



Technology-supported Risk Estimation by Predictive Assessment of Socio-technical Security

Deliverable 2.2.3

TRE_SPASS Technical Data Extraction Tools

Project: TRE_SPASS
Project Number: ICT-318003
Deliverable: D2.2.3
Title: TRE_SPASS Technical Data Extraction Tools
Version: 1.0
Confidentiality: Public
Editor: Mike Osborne
Cont. Authors: M. Osborne, A. Lenin, M. Ford, F. Reis,
M. Nidd, D. Hadziosmanovic, W. Pieters,
A. Tanner
Date: 2016-10-31



Part of the Seventh Framework Programme
Funded by the EC-DG CONNECT

Members of the TRE_sPASS Consortium

1. University of Twente	UT	The Netherlands
2. Technical University of Denmark	DTU	Denmark
3. Cybernetica	CYB	Estonia
4. GMV Portugal	GMVP	Portugal
5. GMV Spain	GMVS	Spain
6. Royal Holloway University of London	RHUL	United Kingdom
7. itrust consulting	ITR	Luxembourg
8. Goethe University Frankfurt	GUF	Germany
9. IBM Research	IBM	Switzerland
10. Delft University of Technology	TUD	The Netherlands
11. Hamburg University of Technology	TUHH	Germany
12. University of Luxembourg	UL	Luxembourg
13. Aalborg University	AAU	Denmark
14. Consult Hyperion	CHYP	United Kingdom
15. BizzDesign	BD	The Netherlands
16. Deloitte	DELO	The Netherlands
17. Lust	LUST	The Netherlands

Disclaimer: The information in this document is provided “as is”, and no guarantee or warranty is given that the information is fit for any particular purpose. The below referenced consortium members shall have no liability for damages of any kind including without limitation direct, special, indirect, or consequential damages that may result from the use of these materials subject to any liability which is mandatory due to applicable law. Copyright 2016 by University of Twente, Technical University of Denmark, Cybernetica, GMV Portugal, GMV Spain, Royal Holloway University of London, itrust consulting, Goethe University Frankfurt, IBM Research, Delft University of Technology, Hamburg University of Technology, University of Luxembourg, Aalborg University, Consult Hyperion, BizzDesign, Deloitte, Lust.

Document History

Authors		
Partner	Name	Chapters
IBM	Mike Osborne	ALL
CYB	Aleksandr Lenin	ALL
CHYP	Margaret Ford	ALL
GMVP	Fatima Reis	ALL
IBM	Axel Tanner	3
IBM	Michael Nidd	3
TUD	Dina Hadziosmanovic	2, 3
TUD	Wolter Pieters	1, 2, 3

Quality assurance		
Role	Name	Date
Editor	Mike Osborne	2016-10-30
Reviewer	Cédric Muller	2016-10-15
Reviewer	Olga Gadyatskaya	2016-10-15
WP leader	Mike Osborne	2016-10-30
Coordinator	Pieter Hartel	2016-10-31

Circulation	
Recipient	Date of submission
Project Partners	2016-10-24
European Commission	2016-10-31

Acknowledgement: The research leading to these results has received funding from the European Union Seventh Framework Programme (FP7/2007-2016) under grant agreement no. 318003 (TRE_sPASS). This publication reflects only the authors' views and the Union is not liable for any use that may be made of the information contained herein.

Contents

List of Figures	iv
List of Tables	v
Management Summary	vii
1 Introduction	1
1.1 Overview	1
1.2 Objectives	2
1.3 Foreground and Background	2
1.4 Approach	2
1.5 Case studies	3
1.5.1 Cloud Infrastructures	5
1.5.2 IPTV in the UK	5
1.6 Vision	6
2 Data to be extracted	8
2.1 Levels of data	8
2.1.1 Global View	9
2.1.2 Case View	9
2.1.3 Model View	11
2.2 Audiences and purpose of data extraction	12
2.3 Legal aspects of data extraction	14
2.4 Data extraction	16
3 Data Discovery and Extraction Process	18
3.1 Case data — based on the Cloud Case Study	18
3.2 Data discovery and extraction	19
3.2.1 Global View	19
3.2.2 Model View	19
3.2.3 Case View - based on the Cloud case study	19
3.2.4 SAVE	20
3.2.5 VMEXTRACT	20
3.2.6 Case View - based on the IPTV case study	20
3.3 Data categorisation	21
3.3.1 Asset categorisation	22
3.3.2 Vulnerability categorisation	23
3.3.3 Attack categorisation	24
3.3.4 Attacker categorisation	25

3.3.5	Social and policy data categorisation	26
3.3.6	Security risk categorisation	26
3.4	Data transformation	26
3.4.1	Log data	27
3.5	Data verification and update	27
3.6	Output	28
4	Conclusions	30
	References	31

List of Figures

1.1	Data management architecture	4
2.1	The hierarchy of data views.	10
3.1	The high level TRE _S PASS process.	19

List of Tables

3.1 A classification of attackers based on motivation	26
3.2 A classification of attackers based on skill	26

Management Summary

Key takeaways:

- This deliverable outlines the purposes, scope, as well as the target audiences for data extraction. It identifies data entities to be extracted, categorises them into 3 views, where each subsequent view narrows down the scope of the previous one.
- Provides a high-level outline of legal aspects of data extraction, as well as some limitations and constraints that these aspects might impose, suggesting solutions to some of them.
- Outlines the steps of data analysis and integrity checking: categorisation, transformation, verification and evaluation. Provides a comprehensive outline of data classification possibilities, as well as defines the approach to the data update cycle.

This deliverable describes the framework within which to handle the TRE_SPASS data management process specifically for technical data, with a focus on data categorisation and legal aspects. Corresponding technical work is covered in the deliverables [The TRE_SPASS Project, D2.2.2 \(2015\)](#) and [The TRE_SPASS Project, D2.4.1 \(2016\)](#).

Through the exploration of the key audiences, purpose and scope of data extraction in a risk context, the nature of the relevant data is identified in order to provide high quality input to the processes developed in other work packages, such as modelling (WP1), analysis (WP3) and visualisation (WP4).

The deliverable identifies the technical data entities which are of interest for the extraction process, as well as introducing the concept of three broadly hierarchical views of the data: global, case and model views. The global view takes in the entirety of the data under consideration. The case view only includes data relating to a single specific environment. Finally, the narrowest view is the model view, which encompasses only those aspects of the case data which are relevant to the TRE_SPASS modelling process. A comprehensive description of data classification options is developed, along with an approach to the handling of the data update cycle.

In this document, we have used the IPTV and the Cloud case study as environments for describing the case view, although the case view does not always have to be as wide-ranging as it is in the IPTV case.

As it constitutes an essential part of any consideration of data handling, we also include a high-level outline of legal aspects of data extraction, as well as some limitations and constraints that these aspects might impose, suggesting possible solutions where appropriate.

Having laid these foundations, we then describe the TRE_SPASS data discovery and extraction process. This process includes a number of steps:

- discovery of potential data sources
- extraction
- analysis and integrity checking
- storage
- output

1 Introduction

1.1 Overview

The role of WP2 within TRE_SPASS is to identify the structures and processes required to achieve effective handling of data in support of risk management. This includes the gathering, processing and output of data. This data includes published data from sources such as industry bodies and governments. However, to ensure practical applicability, data is also sourced from the case studies described in ([The TRE_SPASS Project, D7.1.1, 2013](#)). The data identified within WP2 has been used to support the development of the models in WP1 and further for model checking (WP3) and visualisation (WP4).

The focus of this deliverable is specifically on the final description of the tools used for extracting data in technical environments.

We have investigated technical data in more traditional environments, where virtualisation is not yet widely deployed in ([The TRE_SPASS Project, D2.2.1, 2013](#)). We have investigated Cloud-based infrastructures in ([The TRE_SPASS Project, D2.2.2, 2015](#)) delivered in M36. This report updates the previous two deliverables with the efforts in the last 12 months, including insights from the cloud case study ([The TRE_SPASS Project, D7.2.2, 2016](#)).

One of the goals of automated data extraction in these environments is to obtain, in a fast and easy way (ideally event driven and dynamic) infrastructure data, like physical and virtual hosts, network and storage connectivity, as well as access control of the users of the virtualised infrastructure. This is relevant to include it in the TRE_SPASS model for a cloud as basis for risk analysis in the larger context.

For efficiency, scalability and cost-effectiveness, cloud infrastructures have central management capabilities, independent of the specific underlying virtualisation technology.

In the specific context of this prototype, we work with VMware¹ as the virtualisation environment. VMware has a management component called vSphere ([VMware, 2015](#)) that allows to centrally administer and operate all components of a virtualised infrastructure.

The data extractor for virtualised environments accesses the vSphere administration component via an official API to extract the information and is described in detail in ([The TRE_SPASS Project, D2.2.2, 2015](#)).

