature scheme based on discrete logarithms," Proceedings of the CRYPTO 84 on Advances in cryptology, Santa Barbara, California, US, 1985, pp. 10-18.
- [55] S. A. Weis, "New Foundations for Efficient Authentication, Commutative Cryptography, and Private Disjointness Testing," in *Department of Electrical Engineering and Computer Science*, vol. PhD. Cambridge, MA: Massachusetts Institute of Technology, 2006, pp. 115.

- [56] M. Kuribayashi and H. Tanaka, "Fingerprinting protocol for images based on additive homomorphic property," *Image Processing, IEEE Transactions on*, vol. 14, 2005, pp. 2129-2139.
- [57] C. Gentry, "Fully homomorphic encryption using ideal lattices," Proceedings of the 41st annual ACM symposium on Theory of computing, Bethesda, MD, USA, 2009, pp. 169-178.
- [58] A. Shamir, R. L. Rivest, and L. Adleman, "Mental Poker," in *The Mathematical Gardner*, D. A. Klarner, Ed. Boston, MA: Prindle, Weber & Schmidt, 1981, pp. 37-43.
- [59] S. H. Khayat, "Using Commutative Encryption to Share a Secret," Cryptology ePrint Archive 2008/356, 2008.
- [60] R. Agrawal, A. Efvimievski, and R. Srikant, "Information sharing across private databases," Proceedings of the Proceedings of the 2003 ACM SIGMOD international conference on Management of data, San Diego, California, 2003, pp. 86-97.
- [61] J. L. Massey and J. K. Omura, "Method and apparatus for maintaining the privacy of digital messages conveyed by public transmission," United States Patent and Trademark Office, 1986, United States patent number: 4567600
- [62] A. Reyhani-Masoleh, "Efficient Algorithms and Architectures for Field Multiplication Using Gaussian Normal Bases," *IEEE Trans. Comput.*, vol. 55, 2006, pp. 34-47.
- [63] P. Golle, M. Jakobsson, A. Juels, and P. Syverson, "Universal re-encryption for mixnets," Proceedings of the 2004 RSA Conference, Cryptographer's track, 2004, pp. 163-178.
- [64] W. Zhao, V. Varadharajan, and Y. Mu, "A secure mental poker protocol over the internet," Proceedings of the Proceedings of the Australasian information security workshop conference on ACSW frontiers 2003 - Volume 21, Adelaide, Australia, 2003, pp. 105-109.
- [65] C. Swenson, *Modern Cryptanalysis: Techniques for Advanced Code Breaking*. Indianapolis: Wiley Publishing, Inc., 2008.
- [66] R. Anderson, "Why cryptosystems fail," *Communications of the ACM*, vol. 37, 1994, pp. 32-40.
- [67] L. Cranor, J. Reagle, and M. Ackerman, "Beyond Concern: Understanding Net Users Attitudes About Online Privacy," AT&T Labs-Research Technical Report TR 99.4.3.
- [68] A. F. Westin, "Freebies and Privacy: What Net Users Think," Opinion Research Corporation.
- [69] W. M. Grossman, "New money," net.wars, 2010. [Online]. Available: <http://www.pelicancrossing.net/netwars/privacy/#a000272> [Accessed: 20 Aug 2010].
- [70] R. Jansen, N. Hopper, and Y. Kim, "Recruiting new tor relays with BRAIDS," Proceedings of the 17th ACM conference on Computer and communications security, Chicago, Illinois, USA, 2010, pp. 319-328.
- [71] Microsoft, "Microsoft's Vision for an Identity Metasystem," 2005. [Online]. Available: <http://msdn2.microsoft.com/en-us/library/ms996422.aspx> [Accessed: 12 Mar 2007].

