# 1. Research

## 1.1. Web Crawler

Research has been carried out into improving automatic web crawling. An analysis showed a number of data extraction problems with existing techniques for the kinds of websites relevant to SCIIMS (including online classified advertisements, social networks and job websites). An improved algorithm has been produced with beyond state of the art features. Subsequent experiments have shown that this achieves a success rate in excess of 85%.

Research has also been carried into improved automated Web browsing. Most previous web automation systems use conventional browsers to automate navigation sequences. This is computationally expensive in terms of both CPU and memory, especially with the new breed of websites (Ajax, etc.). This causes scalability problems where many browsers may be executing at the same time, which is the typical case in web automation applications. The approach is based on the insight that, in the case of automated navigation, browsing sequences are known in advance. Therefore it is possible to carry out a test execution to find what elements in the website pages are needed to execute the target navigation sequence and then only load and execute the required elements for subsequent executions. This has proved to be very effective, for instance automated web navigations with Microsoft Explorer are over 5 times slower than with the specialised browser. Figure 1 shows an example of this: to correctly process a click on the greyed A node, only the greyed nodes are needed (detecting the required nodes involves a complex process described in other documents). The other nodes/scripts need not be loaded or executed, resulting in a significant decrease of memory and CPU usage.
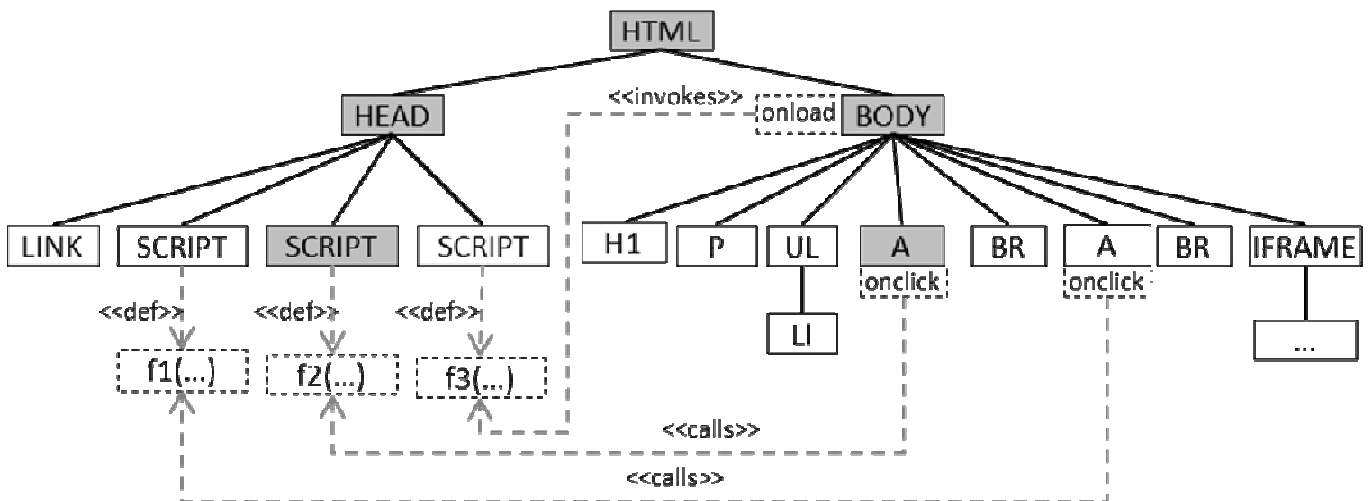


**Figure 1 SCIIMS Web Navigation**

## 1.2. Data Mining

Research has been carried out into Entity Resolution (ER) which covers the problem of identifying distinct representations of real-world entities in heterogeneous databases. This examined the trade-off between the storage needs, performance, and the efficiency of Entity Resolution. Consequently several research papers have been published. An example of this is the paper Infrastructures and Bounds for Distributed Entity Resolution which includes Figure 2 showing the execution time for different ER methods and number of records.



**Figure 2 Entity Resolution Execution Time Against Size Entity Resolution**

Research has been carried out into techniques for the visualisation of heterogeneous data sets as a network of entity nodes with arbitrary connections. As a result improved and new tools for visual analytics have been produced. An example of this from the Prototype/Demonstration System (described in Section 3) is shown in

Figure 3 .