¹VMware <http://www.vmware.com>

1.2 Objectives

The overall objective of this deliverable is on the extraction of technical data as part of the final technical data management process: data discovery, extraction and analysis, to provide input for the TRE_SPASS socio-technical security models. This process describes how the relevant data is acquired and transformed into the right format for input to the models. The objectives are as follows:

- Outline the target audiences for whom data will be extracted.
- Describe and characterise the data types relevant to TRE_SPASS risk calculations.
- Describe the technical data discovery and extraction process.

The aim of this deliverable is to focus on the management aspects of technical data. Key features to be considered within the cloud case study include the hardware, software, configuration and connectivity. Another essential consideration is secure data storage, which is managed within a cloud environment.

Another important area to be considered, which potentially spans both technical and social aspects of a system, is the implementation of access control. Both technical controls and human processes are essential for achieving effective risk management in this respect.

1.3 Foreground and Background

WP2 specifically focuses on the requirements for data handling in support of risk management within the TRE_SPASS project. The aim of the work package is to adopt best practice in managing data, in order to facilitate the innovations being undertaken in related work packages. These include the modelling work being undertaken by WP1, analysis by WP3 and visualisation by WP4.

The priority in this deliverable is to document the final data management structures with a suitable categorisation, as well as related legal issues. In particular, it builds on the approach developed in ([The TRE_SPASS Project, D2.1.1, 2013](#)) to incorporate industry best practice into the TRE_SPASS data management processes.

The data categorisation in the phases of the data extraction steps is considered foreground, as well as the technical tools (like SAVE and VMEXTRACT) that are described in more detail in other deliverables.”

1.4 Approach

The TRE_SPASS approach considers a socio-technical system model (navigator map) and a separate attacker model. The separation of the system model from the attacker model is intended to support a greater level of clarity and flexibility in the risk assessment process,

enabling both aspects to be evaluated individually before contributing to a more integrated view.

As previously identified in the requirements in ([The TRE_SPASS Project, D2.1.1, 2013](#)), it is essential to define the purpose and intended audience prior to undertaking any data extraction. In particular, where data is intended for more than one purpose, the appropriate data types and abstraction levels should be identified for each audience prior to gathering the relevant data. It is also essential that the protection to be afforded to the data on collection, processing and output is considered in advance of initiating any data collection. For the specific details relating to the WP7 case studies we refer the reader to ([The TRE_SPASS Project, D8.4.2, 2016](#)), which contains a detailed register of the types of risk arising from the handling of sensitive data in each of the individual case studies.

The technical data extraction prototype is an important step required to identify and develop the means to extract technical data suitable for supporting risk calculations from an organisation, in line with the TRE_SPASS approach.

The overall architecture of the data management processes is shown in [Figure 1.1](#).

The data management processes gather data derived from sources which are applicable beyond an individual case (global data) e.g. public databases of security information, as well as data relating only to a particular environment (case-specific data). Based on the requirements specification of the navigator map and attacker models, this data is then transformed into the appropriate format for input into the different elements of the modelling process. Under normal circumstances, attacker models will typically form part of the global data. While different attackers may be drawn to different environments, according to their interests and skill sets, it is unlikely that any individual type of attacker will be drawn only to a single specific case.

The navigator map is specific to an individual case (although some elements may be reusable, as discussed in relation to model sharing and reuse in [The TRE_SPASS Project, D5.3.1 \(2013\)](#)). When constructing a navigator map for the case, the user can draw on the data gathered from the case environment, in order to support the model development process.

1.5 Case studies

In order to investigate specific scenarios to offer a practical focus and results to draw on, the project consortium has decided to focus on the IPTV case study in the first instance. This particular case study was selected for the richness and diversity of its social and technical factors, enabling the project tools to be developed in relation to a practical yet challenging scenario. There is also a significant relationship between the IPTV case study and the Cloud case study, in that the IPTV implementation is likely to be dependent on third party cloud implementations. This provides the opportunity to investigate these elements from both the technical and social perspective, with outcomes which will extend significantly across the Cloud case study. Any relevant outcomes may also be used and validated in relation to the Telco case study, where appropriate.

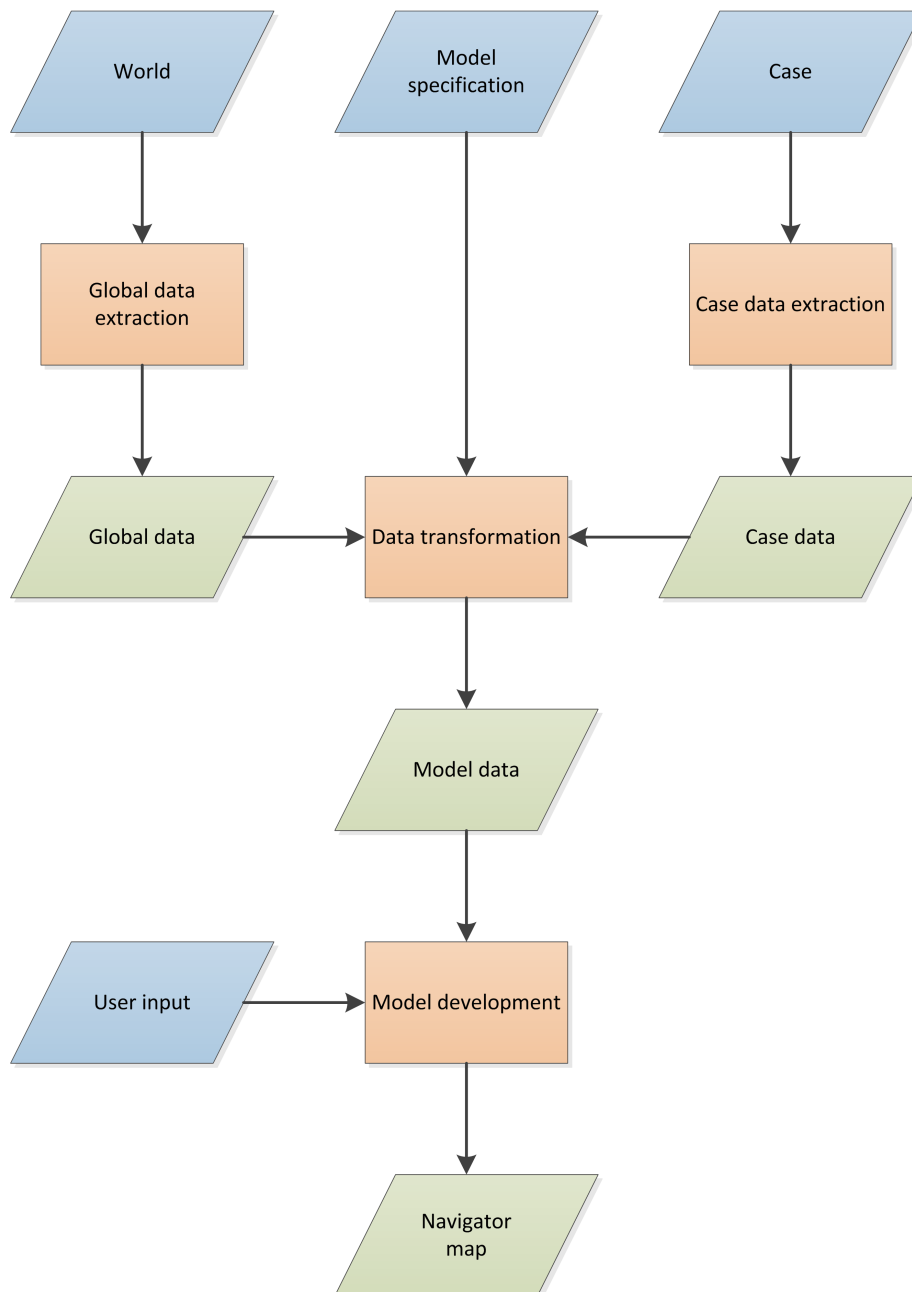


Figure 1.1: Data management architecture

The project is working with the London Rebuilding Society (LRS), a UK-based social enterprise, to test the evolving TRE_sPASS processes and tools by applying them to their home-based payments and money management system, delivered via IPTV set-top box, which is currently under development. The aim of the connected TV system is to enable people who are on low incomes, older people or people with disabilities to be able to manage their own finances securely from within their own homes. In particular, it aims to give them a greater level of autonomy by enabling them to allocate specific funds for

particular purposes (jam-jarring) and also to allocate a particular amount of money for a particular purpose at a particular time to a secondary cardholder. This might include paying by card for groceries at the local supermarket, or withdrawing cash at an ATM.

The development of the system has wide-ranging social aims, including providing people with a greater level of autonomy and flexibility in managing their financial affairs. Another important aspect is to enable users of the system to gain access to the types of services and deals which are only available online. This is an essential feature of the system, since lack of access to online offers and deals can compound the exclusion which people may experience as a result of not being online. This can have profound social, practical and financial implications.

The IPTV case study is described in some detail in deliverable D7.1.1, along with the other case studies. An initial risk assessment was undertaken in line with the CHYP SRA process, in order to identify the range of risks to which the system might be vulnerable. In order to investigate the case study further, the risks associated with the IPTV case study were mapped to create an attack tree. This provided a very comprehensive outline of the types of attack to which the system might be vulnerable. The features of this tree are explored in more detail in ([The TRE_SPASS Project, D3.3.1, 2013](#)).