- [72] D. Chaum, A. Fiat, and M. Naor, "Untraceable Electronic Cash," Proceedings of the CRYPTO, 1990, pp. 319-327.
- [73] B. Aiello, Y. Ishai, and O. Reingold, "Priced Oblivious Transfer: How to Sell Digital Goods," in *Advances in Cryptology — EUROCRYPT 2001*, vol. 2045, *Lecture Notes in Computer Science*, B. Pfitzmann, Ed.: Springer-Verlag, 2001, pp. 119-135.
- [74] K. B. Frikken and M. J. Atallah, "Privacy preserving electronic surveillance," Proceedings of the 2003 ACM workshop on Privacy in the electronic society, Washington, DC, 2003, pp. 45-52.
- [75] J. Biskup and U. Flegel, "Threshold-based identity recovery for privacy enhanced applications," Proceedings of the 7th ACM conference on Computer and communications security, Athens, Greece, 2000, pp. 71-79.
- [76] J. Biskup and U. Flegel, "Transaction-Based Pseudonyms in Audit Data for Privacy Respecting Intrusion Detection," Proceedings of the 3rd International Workshop on Recent Advances in Intrusion Detection, 2000, pp. 28-48.
- [77] M. Kantarcioglu and C. Clifton, "Assuring privacy when big brother is watching," Proceedings of the 8th ACM SIGMOD workshop on Research issues in data mining and knowledge discovery, San Diego, California, 2003, pp. 88-93.
- [78] E. Cristofaro, S. Jarecki, J. Kim, and G. Tsudik, "Privacy-Preserving Policy-Based Information Transfer," Proceedings of the PETS '09: 9th International Symposium on Privacy Enhancing Technologies, Seattle, WA, 2009, pp. 164-184.
- [79] R. Koorn, H. v. Gils, J. t. Hart, P. Overbeek, and R. Tellegen, "Privacy-Enhancing Technologies White Paper for Decision-Makers," 2004. [Online]. Available: http://www.dutchdpa.nl/downloads_overig/PET_whitebook.pdf [Accessed: 01/02/2009].
- [80] F. Bao and R. Deng, "Privacy Protection for Transactions of Digital Goods," in *Information and Communications Security*, 2001, pp. 202-213.
- [81] B. W. Lampson, "A note on the confinement problem," *Commun. ACM*, vol. 16, 1973, pp. 613-615.
- [82] M. Kantarcioglu and J. Vaidya, "An architecture for privacy-preserving mining of client information," Proceedings of the IEEE international conference on Privacy, security and data mining, Maebashi City, Japan, vol. 14, 2002, pp. 37-42.
- [83] C. C. Agarwal and P. S. Yu, "On static and dynamic methods for condensation-based privacy-preserving data mining," *ACM Trans. Database Syst.*, vol. 33, 2008, pp. 1-39.
- [84] R. Agrawal and R. Srikant, "Privacy-preserving data mining," *SIGMOD Rec.*, vol. 29, 2000, pp. 439-450.
- [85] P. Samarati, "Protecting Respondents' Identities in Microdata Release," *IEEE Trans. on Knowl. and Data Eng.*, vol. 13, 2001, pp. 1010-1027.
- [86] A. Yao, "Protocols for Secure Computation," *23th FOCS*, 1982, pp. 160-164.

- [87] O. Goldreich, S. Micali, and A. Wigderson, "How to play ANY mental game," Proceedings of the 19th annual ACM conference on Theory of computing, New York, New York, United States, 1987, pp. 218-229.
- [88] S. Goldwasser, "Multi party computations: past and present," Proceedings of the 16th annual ACM symposium on Principles of distributed computing, Santa Barbara, California, United States, 1997, pp. 1-6.
- [89] A. Shamir, "How to share a secret," *Commun. ACM*, vol. 22, 1979, pp. 612-613.
- [90] B. Chor and N. Gilboa, "Computationally private information retrieval (extended abstract)," Proceedings of the 29th annual ACM symposium on Theory of computing, El Paso, Texas, United States, 1997, pp. 304-313.
- [91] R. Ostrovsky and E. Skeith III, "A Survey of Single-Database PIR: Techniques and Applications," in *Public Key Cryptography*, vol. 4450, *Lecture Notes in Computer Science*, O. Tatsuaki and W. Xiaoyun, Eds. Berlin: Springer 2007, pp. 393-411.
- [92] M. Naor and B. Pinkas, "Efficient oblivious transfer protocols," Proceedings of the 12th annual ACM-SIAM symposium on Discrete algorithms, Washington, D.C., United States, 2001, pp. 448-457.
- [93] L. Shundong, W. Daoshun, D. Yiqi, and L. Ping, "Symmetric cryptographic solution to Yao's millionaires' problem and an evaluation of secure multiparty computations," *Inf. Sci.*, vol. 178, 2008, pp. 244-255.
- [94] Michael O. Rabin, "How to exchange secrets by oblivious transfer," Technical Report TR-81, Aiken Computation Laboratory, Harvard University.
- [95] T. Wen-Guey, "Efficient 1-Out-n Oblivious Transfer Schemes," Proceedings of the 5th International Workshop on Practice and Theory in Public Key Cryptosystems: Public Key Cryptography, 2002, pp. 159-171.
- [96] W. Du and M. J. Atallah, "Privacy-Preserving Cooperative Scientific Computations," Proceedings of the 14th IEEE workshop on Computer Security Foundations, 2001, pp. 273-282.
- [97] S. Goldwasser and Y. Lindell, "Secure Computation without Agreement," Proceedings of the 16th International Conference on Distributed Computing, 2002, pp. 17-32.
- [98] H. Lipmaa, "Verifiable Homomorphic Oblivious Transfer and Private Equality Test," in *Advances on Cryptology --- ASIACRYPT 2003*, vol. 2894 *Lecture Notes in Computer Science*, I. C. S. Laih, Ed. Taipei, Taiwan: Springer-Verlag, 2003, pp. 416-433.
- [99] R. Impagliazzo and S. Rudich, "Limits on the provable consequences of one-way permutations," Proceedings of the 21st annual ACM symposium on Theory of computing, Seattle, Washington, United States, 1989, pp. 44-61.
- [100] E.F. Brickell, D.M. Gordon, and K.S. McCurley, "Method for exponentiating in cryptographic systems," 1994, United States patent number: 5299262
- [101] C. Cachin, "Efficient private bidding and auctions with an oblivious third party," Proceedings of the 6th ACM conference on Computer and communications security - CCS '99, Singapore, 1999, pp. 120-127.

