

## Statement of Requirement (SoR)

---

Reference Number	[REDACTED]
Version Number	1.0
Date	[REDACTED]

1.	Requirement
1.1	Title
	Automated Ingest and Processing of Cyber Data
1.2	Summary
	<p>Dstl is looking to develop a set of automated ingest pipelines to enable the Cyber Data Framework to mature and move towards an active capability.</p> <p>Phase 1 and Phase 2 of this work have seen the development of a standard Cyber Data Framework (CDF). The CDF is a [REDACTED] owned framework to aid in storing, sharing and managing [REDACTED] cyber event data, and provides a platform for research and development of novel analysis tools. The CDF is heavily influenced by the [REDACTED] data structure and also includes a set of specialist additions recommended by previous work.</p> <p>Most of the work up to this point has been primarily focussed on defining and designing the Cyber Data Framework and its hosting solution. The ingestion pipelines of raw data sources into the CDF data structure were not a focus of Phase 1 or 2 - as such many assumptions were made about these processes and their design was left undefined.</p> <p>In this work, attention should be given to how data could be ingested into the CDF in a way that is <u>as automated as possible</u>.</p> <p><b>Research Questions and Objectives</b></p> <p>The following questions form the basis of this work:</p>

Question 1: Given a set of **known and well-structured cyber data** formats, how can an ingestion pipeline be constructed to **automate** the re-structuring of these raw data file types into the Cyber Data Format?

Question 2: Given a **new (unseen) cyber data format** and the CDF data specification, how could an ingestion pipeline be constructed to **automate** the re-structuring of the raw data into CDF data, as much as possible?

The two questions have subtle differences, the first is focussed on engineering software to enable standard data types to be ingested into the CDF. The second is more research focussed and is an exploration into what new, machine-aided techniques exist that could help with this problem.

The primary objectives of this SOR are:

- To use the proposed CDF Data Specification v2.0 and CDF Implementation Specification from Phase 2 to develop a set of ingest pipelines which are able to transform [REDACTED] into CDF structured data.
- To explore the art of the possible for combining [REDACTED] and other machine-aided techniques with Cyber Data Processing, for an automated approach to ingesting raw data into the CDF.

### Task Structure

The two research questions naturally organise the work into two work packages. It is recommended that the work packages are run sequentially as it is expected insights from WP1 will help in WP2, though this is not a required and may be done in parallel:

#### 1) WP1 – Ingest Pipeline Engineering

The first work package aims to answer Question 1 of the research questions. Many of the data sources the CDF expects to be presented with originate from a set of popular open-sourced cyber security tools, which write their output in well-documented, tabular data formats. In order to more rapidly advance the CDF, there must be an automated way to ingest these data sources, reducing the load on the analyst as much as possible. This work package should be focussed on developing software which is able to automatically extract CDF objects from the raw data and input into the CDF Storage.

#### 2) WP2 – Ingest Pipeline Art of the Possible

	<p>The second work package aims to answer Question 2. The landscape and origins of data sources for the CDF will evolve over time, and the CDF aims to be resilient against changes to the Cyber Data domain. When presented with a new data structure, ideally the CDF will be able to automatically ingest as much of this data source as possible. With recent advances in [REDACTED] and Machine Learning in the Cyber security domain, this work package will include a general scoping of how these techniques could be exploited within the CDF Ingestion process, and where appropriate develop a prototype demonstrating the findings.</p>
1.3	<b>Background</b>

This SOR builds upon the following work previously delivered in Phase 1 and 2:

#### Phase 1

Phase 1 of the CDF task was focussed on exploring the feasibility of a standard Cyber Data Framework (CDF) which enables a range of cyber data sources to be condensed into a single data structure which can be used by analysts in Defensive Cyber operations. The output of Phase 1 was an evaluation of existing open source and commercial Cyber Data structures, followed by a proposed Cyber Data Framework v1.0 based on the [REDACTED] object format, with an additional set of custom [REDACTED] objects.

#### Phase 2

Phase 2 builds upon the proposed CDF v1.0 delivered in Phase 1, and considers how the CDF may be implemented in practice. The scope of this phase included an evaluation of the suitability of CDF v1.0, identification of implementation considerations and proposed a hosting system design.

The four work packages comprised:

- WP1 – CDF for Analysis: To test, understand and describe how the CDF can be used to store data to support analysis in support of cyber defence goals.
- WP2 – CDF for Knowledge Transfer: To identify and overcome the practical knowledge management challenges and describe a pathway for the CDF's integration with wider [REDACTED] organisational structures and assets.
- WP3 – CDF for Exploration: To investigate and define how data interrogation in support of cyber defence analysis can be managed effectively using the CDF.
- WP4 – CDF for Implementation: To provide clear guidance on the implementation of the CDF within the Dstl Cyber Research Team and a roadmap for extending this across the desired user community.