The creation of the attack tree was also valuable in identifying the processes which reappeared frequently within the same tree. In some instances, these processes could potentially have been applicable across both IPTV and Cloud case studies (for example compromise of network equipment). This provided the opportunity to develop processes for handling frequently occurring situations (for details we refer the reader to ([The TRE_SPASS Project, D5.3.1, 2013](#)), describing model sharing approaches) and to explore the strengths and weaknesses of the existing modelling approach.

The attack tree structure was chosen on the basis that it is a format widely used and recognised across industry and could be used to interface with other processes as necessary. This would not necessarily be the preferred form of presentation to the user of the eventual TRE_SPASS tools, but provides a well-understood background means of communication in the first instance.

1.5.1 Cloud Infrastructures

The background of the Cloud case study is described in [The TRE_SPASS Project, D7.2.2 \(2016\)](#). The data extraction tool for virtualised environments is presented in detail in [The TRE_SPASS Project, D2.2.2 \(2015\)](#).

1.5.2 IPTV in the UK

The majority of homes in the UK already have the necessary equipment to access the Internet via their television (http://www.deloitte.com/view/en_GB/uk/industries/tmt/media-industry/perspectives-on-the-uk-tv-sector). In many cases, this capability will be via a games console, which may only be used for a relatively narrow range of

activities. Although many new TV sets have integrated networking capabilities, these tend to be at the higher end of the market, making them unaffordable for people on low incomes.

Many of the major TV equipment manufacturers continue to regard their devices as the 'trusted companion' in the living room and this attitude is reflected in the restrictive approach to content in this market. Where it may be acceptable to access the full range of Internet content via computers, tablets and phones, the TV is often still seen as a special case. In particular, the major manufacturers have a strong focus on child protection and avoiding any perceived association between their products and the less respectable side of the Internet.

Manufacturers aim to achieve large volumes in selling their equipment, regarding 10 million shipments as a minimum for a new device. While this market has historically been very conservative, it is gradually becoming more open over time. The requirements for a set top box providing Internet connectivity are fairly basic: support for some kind of I/O device (normally a remote control), network connectivity, processing capability and video output. As an initial platform for development of the IPTV home payments system, the raspberry pi has been identified as providing these features at low cost.

The newer generation of set top boxes, such as the inexpensive Roku streaming boxes (<http://www.roku.com/>) are also potentially of interest, as they provide a much more open platform than many others on the market, have similar capabilities to the raspberry pi, but are hardened in order to support secure provision of content.

Historically, security concerns relating to set top boxes have focused on gaining or reselling access to premium content delivered via the set top box (http://irdeto.com/documents/ds_security_lifecycle_services_en.pdf). However, as the set top box becomes an increasingly standard Internet-connected device, it is appropriate that this focus should shift to encompass the wider range of information security risks.

1.6 Vision

The focus of this document is on the audiences for the data extraction process, the types of data to be extracted, potential methods for data discovery, extraction and classification, with reference to the IPTV case study where appropriate.

By identifying the data required for input to the technical system model, we can move towards a practical implementation of the TRE_SPASS vision and concepts.

In this document we propose the architecture for the data extraction process and describe a prototype designed to demonstrate the functionalities of the proposed approach. Where possible, existing software has been selected for use in the data extraction process. This choice has been made for a number of practical reasons:

- The immediate utility of software which has been proven in the field.
- The associated time saving which allows the project to focus on innovations in the areas of modelling and visualisation.

- The likelihood of wider industry acceptance of TRE_SPASS methods and tools which incorporate existing industry standard software.

The products of the current deliverable are as follows:

- A specification of the data discovery and extraction process.
- A specification of data types needed for TRE_SPASS risk calculations.
- A categorisation of analysed data types.
- Design and implementation of the data extraction prototype.

2 Data to be extracted

For the purposes of WP2 we are currently working with four classifications of data: *(i) physical* (e.g. hardware, locations), *(ii) digital* (e.g. software), *(iii) processes* (e.g. organisational procedures and policies), *(iv) social* (e.g. characteristics of human behaviour). Within the scope of the current task, we consider only *physical* and *digital* data, which in practice relates to the digital infrastructure of an organisation. This includes networks, devices, software and configuration details.

The relation between different classes of data represents a particularly challenging and interesting area for investigation in a socio-technical study of this kind. While *social* data forms an integral part of the TRE_SPASS approach, it will not be considered in any great detail in this deliverable. Research relating to *social* data will be described in greater depth in ([The TRE_SPASS Project, D2.3.1, 2014](#)), due for delivery in M24 of the project.

Depending on the context, different types of data need to be extracted. Firstly, we consider different levels of data with respect to the environment from which the data is extracted (Section 2.1). Secondly, we consider types of data with respect to the audience (i.e., users) that require that data (Section 2.2). We discuss legal implications of extracting data (Section 2.3) as well as constraints and limitations of automated data extraction in (Section 2.4).

2.1 Levels of data

In this task two key inputs guide the process of data extraction: *(i)* the model definition from WP1 (described in [The TRE_SPASS Project, D1.3.1 \(2013\)](#)) and *(ii)* the TRE_SPASS case studies (described in deliverable [The TRE_SPASS Project, D7.1.1 \(2013\)](#)). The model definition specifies a way to model the behaviour of an organisation in the context of the TRE_SPASS project (e.g. decide which variables will be taken into account to describe actors and calculate threat risk). A case study describes a real-world environment that will be used to validate the TRE_SPASS model. The main goal of this deliverable is to identify suitable techniques for data extraction that can accommodate the requirements of both the defined model and the current case study. This is challenging because: *(i)* each case study is different (e.g. consider IPTV as compared to the Telco environment; some data sources required for IPTV will not be relevant to a Telco study) and *(ii)* the TRE_SPASS model is expected to evolve over time.

The process of data extraction in this deliverable mainly concerns the IPTV case study. However, since the extraction process needs to be able to support the context of the other

two case studies (Cloud and Telco), we aim to create an environment that will have the flexibility to accommodate the requirements of these other case studies as well.

To understand the relationship between the data available through case studies and model definitions, we use the notion of data views. We consider three views on the data:

Global The full range of technical and infrastructure data that could be of value as input to the risk modelling and visualisation processes. This includes systems, networks, software and configuration information. This view is independent of any particular case.

Case Technical and infrastructure data that is available within the context of a specific case (e.g. inventory and configuration data).

Model Technical data required for modelling and analysing the organisation in the TRE_SPASS model (e.g. actors and actions).

In Figure 2.1 we visualize the relationships between the different data views.

The *global* view of the data is the superset of the *case* and *model* data. Also, for each case the *case* data must be a superset of the *model* data (since otherwise it would be infeasible to use the case within the tool).

2.1.1 Global View

The global view includes all the technical data that could potentially be extracted and processed from a digital environment for the purposes of risk assessment. This takes a wide view of the range of possible data, without being specific to any individual organisation. For examples as well as for possible data sources we refer the reader to Chapter 5 of ([The TRE_SPASS Project, D2.1.1, 2013](#)).

The purpose of this view is to provide a context within which the requirements of the various case views can be compared and combined, with the aim of identifying the most effective approaches to data collection in different environments.

2.1.2 Case View

The case view contains the subset of the global view data, relevant for the individual case in question. This may be any specific scenario identified by an organisation. For the purposes of TRE_SPASS, this will include the case studies being developed by WP7. For the purpose of this document, we will focus mainly on the IPTV case study.

The case view includes data about the system: the hardware, software, configuration and connectivity data relating to the system implementation. In the context of TRE_SPASS, it is not practical for reasons of privacy, cost and scale to hold some of the more sensitive data (for example personal customer information and payments data) in a central repository. The kind of highly secure, high availability architecture normally required for this purpose would require an exceptionally high degree of maintenance for a project of this kind. For this