- [102] C. Cachin, S. Micali, and M. Stadler, "Computationally Private Information Retrieval with Polylogarithmic Communication," *Proceedings of the EUROCRYPT '99*, 1999, pp. 402--414.
- [103] K. Kurosawa and W. Ogata, "Bit-Slice Auction Circuit," *Proceedings of the 7th European Symposium on Research in Computer Security, Zurich, Switzerland*, 2002, pp. 24-38.
- [104] R. Fagin, M. Naor, and P. Winkler, "Comparing information without leaking it," *Commun. ACM*, vol. 39, 1996, pp. 77-85.
- [105] C. Clifton, M. Kantarcioglu, and J. Vaidya, "Tools for privacy preserving distributed data mining," *ACM SIGKDD Explorations*, vol. 4, 2003.
- [106] M. J. Freedman, K. Nissim, and B. Pinkas, "Efficient private matching and set intersection," 2004, pp. 1-19.
- [107] A. Kiayias and A. Mitrofanova, "Financial Cryptography and Data Security," in *Lecture Notes in Computer Science*, vol. 3570/2005, *Testing Disjointness of Private Datasets*. Heidelberg: Springer Berlin, 2005, pp. 109-124.
- [108] A. Iliev and S. W. Smith, "Protecting Client Privacy with Trusted Computing at the Server," *IEEE Security and Privacy*, vol. 3, 2005, pp. 20-28.
- [109] R. J. Anderson, *Security Engineering: A guide to building dependable distributed systems*, 2nd ed. Indianapolis, In: Wiley Publishing, Inc., 2008.
- [110] R. Anderson and S. Vaudenay, "Minding your p's and q's," *In Advances in Cryptology - ASIACRYPT'96, LNCS 1163*, 1996, pp. 26-35.
- [111] R. Anderson and R. Needham, "Programming Satan's Computer," *in Computer Science Today*, 1995, pp. 426-440.
- [112] R. Anderson, "The Initial Costs and Maintenance Costs of Protocols," 2005.
- [113] B. Kaliski, "RSA Problem," in *ACM SIGKDD Explorations: MIT Laboratory for Computer Science*, 2003, pp. 10.
- [114] R. Silverman and R. Rivest, "Are 'Strong' Primes Needed for RSA," 2001. [Online]. Available: [Accessed: 10 Dec 2009].
- [115] S. Kouichi and S. Hiroki, "A Structural Comparison of the Computational Difficulty of Breaking Discrete Log Cryptosystems," *Journal of Cryptology*, vol. 11, 1998, pp. 29-43.
- [116] Y. Chevalier, R. Kusters, M. Rusinowitch, and M. Turuani, "Deciding the Security of Protocols with Commuting Public Key Encryption," *Electronic Notes in Theoretical Computer Science*, vol. 125, 2005, pp. 55-66.
- [117] A. M. Odlyzko, "Discrete logarithms in finite fields and their cryptographic significance," 1985, pp. 224-314.
- [118] Crown, "Freedom of Information Act 2000," TSO, 2000.
- [119] R. I. M. Dunbar, "Neocortex size as a constraint on group size in primates," *Journal of Human Evolution*, vol. 22, 1991, pp. 469-493.
- [120] J. Scott, "BT tops 15 million broadband users," IT PRO, 2010. [Online]. Available: <http://www.itpro.co.uk/626509/bt-tops-15-million-broadband-users> [Accessed: 31 Aug 2010].
- [121] D. Chamberlin and R. Boyce, "SEQUEL: A structured English query language," *Proceedings of the FIDET '74: Proceedings of the 1974 ACM SIGFIDET (now SIGMOD) workshop on Data description, access and control*, 1974, pp. 249-264.