The work was broken down into 4 sequential sprints:

- 1) Framework Evaluation: Evaluated the appropriateness of CDF v1.0, and suggested suitable extensions - resulting in the CDF Data Specification v2.0
- 2) CDF Use Case definition: Designed a set of use cases and worked examples to illustrate how the CDF Data Specification v2.0 is used in practice.
- 3) System Analysis: Examined the challenges and system requirements of using the CDF Data Specification v2.0 in practice for the proposed Use Cases.

	4) <u>Solution Design</u> : Designed a system architecture which meets the requirements and tackles the challenges identified in the System Analysis
1.4	<b>Requirement</b>
	<p><b>Definition 1.1 (Ingestion Pipeline)</b></p> <p>A CDF <b>Ingestion Pipeline</b> is a set of processes which takes a set of raw cyber data files and re-structures them into the format described in the CDF Data Specification from Phase 2. The Ingestion pipeline <u>contains at least</u> the following processes as proposed in Phase 2:</p> <ul style="list-style-type: none"> <li>• Raw Data Format Identification – Is this a known log type, or unseen?</li> <li>• Data Cleansing – Are there any malformed entries, or missing information?</li> <li>• Data Reformatting to CDF Data Structure</li> <li>• Raw File metadata capture (including Data Classification/Caveats and other protective markings)</li> <li>• Putting Data into CDF Storage with relevant CDF Auditing</li> </ul> <p><b>WP1 – Ingest Pipeline Engineering</b></p> <p>This work package should be primarily focussed on answering Research Question 1:</p> <p><i>‘Given a set of known and well-structured cyber data formats, how can an ingestion pipeline be constructed to automate the re-structuring of these raw data file types into the Cyber Data Format?’</i></p> <p>This should be done by creating software which enables ingest of ‘common’ cyber data into the CDF. ‘Common’ cyber data in this context covers any data format expected by Cyber Security Analysts within [REDACTED] and currently includes (but not limited to) [REDACTED] alerts.</p> <p>1.1 Must use [REDACTED] alerts as the known data structures, these data sources are relevant to [REDACTED] and the CDF’s intended users.</p> <p>1.2 Must propose a suitable software pipeline, which is compatible with the proposed CDF Data Specification v2.0 and integrate with the proposed CDF implementation architecture.</p> <p>1.3 Where possible the proposed pipeline should be resilient to software changes with the CDF data solution, changes with the CD Data specification and resilient to expected future upscaling of the CDF.</p>

- 1.4 The software pipeline must use components licenced suitably for [REDACTED] to distribute with no further cost or imposed restrictions. Specific permitted licences to be agreed between Dstl and the supplier.
- 1.5 Must develop a set of ingest pipelines for each of the common data formats identified in 1.1, which include each of the steps highlighted **Definition 1.1**. Where any of these steps can't be automated, a reasonable manual process should be identified instead.
- 1.6 Must capture any design choices and limitations faced during this process and document in the accompanying report.
- 1.7 Must develop, alongside any software produced, suitably detailed documentation and user guides.
- 1.8 The ingest pipelines must exceed the processing rates expected through manual input. Moreover, performance improvements that take minimal effort to implement should be made
- 1.9 Similarly, whilst no specific error rates are required, the number of errors made by the pipelines must not exceed the expected number of errors made during manual processing.
- 1.10 The pipeline must be able to deliver acceptable performance for injecting one hours' worth of logging (roughly [REDACTED] with an average size of [REDACTED] and a maximum size of [REDACTED] in an hour')
- 1.11 Must define requirements for deploying such applications, and how the CDF data solution should be adapted to account for these additional processes.

## WP2 – Ingest Pipeline Art of the Possible

This work package is focussed on answering Research Question 2 - *“Given a new (unseen) cyber data format and the CDF data specification, how could an ingestion pipeline be constructed to automate the re-structuring of the raw data into CDF data, as much as possible?”*

This work package should explore state of the art Machine Learning techniques that could aid in any of the steps of the ingest pipelines during the processing of unseen cyber data. This could be exploiting [REDACTED] on unseen text-based logs or similar. The required level of automation is not strict, this activity should explore which parts of the ingestion pipeline can be automated. If Machine Learning is not a viable option, then other software solutions should be

	<p>recommended which reduce the load on the operator during the processing of new data structures.</p> <p>2.1 Must perform a literature review and scoping for relevant [REDACTED] techniques or other automation approaches that may be appropriate for this work.</p> <p>2.2 For each stage in the pipeline, where appropriate, identify possible methods of automating the process. This should take into consideration any issues that may arise when using these methods (trust, ethics, compute resources, level of human input required).</p> <p>2.3 Must develop a software prototype which exploits the techniques identified on a new, unseen cyber data source.</p> <p>2.4 In software where Machine Learning is used then the [REDACTED] framework is preferred and full access to the model and training set is required.</p> <p>2.5 Must define requirements for deploying such applications, and how the CDF data solution should be adapted to account for these additional processes.</p> <p><b><u>Security Requirements</u></b></p> <p>The outputs of the project shall be classified no higher than [REDACTED].</p> <p><b><u>Timescales</u></b></p> <p>All work <u>must</u> be delivered in [REDACTED]</p>
<b>1.5</b>	<b>Options or follow on work</b>
	Not applicable