**Figure 3 Graphical View Example**

## 2. Dissemination

### 2.1. Newsletter and Poster

A newsletter and poster have been produced alongside the web site (www.sciims.eu) to communicate the work carried out on the project to interested parties (see Figure 4).



**Figure 4 Newsletter and Poster**

## 3. Prototype/Demonstration System

### 3.1. Modules and Basic Architecture

A working Prototype/Demonstration System has been produced to demonstrate the capabilities of SCIIMS and also for system level experiments. Figure 5 shows the modules and basic architecture of the Prototype/Demonstration system.



**Figure 5 SCIIMS Demonstration System Components**

Supporting this are the:

a) Data Services Layer (Module 5) which combines, integrates and transforms information obtained from the different data sources for use by the SCIIMS modules. This includes data virtualization technology for RDF / OWL which is beyond state of the art,

b) Integration Module (Module 17) which deals with the asynchronous aspects (for instance Data Mining and Web Crawling). To support the Service Orientated Architecture (SOA) an Enterprise Service Bus (ESB) is used.

Not shown on this diagram is the Netica tool which is used to construct graphically Bayesian Belief networks as necessary to support an investigation (see

Figure 6 )

## 3.2. Organising a Criminal Investigation

A number of techniques have been identified for providing computer-assistance to an Investigator and combined into a conceptual system of thinking.

The originating techniques are:

a) **Pirolli-Card** (P-C) – a cognitive model that shows how an investigator typically conducts an Investigation

b) **Goal Structured Notation** (GSN) – a general information-representation technique that can be used in SCIIMS to show a hierarchy of findings, linking evidence to actionable conclusions

c) **Assessment of Competing Hypotheses** (ACH) – an investigative technique that helps to avoid making premature conclusions. The Investigator is supported in maintaining several hypotheses that explain some finding, only excluding those explanations when evidence contradicts them.

Bayesian Belief Networks **(BBN) – a probabilistic technique used to reason with a system of related hypotheses. In SCIIMS it is used to process the hypotheses of an ACH when the ACH would otherwise be too difficult to consider mentally (see**
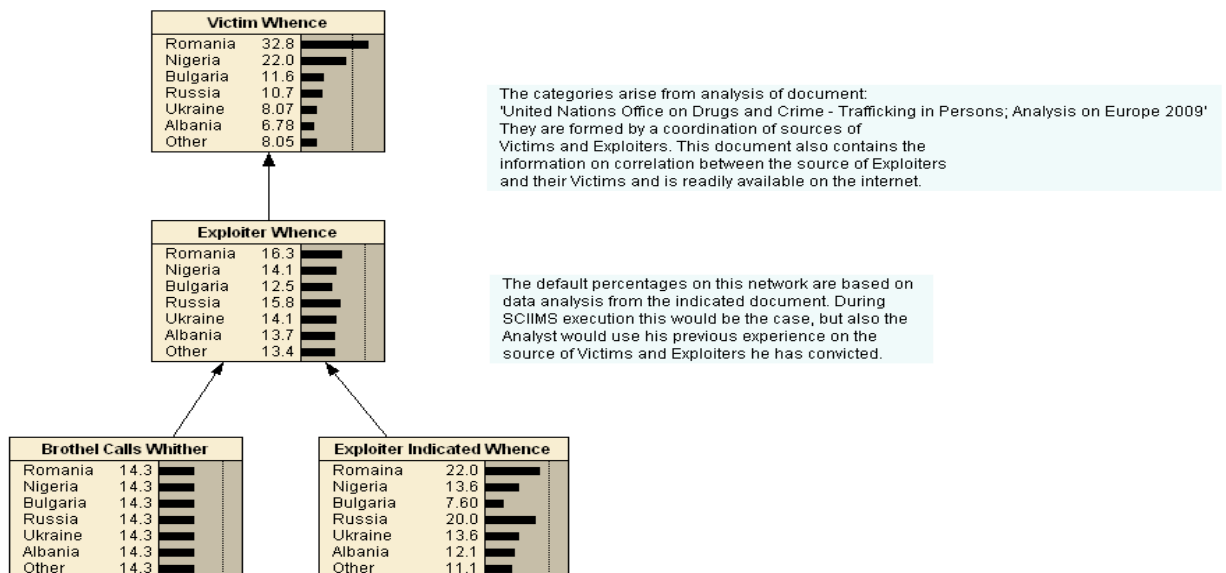
d) Figure 6)