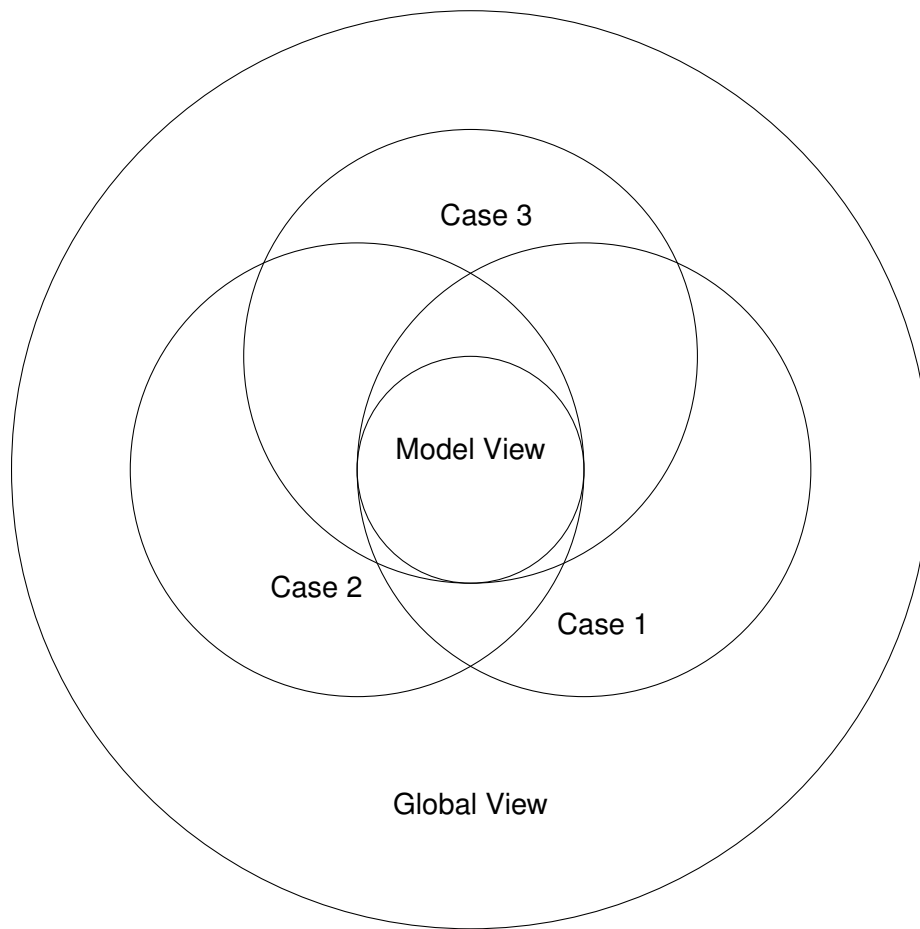


Figure 2.1: The hierarchy of data views.

reason, data held by the project will not include the most sensitive personal data of system users, but rather generalised data derived from data handled in this type of environment.

While we can identify the types of technical data which will need to be handled by the system with a reasonable degree of confidence at this stage, the deeper social investigations undertaken in preparation for ([The TRE_SPASS Project, D2.3.1, 2014](#)) will inevitably influence our understanding of the relationship between the technical systems and the people involved with them (in a range of roles such as end user, support technician, attacker).

The types of system data which would be required to investigate the IPTV system are described in some detail in ([The TRE_SPASS Project, D7.1.1, 2013](#)) and explored further through the attack tree developed in ([The TRE_SPASS Project, D3.3.1, 2013](#)). Details of the handling of sensitive data in relation to the case studies are included in ([The TRE_SPASS Project, D8.4.1, 2013](#)).

A range of system data will be of interest when investigating the potential risks associated with a specific environment, as mentioned above:

- Hardware e.g. model details of set top box and components, including known vulnerabilities and countermeasures.
- Software e.g. operating system, application software, browser software (if application is browser-based), software versions and patch levels, including known vulnerabilities and countermeasures.
- Configuration e.g. system settings, implementation of accounts and access policies e.g. separation of duties through access control.
- Connectivity, to include both type and capacity of connection e.g. permanent wireless broadband connection of 8Mbps.

It is important to identify which data elements for the target case study, the IPTV system, will be applicable beyond this individual scenario. This is an essential feature to identify in relation to activities explored in other work packages e.g. model sharing in WP5. There may be significant risks and little value in sharing data elements which are specific to a single case. However, elements which are potentially applicable to a range of cases can usefully be shared more widely. Under most circumstances, these types of data will be more generic and therefore present a lower risk of exposing sensitive information when shared with other organisations.

The organisation-independent data to be extracted should include:

- Information about known vulnerabilities and security incidents available from public databases, relevant for the target system.
- Information about typical assets / asset groups for particular types of enterprise.
- Publicly available information about organisations which have been the subject of analysis.

The availability of organisation-specific data is limited by the lack of a live system in the case study at present, but the data to be extracted should include:

- Proposed network topology.
- Planned configuration of IPTV device.
- Asset / asset groups.
- List of known vulnerabilities (per asset). This will require the identification of published vulnerabilities associated with the specific assets held by the organisation.

2.1.3 Model View

The model view consists of the subset of the global view data sufficient to serve as input for the TRE_SPASS model for successful processing.

As described in Chapter 3 of ([The TRE_SPASS Project, D1.3.1, 2013](#)), the current TRE_SPASS model operates with the following data:

- Location domain

- Domain identifier
- Location
 - Location identifier
 - Connected locations
 - Risk of detection at this location
 - Domain location belongs to
 - Access policies (required and enabled components)
- Actor
 - Actor identifier
 - Behaviour or behaviour class
 - Role
 - Likelihood of social-engineering attack to be successful
 - Risk appetite of the actor
 - Goal of the actor (minimising risk of detection or time for actions)
- Action
 - Action identifier
 - Time to perform an action
 - Risk of detection when performing an action
 - Cost of performing an action
- Data
 - Credentials
 - Identity
 - Other data
 - Access policies (required and enabled components)
 - Owner

2.2 Audiences and purpose of data extraction

TRE_SPASS tools will have a number of audiences which should be considered. This section outlines the different aspects of data to be taken into consideration when designing the data extraction processes for these audiences.

In this section, we specifically identify the groups for whom the TRE_SPASS tools are likely to be relevant. There are, of course, also groups involved in the development of tools within

the other work packages, especially work packages 1, 3, 4 and 5 for whom these data will be particularly important. Every effort has been made to align the data requirements of these work packages with the activities of WP2 throughout the course of the project.

One of the most prominent audiences for the TRE_SPASS tools is the community of *security practitioners*. These may include CISOs, risk managers within organisations, security consultants and auditors. The purpose of data extraction in this context is to provide input to the TRE_SPASS models, with the aim that the resulting outputs and visualisations should provide professionals with a decision support tool, enabling them to save time and have increased confidence in the outcome when evaluating scenarios. Convenience, flexibility, transparency and repeatability are important qualities for tools to provide utility to professionals working in this type of environment.

For example, a security officer within a large banking organisation may be considering a number of possible measures to be implemented in response to possible risk scenarios which they have identified within their organisation. The gathering of data in relation to a limited number of those scenarios and the subsequent application of TRE_SPASS tools, could enable them to decide on the relative importance to their own operational continuity of implementing specific countermeasures. In this respect, the tools would not necessarily need to provide a specific probability for each scenario, but rather a relative likelihood, which would enable the security professional to identify immediate operational priorities for investment.

Another group, which has been identified as an important driver of growth in difficult economic times, is the *SME (Small and Medium-sized Enterprise)*. In this environment, it is highly unlikely that an individual risk or security specialist will be employed. However, small businesses are known to be vulnerable to cybercrime and eager to find ways to combat it. Given that companies of this size are unlikely to be in a position to hire specialist consultants, a less complex set of tools to evaluate the most prominent risks to their organisation could provide a valuable safeguard. These would need to be intuitive and easy to use, since the business owner or employee will have a limited amount of time to spend on this type of activity. Improving the level of resilience across SMEs has the potential to produce wider benefits by limiting the level of fraud across the economy as a whole, resulting in a healthier overall business environment.

Depending on the context, different audiences may have different requirements from the system. We discuss these requirements by using the concept of use cases. A use case describes the processes which an individual may employ in order to achieve particular goals. In the TRE_SPASS context, this implies the use of the TRE_SPASS tools for a defined purpose. A use case should not be confused with case studies, which, in this project, provide a particular target for validating the tool. We consider the following use cases (detailed descriptions can be found in [The TRE_SPASS Project, D6.2.1 \(2013\)](#)):

Security investment A security consultant uses the TRE_SPASS navigator to decide on effective investments in information security.

Audit A security auditor assesses whether an organisation complies with a certain security level.

Innovation An IT manager analyses the security weaknesses of a newly designed product.

Product-service system An IT manager analyses whether a newly designed product provides a business case for fraudsters.

For each use case we now briefly describe how the data extraction requirements relate to the practical goals of the specific use case. In this context, the aim of the TRE_SPASS tool is to gather and process information relevant for performing the task of the use case. Inputs to the TRE_SPASS tool, gathered during data extraction, provide information that is relevant for calculating risk.

During the first use case (*security investment*), the user needs the following knowledge to estimate the most promising security investments: the overall security level of the organisation, present threats, the most likely paths of malicious activity, etc. This knowledge can be calculated by using the inputs such as: information about asset costs (i.e., used as a component to calculate the impact of an attack), common access patterns (i.e., reveal the most likely penetration path and characterise a threat), configuration data (i.e., use the specifics of configuration to evaluate the vulnerability of the system).

During the second use case (*audit*), the user needs to evaluate the overall security level within the organisation. Various items of infrastructural information might contribute to this evaluation. For example, technical information such as the implementation of firewalls, encryption, network segmentation, as well as access control details including password policies and physical access control measures will provide important technical input. Active scanning of the system (i.e., to check for the existence of vulnerabilities) and other data gathering techniques can also provide a more detailed understanding of the environment being audited.

During the third and fourth use cases, the users need to analyse the security of the new product. For this, knowledge about the most common vulnerabilities and attack paths (e.g. from publicly available databases) can, together with technical product specifications, be used to evaluate the product in the early design phase.