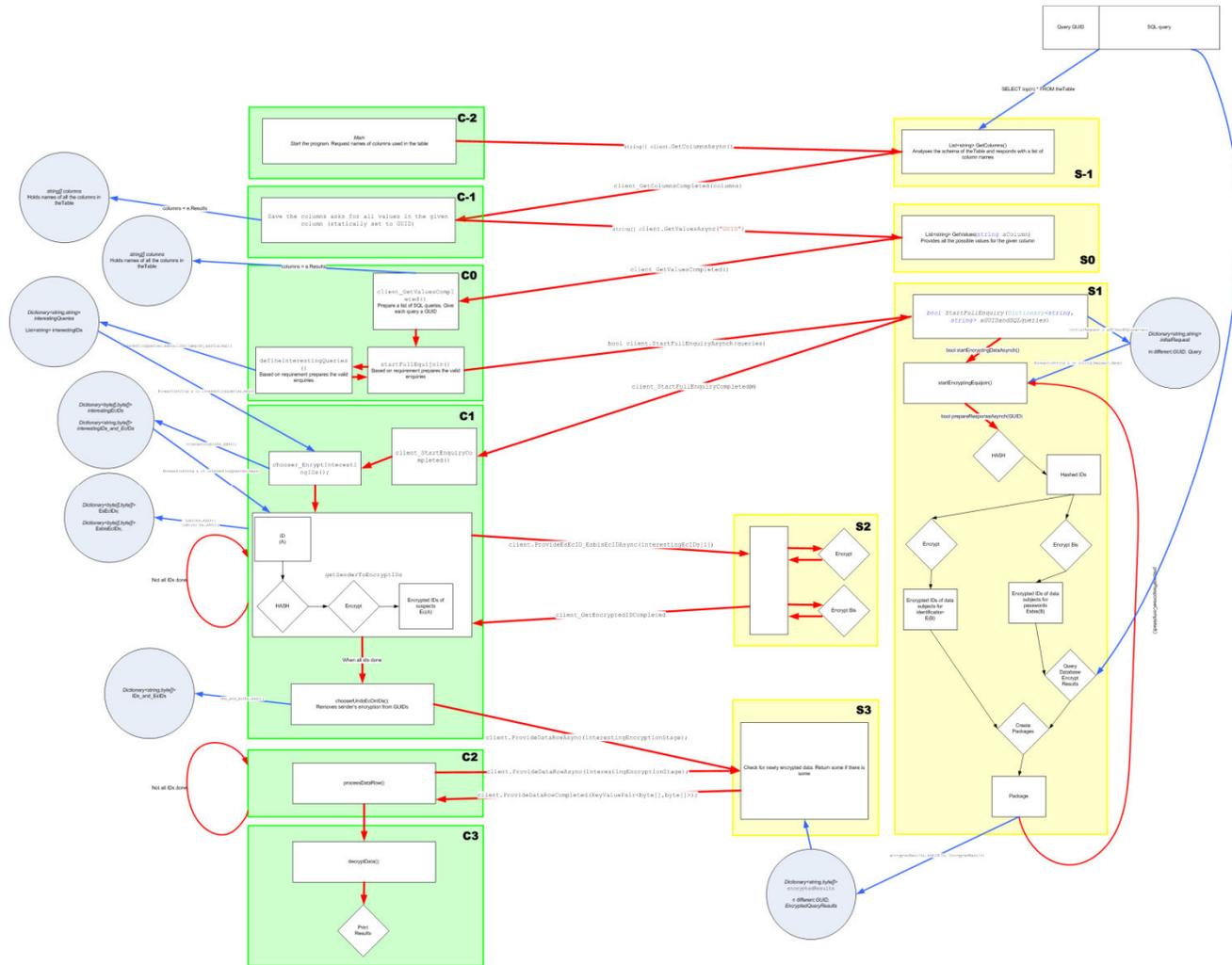
- [122] "Bouncy Castle C# AP," 2007. [Online]. Available: <http://www.bouncycastle.org/csharp/> [Accessed: 20 Aug 2007].
- [123] "Soham trial: 'Crucial' phone evidence," BBC News, 2003. [Online]. Available: <http://news.bbc.co.uk/1/hi/england/cambridgeshire/3246111.stm> [Accessed: 6 November 2003].
- [124] D. Asonov and J.-C. Freytag, "Almost Optimal Private Information Retrieval," in *Privacy Enhancing Technologies*, 2003, pp. 239-243.
- [125] Z. Kwecka, W. J. Buchanan, and D. Spiers, "Privacy-Preserving Data Acquisition Protocol," Proceedings of the IEEE Region 8 International Conference on Computational Technologies in Electrical and Electronics Engineering (SIBIRCON), Irkutsk, vol. 1, 2010, pp. 131-136.
- [126] Information Commissioner, "Data Protection Guidance Note: Privacy enhancing technologies," 2.0 ed. London: Information Commissioner's Office, 2007.
- [127] Microsoft, "Large Data and Streaming," in *MSDN*, vol. 2009: Microsoft, 2009.
- [128] Z. Kwecka, "Future Investigative Data Acquisition Techniques survey," 2010. [Online]. Available: <http://www.evidence-acquisition.org/survey.aspx> [Accessed: 04 Apr 2010].
- [129] F. Olumofin and I. Goldberg, "Privacy-preserving Queries over Relational Databases," in *Privacy Enhancing Technologies Symposium*, vol. 6205, *Lecture Notes in Computer Science*, M. Atallah and N. Hopper, Eds. Berlin: Springer, 2010, pp. 75-92.
- [130] "New credit card protection agreed," BBC News, 2010. [Online]. Available: <http://news.bbc.co.uk/1/hi/business/8567677.stm> [Accessed: 21 Dec 2010].