**Figure 6 Example BBN**

SCIIMS allows a simple investigation to be shown and then enhanced using GSN, ACH and BBN techniques as the need arises. This is depicted in
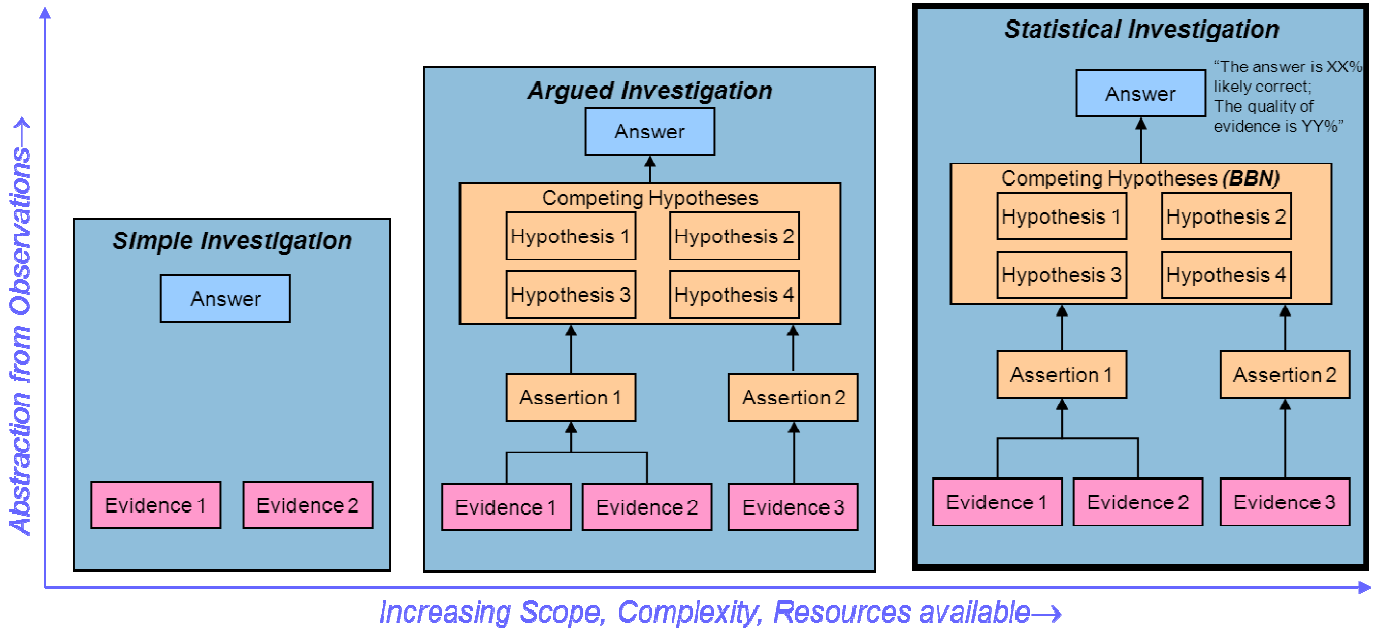
Figure 7 below.