2.3 Legal aspects of data extraction

Laws and regulations governing data storage and processing form a hierarchy, with international law at the highest level. In describing the regulatory environment, it is important to identify the competent authorities and the way in which they influence practice in a particular context. A number of principles have a degree of international recognition. For example, the protection of personal data from misuse and the appointment of a designated authority to control data handling. Frequently, the relevant legislation is implemented at country level on the basis of Directives agreed at a higher level.

The following conventions of international law must be respected in relation to data handling:

1. Convention for the protection of human rights and fundamental freedoms (<http://conventions.coe.int/treaty/en/Treaties/Html/005.htm>).

2. Convention for the protection of individuals with regard to automatic processing of personal data (<http://conventions.coe.int/Treaty/Commun/QueVoulezVous.asp?NT=108&CM=8&DF=7/24/2008&CL=ENG>).

EU directives on data processing provide an intermediate level in this hierarchy. For the full list of corresponding directives we refer the reader to http://ec.europa.eu/justice/data-protection/law/index_en.htm.

The most prominent include:

- 95/46/EU – protection of individuals with regard to automatic processing of personal data.
- 2002/58/EU – processing of personal data and protection of private life in the ICT sector.

Any processing of personal data requires adherence to the relevant legislation, including Directives as transposed into local law at the country level.

The EU regulations place the following obligations on those responsible for data collection, storage and processing:

- Data collection is legal and fair only with the explicit consent of the person concerned.
- The person whose data is collected and processed has to be aware of which data is collected and why, as well as how it will be processed.
- Individuals must have the opportunity to query the data stored about them at any time. Data processing facilities must provide individuals with full access to data held about them.
- The data processing entity must stop data collection and processing immediately upon the request of the person concerned (individuals must have the opportunity to request that no more data be processed or stored about them).
- There must be a clearly defined data disclosure policy, where all the persons who get access to data (e.g. employees for business duties) need to be explicitly listed.
- There must be a strict access control policy, describing the minimum required controls for those wanting to gain access to the data.
- There must be a strong non-repudiation mechanism. Each read/write operation must leave a unique trace to allow identification of the person who accessed the data and the time when it was accessed.
- Appropriate levels of confidentiality, integrity and availability must be adhered to, minimizing the risk of activities such as data tampering and disclosure. Steps must also be taken to mitigate the insider threat.

The EU legal framework on the protection of personal data is currently undergoing significant reforms. Developments in this area, as well as the proposed Directive on Network and Information Security should inform the decisions of the project in relation to data handling. In particular, the “General Data Protection Regulation” (http://europa.eu/rapid/press-release_IP-12-46_en.htm) is expected to supersede EU Directive 95/46/EC in

due course. This is a potentially significant change, as Regulations of this kind have immediate legal force and do not require transposition into law at the member state level.

Contrastingly, EU Directives are only written into law at the national level and may have significantly different outcomes in different jurisdictions. For example, due to the difference in transposition of the Directive, the UK and Ireland do not treat IP addresses as personal data, although most other EU member states do. Due to the diversity of implementations at national level, it is not practical for the TRE_SPASS project to tailor its processes to the requirements of individual countries.

Moreover, tools for data extraction can have legal and indeed criminal implications. In 2001 the Council of Europe developed the “Convention on Cybercrime” (ETS-No.: 185, <http://conventions.coe.int/Treaty/Commun/QueVoulezVous.asp?NT=185&CL=ENG>) which was passed by the council in 2004. Many members and non-members have already adopted the convention into their legal systems. While Germany was in the process of adopting it in 2007 (§ 303b StGB) a widespread discussion arose in the security community regarding the definition of so-called “Hacker-Tools”. Both development and publication of these would be prohibited. There was significant controversy over the dual-use character of most security tools like password crackers or network scanners for penetration tests. The interpretation of this ambiguous law later became more distinct with the help of journalists who turned themselves in, publicly confessing to publishing “Hacker-Tools”. In this context, we must bear in mind that our data extraction tools could be interpreted in some legislations as malicious attack tools.

In summary, for the data storage/protection aspects of the TRE_SPASS tools we should take both international law and the EU directives and regulations into consideration. More detailed legal requirements will need to be managed locally by those organisations who choose to implement the TRE_SPASS tools. In certain cases, adherence to national laws may add additional constraints and require additional security measures and precautions to be adopted.

2.4 Data extraction

In this deliverable we are focusing on the development of a workflow for data extraction. Under particular circumstances, a more automated approach may be desirable, however it will be necessary to take a range of issues and constraints into consideration when making design choices of this kind. Some of the possible issues and constraints are listed here:

Increased complexity. Solutions for automated data extraction will make the TRE_SPASS tools more independent and self-sustained, but at the cost of increased system complexity, which adds additional limitations on deployment and maintenance.

Additional points of failure. With increased complexity, use of solutions for automated data extraction introduces additional points of failure, which in turn increases maintenance costs and requires more highly skilled staff.

New threats and attack vectors. Automation itself introduces potential threats relating to data integrity and the risk of data leakage:

Spoofing A malicious actor on the network may send malformed output to the scanning facility thus providing the TRE_SPASS tools with false results.

DoS A malicious actor on the network may swamp the scanning facility with service requests, thereby cutting input to the TRE_SPASS tools.

Data leakage Automated data acquisition requires devices to communicate to the scanning facility and reveal sensitive information. A malicious actor impersonating the scanning facility could potentially gain access to sensitive data.

3 Data Discovery and Extraction Process

The data discovery and extraction process is a workflow which encompasses several steps: data discovery, extraction, categorisation/classification, transformation, evaluation, storage and output. Data discovery involves the exploration of an environment in order to identify the range of available data. Extraction involves gathering the identified data into a designated location. Data characterisation aims at classifying the extracted information in clusters that will enable measurable calculation of risk concepts. Data transformation translates the extracted data into concepts that are suitable for calculating risk. Data evaluation performs checks on the integrity and quality of the extracted data.

We now describe possible data sources and then explain each step of the process in more detail.

3.1 Case data — based on the Cloud Case Study

This case study uses and evaluates TRE_SPASS tools in the context of the TRE_SPASS process. A high level diagram of the TRE_SPASS process is shown in the figure below. The standard TRE_SPASS process described by WP5 is split into 4 steps, Data Collection, Modeling, Scenario Analysis and Visualisation. For the purposes of highlighting the tools, and highlighting the integration of the very first activities ('Stage 0') we have split the data collection into two parts. The process is thus displayed as 5 steps as follows :

- Data Collection: Situational Analysis that is used to gather information in preparation for a risk analysis.
- Data Collection: Technical Analysis that is used to collect the necessary data for a TRE_SPASS based risk analysis.
- Modeling: where the technical and social data is brought together to model the target of the risk analysis.
- Scenario Analysis: where a particular scenario is applied to the system model and analysed.
- Visualization: where output of the analysis is visualized.

During the course of the TRE_SPASS project, a number of tools and process have been developed. A subset of the tools that were relevant and available for the cloud case we evaluated. This subset of tools is shown in the above figure categorised against the process steps. The remainder of this section introduces those tools that were used in one of the iterations of the cloud case study.

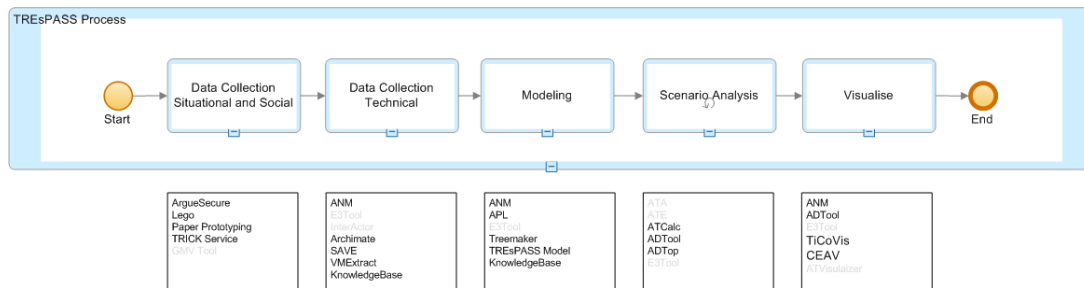


Figure 3.1: The high level TRE_SPASS process.

3.2 Data discovery and extraction

We now present techniques that can be applied to extract technical data from different environments.

3.2.1 Global View

As noted earlier, global data represents all data that can be extracted in an ideal setup of an organisation (e.g. having data available on network behaviour, logs, process policies, etc.). For the description of tools suitable for extracting the data we refer the reader to Section 6.5.1 “Mechanical Retrieval of Infrastructure Data” in ([The TRE_SPASS Project, D2.1.1, 2013](#)).

3.2.2 Model View

The model view comprises data that can be leveraged to build the TRE_SPASS model. More specifically, the current model definition includes data on location (e.g. connections to other locations, access policies), actors (e.g. role, likelihood of social engineering), action (e.g. actor, time to perform an action, cost of performing an action), etc. Different general-purpose tools can be used to extract the model data. The choice of the tool mainly depends on the conditions in the particular case organisation.