- [131] S. Mason, *Electronic Evidence: Disclosure, Discoverability & Admissibility*. London: LexisNexis Butterworths, 2007.
- [132] A. Gkoulalas-Divanis, P. Kalnis, and V. S. Verykios, "Providing K-Anonymity in location based services," *SIGKDD Explor. Newsl.*, vol. 12, 2010, pp. 3-10.
- [133] Amazon.com, "Amazon Elastic Compute Cloud (Amazon EC2): Pricing," in *Amazon Web Services*, 2011.
- [134] Chad Perrin, "Welcome to the future: cloud-based WPA cracking is here," TechRepublic, 2010. [Online]. Available: <http://www.techrepublic.com/blog/security/welcome-to-the-future-cloud-based-wpa-cracking-is-here/4097> [Accessed: 10 Dec 2010].
- [135] Z. Kwecka and W. J. Buchanan, "Minimising Collateral Damage: Privacy-Preserving Investigative Data Acquisition Platform.," *International Journal of Information Technologies and Systems Approach (IJITSA): Special issue on Privacy and Security Issues in IT*, vol. 4, 2011.
- [136] B. Schneier, *Secrets and Lies: digital security in a networked world*. Indianapolis, Indiana: Wiley Computer Publishing, Inc., 2004.
- [137] W. Thomas, "The answer is 42 of course," *Queue*, vol. 3, 2005, pp. 34-39.
- [138] S. Kubrick, "Spartacus," Produced by E. Lewis and K. Douglas. United States: Universal Pictures, 1960.
- [139] AP, "Online outrage after judgement of Twitter airport bomb threat joke," The Sunday Morning Herald, 2010. [Online]. Available: <http://www.smh.com.au/technology/technology-news/online-outrage->