**Figure 7 Increasingly Complex Investigations**

Figure 8 shows how an argument can be represented on the SCIIMS Prototype/Demonstration System using GSN.
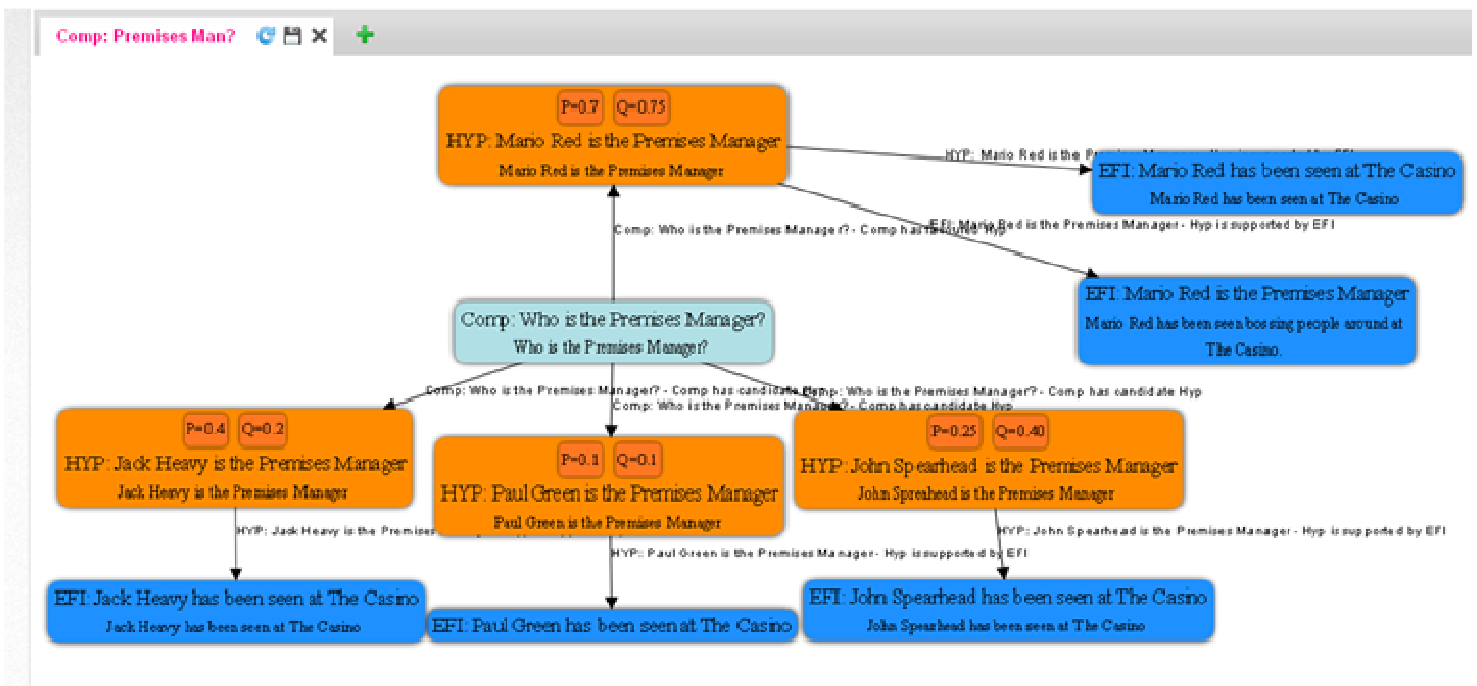
**Figure 8 Example of an Analyst using the ACH Investigation Structure in the Ontology**

## 3.3. Ontology

Central to the SCIIMS Prototype/demonstration System is the Ontology which has the following major sections (see

Figure 9):

a) Operational - aligned with the Adversaries Model and containing persons, places, relationships, events etc.

b) Process – for the Pirolli and Card business process, user identification etc.