We use the definition of the model data as a guideline for choosing the extraction techniques in the specific case study.

3.2.3 Case View - based on the Cloud case study

The goal of automated data extraction in these cloud environments is to obtain, in a fast and easy way (ideally event driven and dynamic) infrastructure data, like physical and virtual hosts, network and storage connectivity, as well as access control of the users of the virtualised infrastructure. For efficiency, scalability and cost-effectiveness, cloud infrastructures have central management capabilities, independent of the specific underlying

virtualisation technology. This section describes the two tools that were developed as part of the work package 2

3.2.4 SAVE

Save is a data extraction and policy analysis tool that was developed as part of the EU FP7 TClouds project. The policy analysis required a simple network topology that the SAVE data collection engine built from information extracted from a number of cloud operating systems. During the TRE_SPASS project the data collection engine was extended to capture a richer set of information better suited to the requirements of the TRE_SPASS modelling language developed in WP1. The extended data extraction capabilities, in particular the ability to capture a consistent snapshot of a virtualized infrastructure were later transferred to an IBM product. During this case study these capabilities were evaluated during an internal IBM process described later in this document.

3.2.5 VMEXTRACT

In the specific context of this prototype, we work with VMware as the virtualisation environment. VMware has a management component called vSphere that allows to centrally administer and operate all components of a virtualised infrastructure.

The data extractor described here accesses the vSphere administration component via an official API to extract the information.

The overall goal of WP2 is to extract physical, social, and IT infrastructure data to feed the TRE_SPASS model. Task T2.2 especially is concerned with the extraction of technical data. This deliverable describes the process to extract technical data from a given, virtualised IT infrastructure in an automated fashion.

As IT infrastructures, especially today's cloud infrastructures, can be very large and are highly dynamic, such an automated process is an important aspect of data extraction for WP2, most specifically in connection with the task T7.2: Case Study A: Cloud infrastructure as critical infrastructure (see ([The TRE_SPASS Project, D7.1.2, 2015](#), Chapter 3)), lead by partner IBM.

This section describes a tool for data extraction from virtualised environments, here specifically for centralised VMware infrastructures. Given a user id with suitable access rights, allowing read access to all items of interest in the vSphere environment, the script can extract data and will keep it as an internal RDF representation. This RDF representation can be exported as XML that can be used to import information about the virtualised environment into the TRE_SPASS model of a cloud environment.

3.2.6 Case View - based on the IPTV case study

The IPTV case study represents a distributed computer architecture which is based on several locations (e.g. account holder's home, service provider premises). Due to the

fact that part of the system is located in a location that is outside the company premises (i.e. in the user's home), remote data extraction will have some limitations. For example, active scanning for vulnerabilities will not be the preferred option, as it could influence system performance. Since the IPTV software is implemented on Raspberry Pi, tools for general-purpose host-based analysis may not be an option (due to the embedded character of the Raspberry Pi). This however, should not influence the usage of network-based tools for data extraction. Also, the provider is planning to use the cloud service. This somewhat limits (and simplifies) the analysis of network communication in the system.

We now briefly summarize the basic tools that are used for data extraction:

- Vulnerability scanning and management solutions (e.g. Nessus¹, OpenVAS²)
- Tools to gather information about operating systems, running software, open ports, etc. (e.g. Nmap³)
- Tools to collect systems configurations and settings. I.e. using the Simple Network Management Protocol (SNMP) it is possible to collect a wide range of configuration information. Use of dedicated vendor-specific tools and protocols may also be useful here.

Transactional data is highly sensitive, as it contains clients' personal data, e.g. first name, last name, physical address, property-related data, bank accounts and transactions. Opportunities to collect such data will greatly depend on the implementation of data storage and the design of the IPTV system. As the IPTV system is not live yet, we will consider manual data source discovery and input at this stage, with the potential for vendor-specific tools once the system goes live.

3.3 Data categorisation

Data categorisation aims at linking data from data sources to relevant features in the TRE_SPASS tools and processes, based on taxonomies. This data categorisation has important roles in: (i) data representation and (ii) interpretation of risk elements. Firstly, different categories of a specific data type will each have a tailored representation in the TRE_SPASS model. For example, the representation of different categories of one resource type (e.g., resource: asset; categories: *people*, *hardware*) will differ significantly (e.g., *hardware* will be described by technical aspects while *people* will be described by social aspects). Secondly, data categories are fundamental for establishing links to risk concepts, which will further be used for building the model. For example, different types of assets (e.g., hardware vs. people) will have different types of vulnerabilities (e.g., buffer overflow vs. phishing) that influence the total security risk, as well as different countermeasures / controls (e.g., patching vs. education).

Categorisation is important not only to provide appropriate input to the model, but also to support the development of visualisations, with a range of different views. Through the

¹<http://www.tenable.com/products/nessus>

²<http://www.openvas.org/>

³<http://nmap.org/>

course of the task we will define suitable categorisations. As an initial step, we considered the existing characterisations from the literature. Although there are many proposed taxonomies in the literature, there is no standard or universally accepted taxonomy for attacks, vulnerabilities and attacker profiles (e.g., see critics of various taxonomies in [Bishop and Bailey \(1996\)](#)). Depending on the purpose of the taxonomy, researchers focus on different aspects (e.g. extension to the application domain, organisation management, reporting). For a comprehensive survey we refer the reader to [Ilgure and Williams \(2008\)](#). To choose a suitable taxonomy for the TRE_sPASS context, we set two general requirements:

- The taxonomy needs to enable an intuitive link to risk concepts.
- The taxonomy needs to support the socio-technical application domain.

Firstly, since one of the fundamental goals of the TRE_sPASS project is risk evaluation, the chosen taxonomy of data sources has to support easy translation to risk concepts. For example, the chosen categories in data classification have to be measurable in relation to basic risk categories (e.g., evaluate impact of a DoS attack). According to [Bishop and Bailey \(1996\)](#), a taxonomy should assist in decision-making regarding resource investment. Therefore, we focus on the existing taxonomies that are designed with a specific application in mind: risk assessment.

Secondly, the selected taxonomy should be specific to the context (e.g., a taxonomy of vulnerabilities in operating systems is of little use when conducting a security assessment of a cryptographic protocol ([Ilgure & Williams, 2008](#))). Since TRE_sPASS is concerned with the socio-technical aspects of the system, the suitable taxonomy will focus on aspects like: system architecture, physical and digital infrastructure configuration, security policies, user behaviour, and update management.

We now present basic categorisations for different data sources that fit the stated requirements.

3.3.1 Asset categorisation

The International Organisation for Standardization ("[ISO/IEC 13335-1](#)", 2004) classifies organisational assets as:

- physical assets (e.g. computer hardware, communications facilities, buildings),
- information / data (e.g. documents, databases),
- software,
- the ability to produce some product or provide a service,
- people,
- intangibles (e.g. good will, reputation).

The categorisation of assets is relevant for two purposes. Firstly, different properties may be assigned to different types of assets in the model. For example, information can easily be copied, but physical assets cannot. Secondly, depending on the attacker profile, different types of assets may have different attractiveness to the attacker. Although the value of

the asset to the organisation can be estimated by the organisation itself, the value to the attacker will be based on a combination of asset properties and attacker properties.

3.3.2 Vulnerability categorisation

We distinguish two relevant concepts regarding vulnerability categorisation: vulnerability classification and vulnerability evaluation. A vulnerability classification allows us to characterise vulnerabilities in a systematic way. In the context of the TRE_SPASS project, classification is important as it helps to break down the system state from different perspectives (e.g. social, technical, physical). Vulnerability evaluation analyses the severity of a vulnerability. Evaluations are important as they provide comparable outputs which can be used to calculate risk concepts. We now present some common classification and evaluation approaches.

Most of the vulnerability classifications in the literature relate to software vulnerabilities. For the TRE_SPASS context, we need a more general classification framework. We now present basic taxonomies for classifying a vulnerability. A general taxonomy by Howard (John D. Howard, 1998) is a well accepted, and often cited classification that introduces three main categories of vulnerabilities:

- implementation,
- design,
- configuration.

As a more extensive taxonomy, Bishop and Bishop (1995) propose six axes for classifying a vulnerability:

- nature: describing a type of flaw according to "Protection analysis" categories,
- time of introduction: when the vulnerability was introduced,
- exploitation domain: what is gained through the exploitation,
- effect domain: what can be affected by the vulnerability,
- minimum number: the minimum number of components necessary to exploit the vulnerability,
- source: the source of identification of the vulnerability.

A common methodology for vulnerability evaluation is the Common Vulnerability Scoring System (CVSS) (<http://www.first.org/cvss/>). This approach uses three different metrics to evaluate the severity of a vulnerability: (i) base, (ii) temporal and (iii) environment. Each type of metric has different aspects for evaluation. First, **base** metrics evaluate the fundamental characteristics of a vulnerability: *access vector* (measures how close an attacker must be to attack the target), *access complexity* (measures the complexity of the attack required to exploit the vulnerability), *authentication* (measures the number of times an attacker must authenticate to exploit the vulnerability), *confidentiality impact* (measures the impact on confidentiality of a successful exploit of a vulnerability), *integrity impact*

(measures the impact on integrity of a successful exploit) and *availability impact* (measures the impact on availability of a successful exploit).