after-judgement-of-twitter-airport-bomb-threat-joke-20101115-17t68.html
[Accessed: 21 Nov 2010].

Appendix A

Simplified Operation of IDAP





Full code produced during this thesis is available upon request.

Appendix B

Empirical Evaluation Results



$m = 1, 50 < n < 1 \text{mln}$

m	n	OT	OT	PE	PE
1	50	39.724083	3.542627	2.457119	2.435369
1	100	79.339033	6.990377	4.723569	4.701819
1	500	396.258633	34.572377	22.85517	22.83342
1	1000	792.408133	69.049877	45.51967	45.49792
1	5000	3961.604133	344.869877	226.8357	226.8139
1	10000	7923.099133	689.644877	453.4807	453.4589
1	50000	39615.05913	3447.844877	2266.641	2266.619
1	100000	79230.00913	6895.594877	4533.091	4533.069
1	500000	396149.6091	34477.59488	22664.69	22664.67
1	1000000	792299.1091	68955.09488	45329.19	45329.17

$1 < m < 50000, n = 1 \text{mln}$

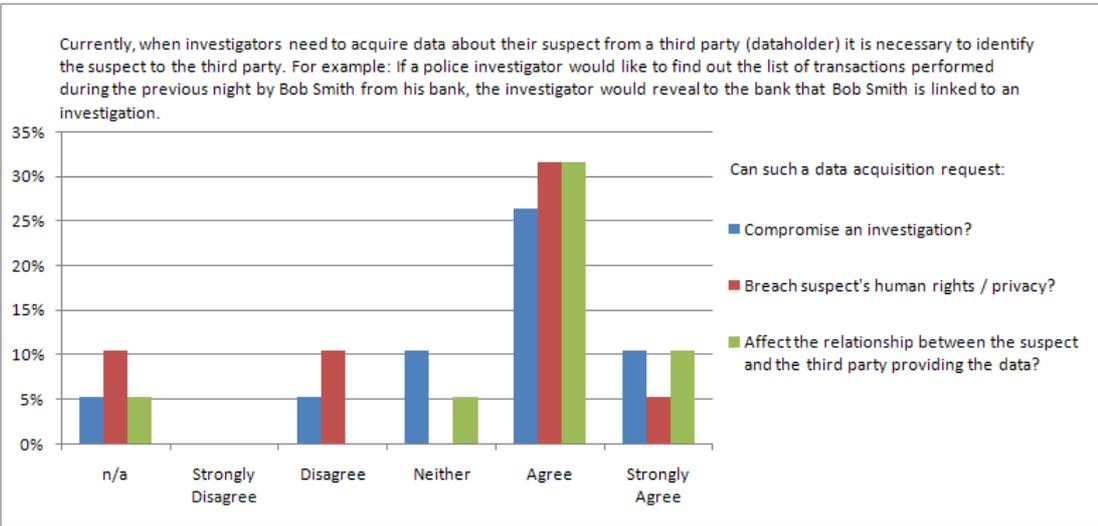
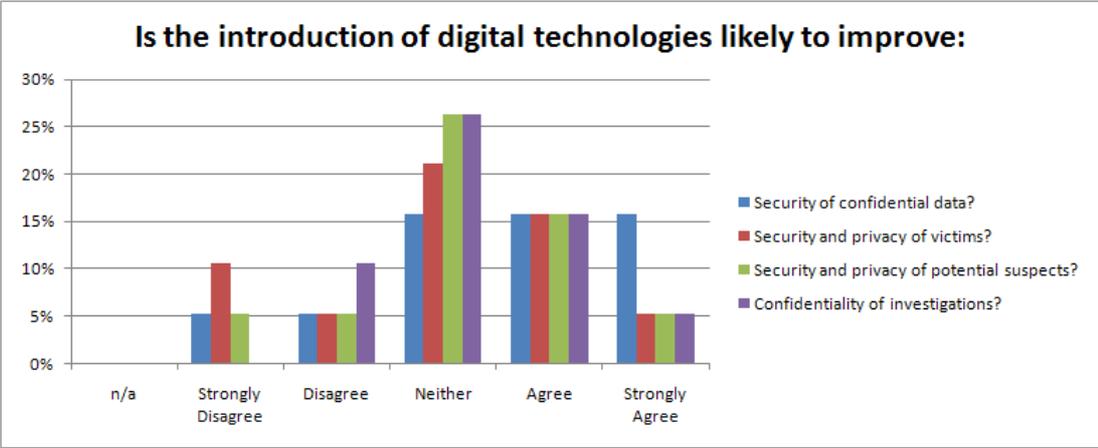
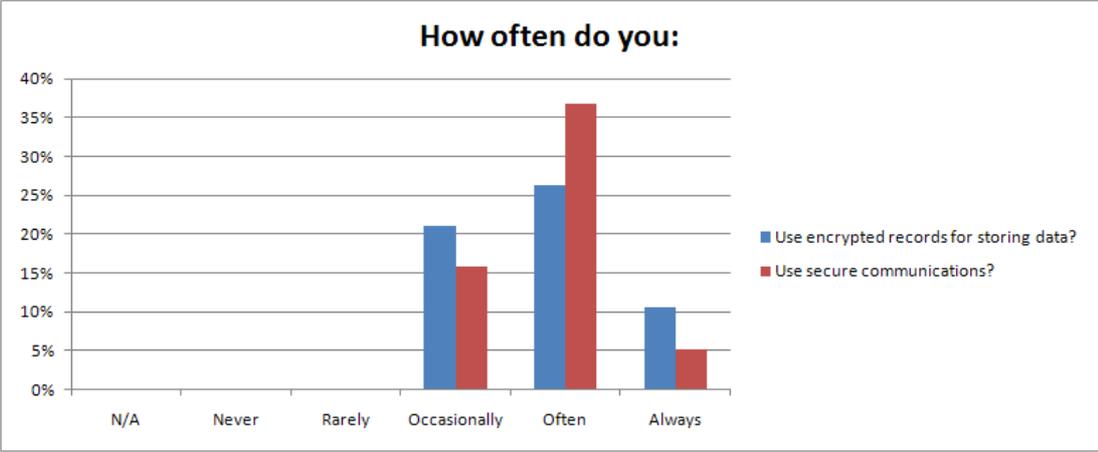
m	n	OT	OT	PE	PE
1	1000000	792299.1091	68955.09488	45329.19	45329.17
5	1000000	983619.4886	258019.4721	45329.87	45329.84
10	1000000	1222769.963	494349.9437	45330.71	45330.69
50	1000000	3135973.758	2384993.716	45337.47	45337.45
100	1000000	5527478.502	4748298.432	45345.91	45345.89
500	1000000	24659516.45	23654736.16	45413.48	45413.46
1000	1000000	48574563.89	47287783.31	45497.94	45497.92
5000	1000000	239894943.4	236352160.6	46173.62	46173.6
10000	1000000	479045417.8	472682632.1	47018.21	47018.19
50000	1000000	2392249213	2363326405	53774.97	53774.95

Results from the simulations are not included with this thesis, as it is possible to replicate them based on the complexity tables.