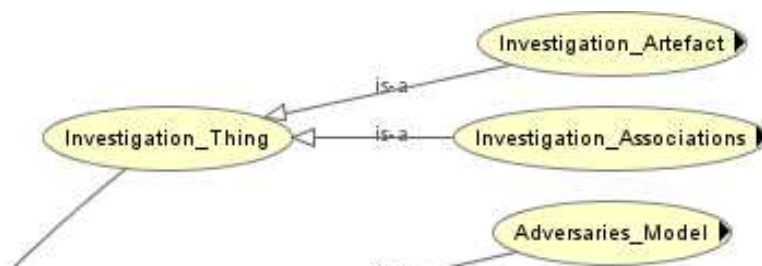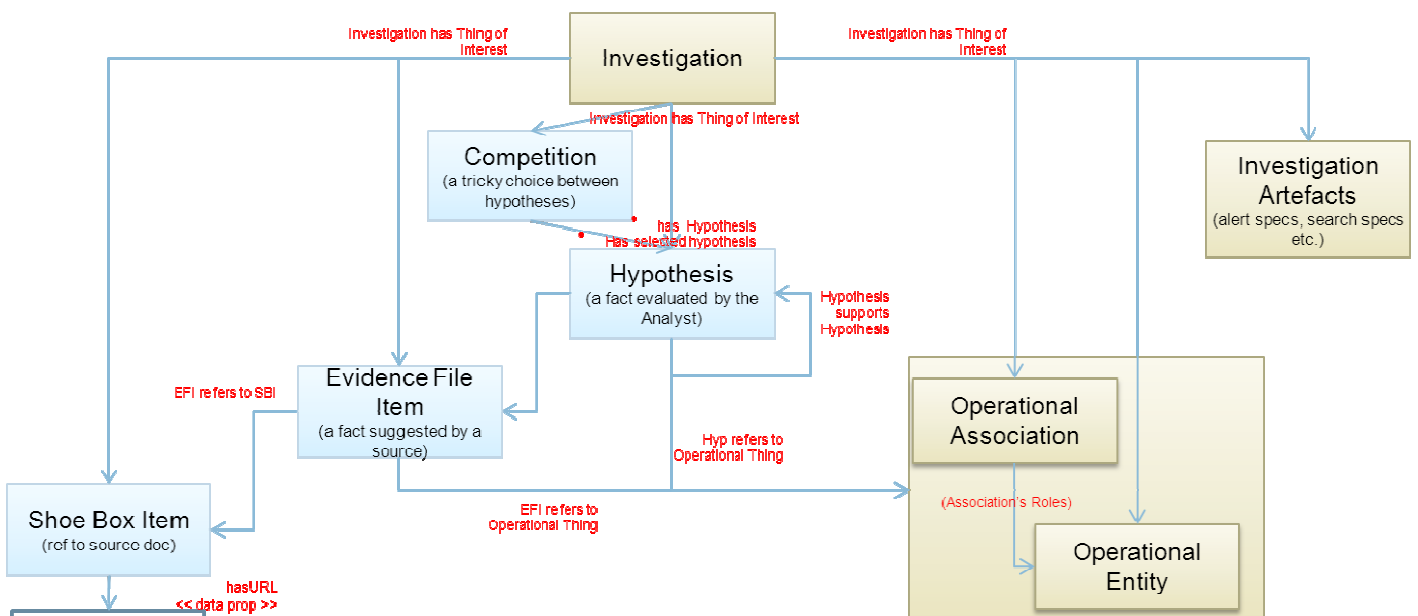c) Investigation – for shoebox data, searches, alerts etc.

d) HCI

e) BBN

**Figure 9 SCIIMS Ontology**

Figure 10 illustrates how an investigation is structured within the Ontology both for entities for hypotheses, evidence file items, shoebox items etc. and their associations.

**Figure 10 Investigation Structure in the Ontology**

Associated with the Ontology is the Graphical Query Search Module.   This can be used to construct detailed searches both graphically and using SPARQL for information available from the Ontology.  Figure 11 shows an example of the module and the results.