Second, **temporal** metrics measure time-dependent qualities of a vulnerability, such as *exploitability* (measures how complex it is to exploit a vulnerability in the target system), *remediation level* (measures the level of an available solution that can mitigate the problem), *report confidence* (measures the degree of confidence in the existence of the vulnerability)

Third, **environment** metrics measure the environment-specific qualities of a vulnerability: *collateral damage potential* (measures the potential for a loss of life or physical assets), *target distribution* (measures the percentage of systems which are vulnerable), *security requirements* (measures the loss with respect to basic security requirements– availability, integrity, confidentiality).

3.3.3 Attack categorisation

Several works base the attack categorisation on either *impact of attack* or *attack consequences* (Lindqvist & Jonsson, 1997; Hansman & Hunt, 2005). Here we present the taxonomy proposed by Lindqvist (Lindqvist & Jonsson, 1997). The authors consider two “dimensions” to categorise attacks: technique and result. According to the technique, the attacks are categorised as follows:

- bypassing intended controls,
- active misuse of resources,
- passive misuse of resources.

Each dimension is broken down further. For example, *bypassing intended controls* can imply three types of attacks: password attacks, spoofing privileged programs or utilising weak authentication. In addition, password attacks can refer to: password capturing or password guessing. *Active misuse of resources* is further divided into two subclasses: exploiting inadvertent write permission (e.g. write to files that the user is normally not supposed to edit) and resource exhaustion (e.g. perform a denial-of-service attack). *Passive misuse of resources* implies an unauthorised use of read permissions (e.g. a user accesses documents that he is normally not supposed to read).

According to the result, the attacks are categorised as:

- exposure,
- denial of service,
- erroneous output.

The category *exposure* can be further divided into (i) disclosure of confidential information and (ii) service to unauthorised entities. The category *denial of service* can be divided into (i) selective and (ii) unselective. The category *erroneous output* implies attacks which cause discrepancy between the actual status and what is shown (e.g., on user interface or through network communication).

The Common Attack Pattern Enumeration and Classification database (CAPEC, <http://capec.mitre.org/>) is a free publicly-available community developed database of attack patterns along with a comprehensive schema and classification taxonomy. The latest release consists of over 400 attack patterns classified in several views (e.g., attack mechanism, standard abstraction). CAPEC actively leverages other standard databases such as Common Weakness Enumeration (<http://cwe.mitre.org/>), Malware Attribute Enumeration and Characterization (<http://maec.mitre.org/>). This is valuable in that it enables easy and consistent data sharing and extraction. In ([The TRE_SPASS Project, D5.3.1, 2013](#)) we outline our use of CAPEC as a taxonomy for attack classification and sharing. However, CAPEC has mainly been developed for classifying attacks which target software. It is therefore necessary, in order to provide support for the socio-technical conditions of our context (e.g., social relationship, physical barriers), for an extension to the original scheme to be developed. We will investigate practical approaches to achieving this greater functionality through the course of the project.

3.3.4 Attacker categorisation

The characteristics of the attacker are relevant for calculating the likelihood of attack realisation (e.g., consider the attacker skill level and system vulnerability to calculate penetration likelihood). The most general categorisation of the attacker profile considers three aspects ([Raymond, 2000](#)):

- access,
- resources,
- activity.

Based on the access capabilities, an attacker can be external or internal. Based on type of resources, an attacker can be *(i)* static (an adversary chooses resources before starting the attack and cannot change the resources during the course of the attack) and *(ii)* adaptive (an adversary can change resources during the course of the attack). Based on the activity, an attacker can be active or passive. An active adversary can arbitrarily modify the data (e.g., tamper message), while a passive attacker can only listen.

Further, the attacker can be classified according to different aspects such as motivation ([Table 3.1, NIST \(2002\)](#)), skillset ([Table 3.2, Meyers, Powers, and Faissol \(2009\)](#)), etc. For TRE_SPASS, we aim at extending common taxonomies to fit the context. More specifically, we consider the following aspects: resource, skill, social context. First, we characterise attackers based on two aspects of resource: *(i)* the character of the resource (e.g., time, money) and *(ii)* the type of use of the resource (e.g., static or adaptive). Second, we foresee using attacker classification with different skill levels, such as those in [Table 3.2](#). Finally, to accommodate the socio-technical conditions, we plan to include the characteristics of the social relationship (e.g., the attacker's relationship to the victim). This will be further discussed in ([The TRE_SPASS Project, D2.3.1, 2014](#)).

The process of calculating quantitative parameters of risk with respect to attacker profiling was undertaken. For further details we refer the reader to Section 3.2.3 “Multi-parameter computations with attacker profiling” of ([The TRE_SPASS Project, D3.3.1, 2013](#)).

Table 3.1: A classification of attackers based on motivation

Attacker profile	Hacker	Computer criminal	Terrorist	Industrial espionage	Insider
Motivation	challenge, rebellion	destruction of information, monetary gain	blackmail, revenge	competitive advantage	prestige, ego

Table 3.2: A classification of attackers based on skill

Attacker profile	Script kiddies	Hacktivist	Cyber punk	Insider	Coder	White hat hacker	Black hat hacker	Cyber terrorist
Skill	very low	low	low	moderate	high	high	very high	very high

3.3.5 Social and policy data categorisation

The categorisation of social and policy data will be discussed in task 2.3 and ([The TRE_SPASS Project, D2.3.1, 2014](#)), due in month 24.

3.3.6 Security risk categorisation

[The TRE_SPASS Project, D5.1.1 \(2013\)](#) has identified The Risk Taxonomy of The Open Group ([The Open Group, 2009b](#)) based on the Factor Analysis of Information Risk (FAIR) framework ([The Open Group, 2009a](#)) as the basic conceptual framework for the risk management context. FAIR describes the factors which form risk and how they are related to each other, data types required to represent them and where to get the data from. One key element is the combination of defensive properties and attacker properties to estimate risk. For example, a highly skilled attacker will have a higher probability of success in an attack step than a less skilled one. Alternatively, a highly skilled attacker will need fewer resources to achieve the same probability of success. In their simplest form, such estimations will be based on a look-up table, as in the FAIR framework. For example, a highly skilled attacker trying to exploit a highly difficult vulnerability will have a medium probability of success. For a broader overview of the FAIR taxonomy we refer the reader to ([The TRE_SPASS Project, D5.1.1, 2013](#)), ([The TRE_SPASS Project, D3.1.1, 2013](#)) and Chapter 4 of ([The TRE_SPASS Project, D6.2.1, 2013](#)).

3.4 Data transformation

Data transformation represents the process of translating the extracted data into concepts that are suitable for calculating risk. Task 2.4 focuses on the transformation process. For completeness, in this deliverable we only briefly present the concept and illustrate the process through a few simple examples.

In our context, the transformation process defines the relationship between *case* data and *model* data. In particular, by performing this transformation we translate data extracted

from a particular case into the appropriate format for input to the model. Here are some examples of such transformations:

- transform the output of network reconnaissance tools into model data (e.g., parse *traceroute*, *nmap* output to extract *actors* and *locations* of the system),
- aggregate communication statistics to characterise *connections to other locations*,
- aggregate activities per actor in network or log data to calculate *time to perform an action*,
- extract different functionalities from log data to enumerate *actions*.

3.4.1 Log data

The aim of gathering log data is to have access to detailed, synchronised records of actions taking place across the environment under investigation. These actions relate specifically to the requirements for input to the TRE_sPASS model described in ([The TRE_sPASS Project, D1.3.1, 2013](#)). Logging records should be transmitted and stored securely, for reasons of both confidentiality and integrity. A time synchronisation protocol such as ntp is also essential, to ensure that when logs from disparate sources are combined, they can be interpreted in a meaningful way.

Depending on the nature of the individual case, it will be necessary to determine different levels of completeness required for the data extracted. This may vary according to the security levels required within the environment, the quantity of data generated and the quality of network links. The interfaces, update mechanisms and frequency of updates will need to be identified for each data source, or equally a standard approach may be adopted across a particular environment, or even for particular types of environment.

It is envisaged that log analysis will form a particularly important element of the telco case study, as this focuses very specifically on measures to alleviate the risk of financial loss due to fraud. This will require the identification of incidents, interpretation of system logs, troubleshooting and tracing of the state of the system over time. Transaction logs will also provide an important tool for measuring how much any individual has used specific services.

3.5 Data verification and update

There must be a verification process to check the integrity of the data extracted. Data evaluation and output will determine the relevance, applicability and quality of the extracted data, which will be reflected in the data output. Data storage is one of the output aspects that will be determined by evaluation of the extracted data and will depend on:

- The security classification attached to the data.
- The confidentiality, integrity and availability requirements for the data.

- The architecture of the TRE_SPASS tools and models.