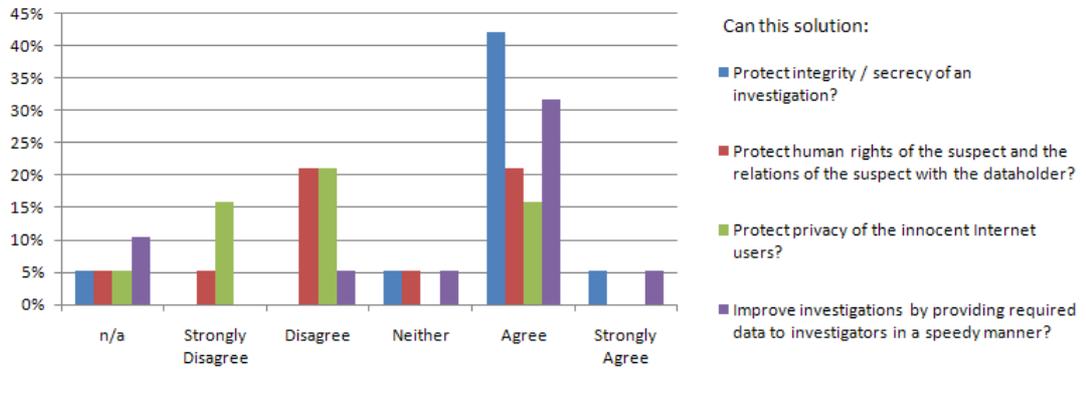
Appendix C

Survey

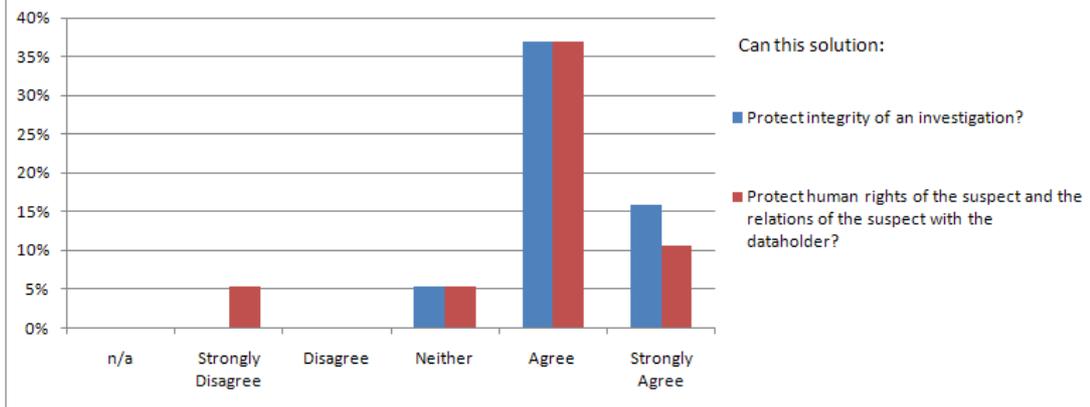




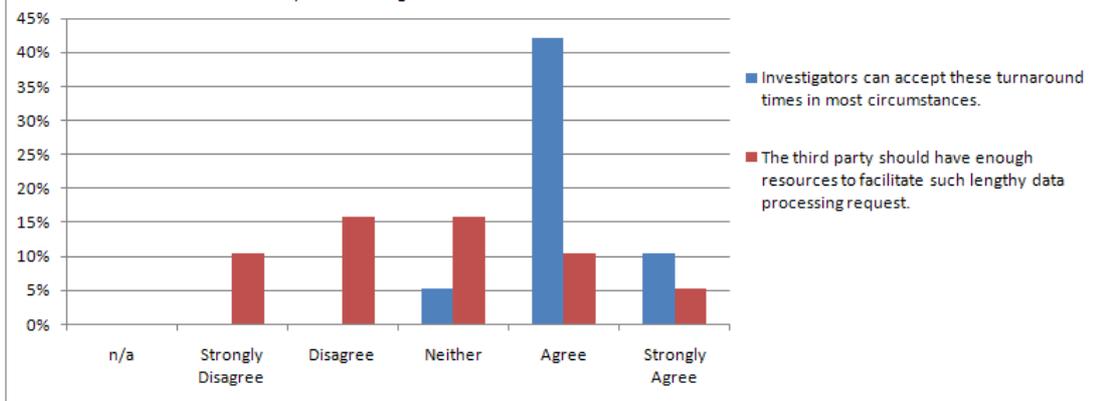
The government plans to introduce a centralised system allowing investigators access to data collected by Internet Service Providers (as well as some other Content Service Providers). In such a system investigators will be able to access required data fast and in a secure manner. Additionally the ISPs will not be able to monitor enquires made by the police.