**Figure 11 SCIIMS Graphical query Module**

The Entity Editor allows a user to navigate an Ontology (a graphical structure is displayed - see Figure 12) and provides direct read and write access.
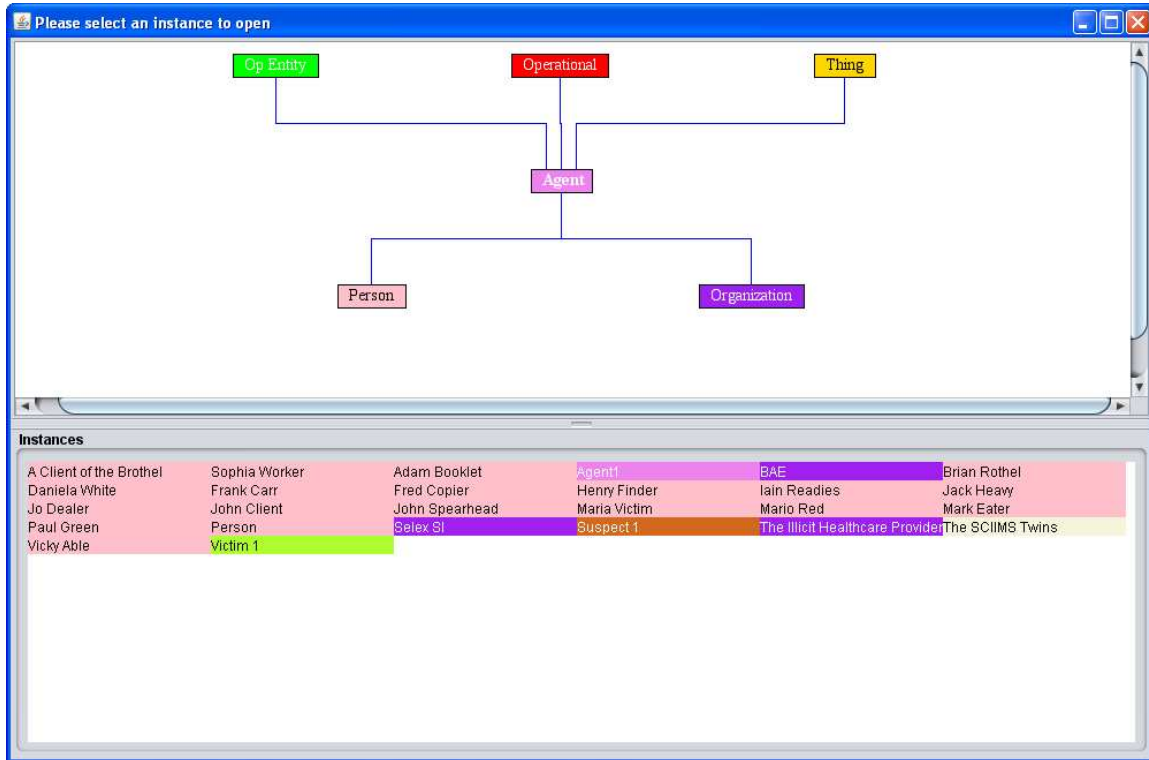


**Figure 12 Entity Editor**

.

## 3.4. HCI

Figure 13 shows an example of the SCIIMS Prototype/Demonstration System HCI, in this case the HCI for the classification of advertisements.
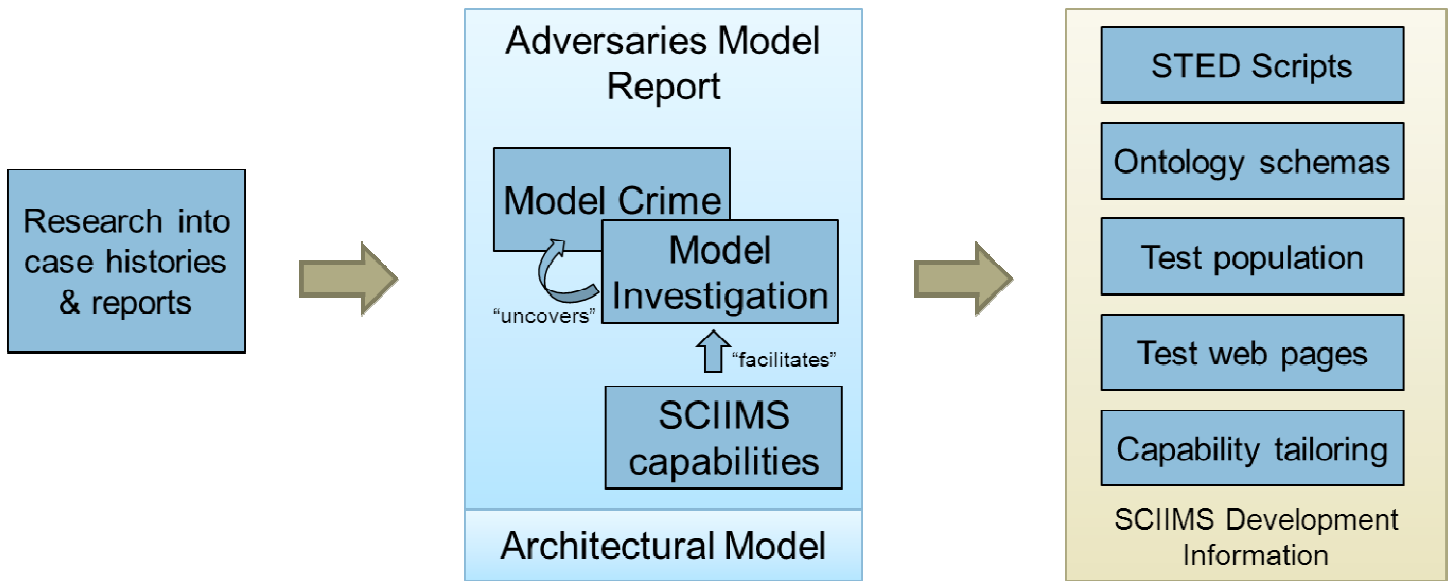
**Figure 13 Example Prototype/Demonstration System HCI Screen**