Unfortunately, any particular instance of the model can be expected to be incorrect and incomplete (to some degree) by the time analysis is run. In TRE_SPASS we support an update cycle that can iteratively improve the completeness and accuracy of the model, allowing the data gathering effort to focus on those elements that can be expected to offer the greatest value. This intelligent update cycle can be supported by augmenting the data gathered with metadata qualifications of data *quality* and *volatility*. For now, the exact units or measurement scale for these values remain a point for discussion, but the principles are clear.

The first pass of data collection will almost certainly be an interview, in which people familiar with particular aspects of the target environment will describe hiring practices, building layout, and so on. The notes from this discussion represent the first iteration of the model, and can be described as having fairly low quality. A traditional security reviewer might look at the results of this first meeting, notice a likely attack vector, and ask to visit a particular room to improve both the quality and the completeness of data relevant to that vector.

The equivalent cycle in TRE_SPASS is for the data to be collected, modelled, and analysed. If an important attack vector is identified, but relies on data that is considered to be of low quality, then that discovery might justify a more thorough study of the components involved. It would similarly be useful to know which attacks were not feasible because they were blocked by a model component with a low quality rating. Extending the analysis to accomplish this may not currently be practical; further research in this direction is required.

Volatility would also inform an update cycle for the data in the model. High volatility data is different from low quality data because the measurement can be considered reliable for the time of measurement. An example of a volatile measurement might be the identity of contracted cleaning staff who may be replaced without notice due to vacation or illness. How to differentiate between volatile data and statistically distributed data is not yet clear. For example, it might be more likely for a door to be left open during the summer than during the winter; should this be modelled as a seasonally-volatile chance that the door is left open, or should it be modelled as a less-volatile probability curve that includes the date as a parameter?

Data verification, integrity checking, and revision is outside the scope of the current deliverable and will be addressed in ([The TRE_SPASS Project, D2.5.1, 2016](#)), due in M48.

3.6 Output

Being the final phase of the data extraction process, data output is responsible for providing all the other work packages and models with the necessary operational data. Data is stored in the WP2 TRE_SPASS Information System and may be queried for extraction and transmission using http-based API calls.

The usage of an API adds an abstraction layer between the data and the underlying storage and enables data to be queried without the need to reveal the underlying structure of the data storage in use.

Ideally, the data output should enable the data to be queried in the global view and also any of its subsets. However in the scope of this deliverable we outline the output of the model view specifically, as the data in the model view is the input data for the TRE_SPASS model and it is required in the first instance for validating the operation of the TRE_SPASS toolchain.

The TRE_SPASS model is designed in such a way that the expected input format is a single file containing all relevant data. The format and structure of the input file are described in ([The TRE_SPASS Project, D1.3.1, 2013](#)).

Details of the TRE_SPASS Information System can be found in ([The TRE_SPASS Project, D2.4.1, 2016](#)).

4 Conclusions

This deliverable presented a framework for the TRE_SPASS technical data extraction tools and a proof-of-concept tool suite for the cloud case study. We discussed the data extraction process, data analysis, data categorisation, and data storage. The prototype is based on the use of a limited set of data extraction tools as input for the TRE_SPASS navigator maps.

The work performed in this deliverable follows to a large extent the requirements specified in [The TRE_SPASS Project, D2.1.2 \(2015\)](#). For example, **IR2.7** requires the identification of audiences and purposes of data extraction prior to the extraction. By discussing the purpose and suitable audiences, we address this requirement in Section [2.2](#) of the current deliverable. **IR2.9** requires that the extraction methods must be designed with the appropriate terms and conditions of an organisation. We address this requirement by introducing the notion of *case view* in Section [2.1](#) of the current deliverable and the legal aspects of the data in Section [2.3](#). The *case view* on the potential data considers the conditions present in the specific organisation to adjust the extraction process, corresponding to **IR2.8**.

Our discussion of data classification and categorisation fulfils requirements **IR2.10**, **R09** and **R41**.

References

- Bishop, M., & Bailey, D. (1996). *A critical analysis of vulnerability taxonomies* (Tech. Rep.). Davis, CA 95616-8562: University of California at Davis. Retrieved from <http://seclab.cs.ucdavis.edu/projects/vulnerabilities/scriv/ucd-ecs-96-11.pdf>
- Bishop, M., & Bishop, M. (1995). *A taxonomy of UNIX system and network vulnerabilities* (Tech. Rep.). Davis, CA 95616-8562: University of California at Davis.
- Hansman, S., & Hunt, R. (2005). A taxonomy of network and computer attacks. *Computers & Security*, 24(1), 31-43. Retrieved from <http://dblp.uni-trier.de/db/journals/compsec/compsec24.html#HansmanH05>
- Igure, V., & Williams, R. (2008). Taxonomies of attacks and vulnerabilities in computer systems. *Commun. Surveys Tuts.*, 10(1). Retrieved from <http://dx.doi.org/10.1109/COMST.2008.4483667> doi: 10.1109/COMST.2008.4483667
- ISO/IEC 13335-1 [Computer software manual]. (2004).
- John D. Howard, T. A. L. (1998). *A common language for computer security incidents* (Tech. Rep.). Eubank, Albuquerque, United States: Sandia National Laboratories.
- Lindqvist, U., & Jonsson, E. (1997). How to systematically classify computer security intrusions. In *Proceedings of the 1997 IEEE Symposium on Security and Privacy* (pp. 154–). Washington, DC, USA: IEEE Computer Society. Retrieved from <http://dl.acm.org/citation.cfm?id=882493.884387>
- Meyers, C., Powers, S., & Faissol, D. (2009). *Taxonomies of cyber Adversaries and Attacks: a Survey of Incidents and Approaches* (Tech. Rep. No. 419041). 7000 East Ave, Livermore, CA 94550, United States: Lawrence Livermore National Laboratory. Retrieved from <https://www-eng.llnl.gov/pdfs/taxonomies.pdf>
- NIST. (2002). *Risk management guide for information technology systems, sp 800-30* (Tech. Rep.). Gaithersburg, Maryland, United States: National Institute of Standards and Technology.
- Raymond, J.-F. (2000). Traffic analysis: Protocols, attacks, design issues, and open problems. In H. Federrath (Ed.), *Workshop on design issues in anonymity and unobservability* (Vol. 2009, p. 10-29). Springer. Retrieved from <http://dblp.uni-trier.de/db/conf/diau/diau2000.html#Raymond00>
- The Open Group. (2009a). FAIR (Computer software manual No. C081). Retrieved from <http://pubs.opengroup.org/onlinepubs/9699919899/toc.pdf>
- The Open Group. (2009b). Risk taxonomy (Computer software manual No. C081). Retrieved from www.opengroup.org/pubs/catalog/c081.htm
- The TRE_SPASS Project, D1.3.1. (2013). *Initial prototype of the socio-technical security model*. (Deliverable D1.3.1)
- The TRE_SPASS Project, D2.1.1. (2013). *Initial requirements for the empirical models*. (Deliverable D2.1.1)
- The TRE_SPASS Project, D2.1.2. (2015). *Final requirements for the empirical models*. (Deliverable D2.1.2)
- The TRE_SPASS Project, D2.2.1. (2013). *Technical data extraction prototype*. (Deliverable D2.2.1)

- The TRE_SPASS Project, D2.2.2. (2015). *Data extraction from virtualized infrastructures*. (Deliverable D2.2.2)
- The TRE_SPASS Project, D2.3.1. (2014). *Social data and policy extraction prototype*. (Deliverable D2.3.1)
- The TRE_SPASS Project, D2.4.1. (2016). *TRE_SPASS information system*. (Deliverable D2.4.1)
- The TRE_SPASS Project, D2.5.1. (2016). *TRE_SPASS information testing and degradation tools*. (Deliverable D2.5.1)
- The TRE_SPASS Project, D3.1.1. (2013). *Initial requirements for quantitative analysis tools*. (Deliverable D3.1.1)
- The TRE_SPASS Project, D3.3.1. (2013). *First report on stochastic analysis methods*. (Deliverable D3.3.1)
- The TRE_SPASS Project, D5.1.1. (2013). *Initial requirements for process integration*. (Deliverable D5.1.1)
- The TRE_SPASS Project, D5.3.1. (2013). *Abstraction levels for model sharing*. (Deliverable D5.3.1)
- The TRE_SPASS Project, D6.2.1. (2013). *Initial refinement of functional requirements*. (Deliverable D6.2.1)
- The TRE_SPASS Project, D7.1.1. (2013). *Initial requirements for implementation of case studies*. (Deliverable D7.1.1)
- The TRE_SPASS Project, D7.1.2. (2015). *Final requirements for implementation of case studies*. (Deliverable D7.1.2)
- The TRE_SPASS Project, D7.2.2. (2016). *Final report case study a*. (Deliverable D7.2.2)
- The TRE_SPASS Project, D8.4.1. (2013). *Initial case study register*. (Deliverable D8.4.1)
- The TRE_SPASS Project, D8.4.2. (2016). *Final case study register*. (Deliverable D8.4.2)
- VMware. (2015). *VMware vSphere*. <http://www.vmware.com/products/vsphere-operations-management>. (Online; accessed 2015-10-26)