We have developed a system that can be used by the investigators to hide the identity of the suspects during the data acquisition process from ISP or any other third party. Can hiding of the suspect's identity:

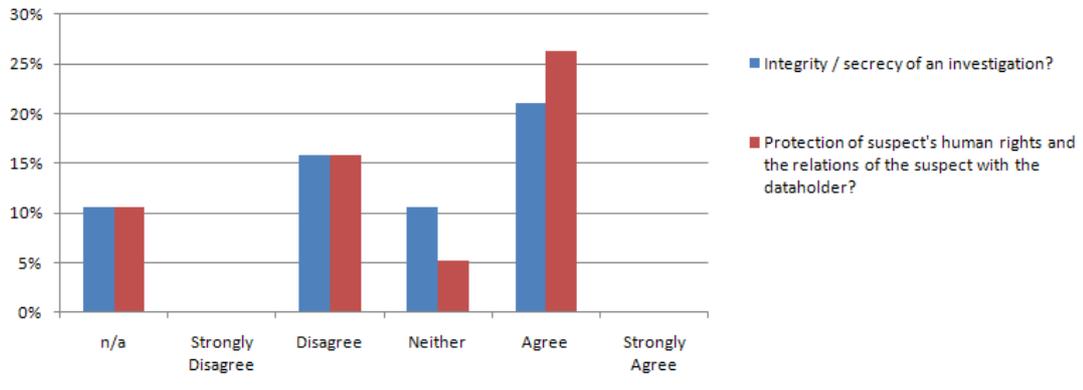


There are some drawbacks introduced by our system. The time required for the enquiry depends on the number of records in the third party database. For example: Data acquisition from a database of a thousand employee records would take approx. half a minute, whereas notice served to an ISP with around million of users can take approx. 500min (8 hours) to be processed by an ordinary computer. Do the benefits of this system outweigh this drawback?

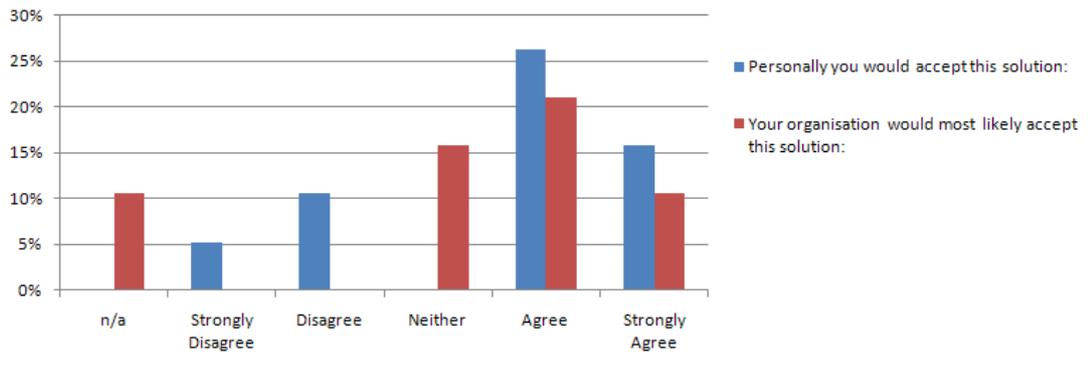


It is possible to lower the turnaround time for processing data from large databases by operating on the subset of the database. In this way the processing time can be minimised, however, it would be easier for the third party to identify the suspect. Still the probability of the third party finding out the identity of the suspect would be 1 in 1000 or 1 in 10000.

Do you think that this solution is acceptable from the perspective of:



In order for the privacy to be maintained the system encrypts a large chunk of the third party database with techniques similar to these used to transfer data securely over the Internet and transfers it to the investigators. From these the investigators are then able to extract only the data of the suspect specified on encrypted acquisition request and nobody else. It can be proven that it is not possible for the investigators to extract any other information from the third party database. Do you think that the benefits of the added privacy and security, as well as secrecy of the investigation outweigh this drawback?



Some critics worry that the government agencies are capable of breaking encryption. In order to guarantee to the public that the investigators will only have access to the data requested we propose to introduce another party into the system. An independent body, such as Information Commissioner's Office could filter the encrypted responses from the dataholders, passing on only the data requested by the investigators. This party would not be aware of the investigation context, ID of the suspect, or the content of the response from the dataholder. And as this party would have no intentions in finding these information, nor have the processing power to do so, it could be trusted not to attempt to break the encryption.