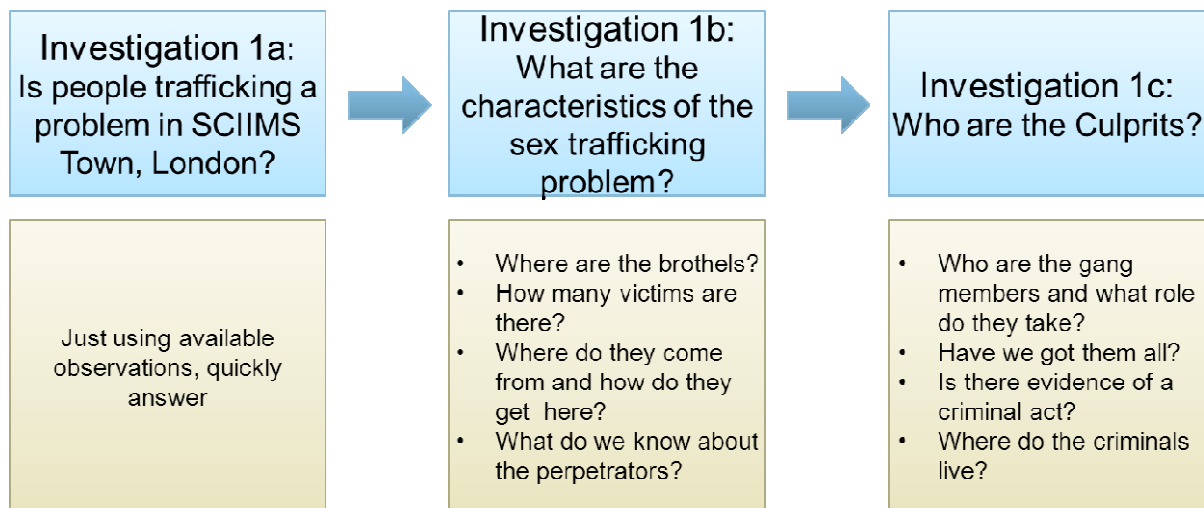## 4. System Test, Demonstrations and Experiments (STED)

A number of demonstrations and system level experiments (where "human in the loop" experiments are carried out).using the Prototype/Demonstration System have been conducted to gain feedback on the various SCIIMS capabilities The overall concept was to produce a test script based on the Adversaries Model Scenario (this describes in detail an example of an investigation into people trafficking ) and use it as a basis for the System Tests, Demonstrations (which uses a subset of the tests) and experimentation   Figure 14 illustrates how the Adversaries Model is used as a basis for the system test, demonstration and experimentation scripts as well as defining the test data for the Prototype/Demonstration System relational databases, ontology and web pages.

- System test, experimentation and demonstration coordinated through architecture
- Basis for future rapid development process

**Figure 14 Innovative Process for Developing SCIIMS**

Figure 15 summarises the three stages to the Adversaries Model investigation. The separate test scripts were written to cover phases 1a, 1b and 1c which separate the scripts covering the foraging parts (e.g. web crawling and data mining) and the sense making scripts.



- *This sequence occurs over several weeks in August 2016 as the criminal details emerge*

**Figure 15 Adversaries Model Investigation Scenario**

## 4.1. Demonstration Results

Figure 16 shows an example of the results from the demonstrations. Anything above 0 is a useful capability.
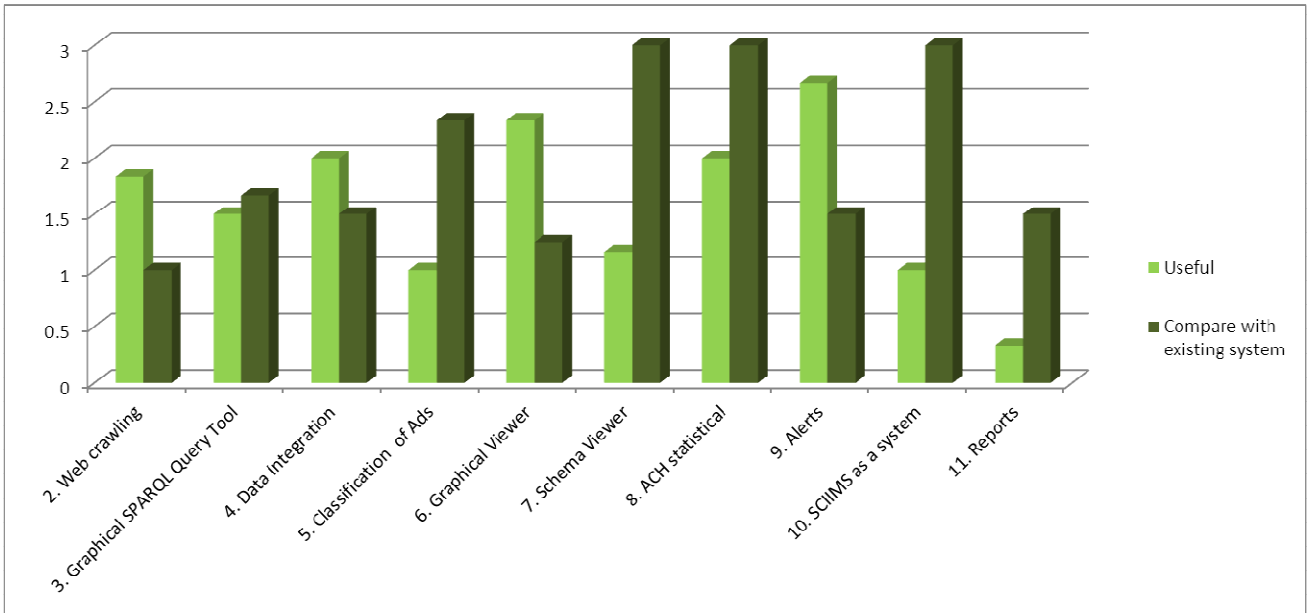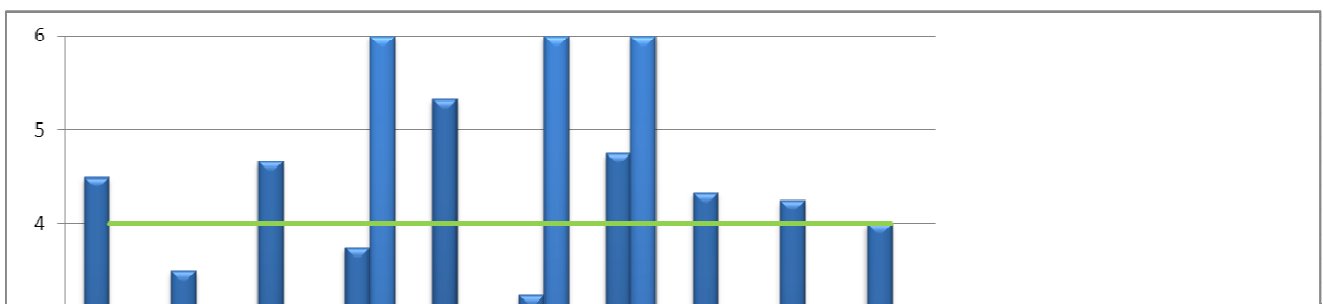
**Figure 16 Usefulness and comparison with existing system**

## 4.2. Experimentation Results

Figure 17 shows an example of the results from the experiments. Anything at or above the green line is easy to use.

**Figure 17 Ease of Use and Comparison with Existing Systems**