1.6	Deliverables & Intellectual Property Rights (IPR)						
Ref.	Title	Due by	Format	TRL*	Expected classification (subject to change)	What information is required in the deliverable	IPR DEFCON/ Condition
D1	CDF Ingestion Pipelines & Documentation	REDACTED]	Software and Documentation	3	[REDACTED]	Software and associated documentation and instructions for the outputs of WP1: <ul style="list-style-type: none"><li>• User Guides</li><li>• Installation Process</li></ul> Where appropriate, libraries, source code, and other artefacts needed to build and distribute the pipeline will be delivered	DEFCON 703
D2	CDF Ingestion Pipeline Report	[REDACTED]	Report		[REDACTED]	Report describing the design process of the pipelines for WP1, including: <ul style="list-style-type: none"><li>• Platform evaluation process</li><li>• Pipeline design considerations and limitations</li><li>• Results from testing</li></ul>	DEFCON 703



						<ul style="list-style-type: none"> <li>Resilience against future software changes/upscaling</li> </ul>	
D3	Automated Cyber Data Processing Report	[REDACTED]	Report		[REDACTED]	<p>Report describing the research performed in WP2, including:</p> <ul style="list-style-type: none"> <li>Literature Review</li> <li>Art of the possible mapping to ingest pipeline stages</li> <li>design and performance analysis</li> <li>Wider CDF solution considerations</li> </ul>	DEFCON 703
D4	CDF Automated Ingest Pipelines & Documentation	[REDACTED]	Software and Documentation		[REDACTED]	<ul style="list-style-type: none"> <li>Software (source code) and associated documentation for the outputs of WP2, (including documentation, and ML models and training set)</li> </ul> <p>Where appropriate, libraries, source code, and other artefacts needed to build and distribute the pipeline will be delivered</p>	DEFCON 703

**\*Technology Readiness Level required**

1.7	<b>Standard Deliverable Acceptance Criteria</b>
	<p><i>As per Framework T&amp;Cs</i></p> <p>           All Reports included as Deliverables under the Contract e.g. Progress and/or Final Reports etc. must comply with the <a href="#">Defence Research Reports Specification (DRRS)</a> which defines the requirements for the presentation, format and production of scientific and technical reports prepared for MoD.         </p> <p>           Interim or Progress Reports: The report should detail, document, and summarise the results of work done during the period covered and shall be in sufficient detail to comprehensively explain the results achieved; substantive performance; a description of current substantive performance and any problems encountered and/or which may exist along with proposed corrective action. An explanation of any difference between planned progress and actual progress, why the differences have occurred, and if behind planned progress what corrective steps are planned.         </p> <p>           Final Reports: shall describe the entire work performed under the Contract in sufficient detail to explain comprehensively the work undertaken and results achieved including all relevant technical details of any hardware, software, process or system developed there under. The technical detail shall be sufficient to permit independent reproduction of any such process or system.         </p> <p>           All Reports shall be free from spelling and grammatical errors and shall be set out in accordance with the Statement of Requirement (1) above.         </p> <p>           Failure to comply with the above may result in the Authority rejecting the deliverables and requesting re-work before final acceptance.         </p>
1.8	<b>Specific Deliverable Acceptance Criteria</b>
	<p>           All reports shall be delivered as both a Microsoft Word (.docx) and Portable Document Format (PDF) (.pdf) document unless specified below (or agreed with the Authority prior to delivery).         </p>

2.	<b>Quality Control and Assurance</b>
2.1	<b>Quality Control and Quality Assurance processes and standards that must be met by the contractor</b>
	<input checked="" type="checkbox"/> <b>ISO9001</b> (Quality Management Systems)  <input type="checkbox"/> <b>ISO14001</b> (Environment Management Systems)  <input type="checkbox"/> <b>ISO12207</b> (Systems and software engineering — software life cycle)  <input type="checkbox"/> <b>TickITPlus</b> (Integrated approach to software and IT development)  <input type="checkbox"/> <b>Other:</b> (Please specify below)
2.2	<b>Safety, Environmental, Social, Ethical, Regulatory or Legislative aspects of the requirement</b>
	Not Applicable

<b>3.</b>	<b>Security</b>	
<b>3.1</b>	<b>Highest security classification</b>	
	<b>Of the work</b>	[REDACTED]
	<b>Of the Deliverables/ Output</b>	[REDACTED]
<b>3.2</b>	<b>Security Aspects Letter (SAL)</b>	
	Yes If yes, please see SAL reference- [REDACTED]	
<b>3.3</b>	<b>Cyber Risk Level</b>	
	Choose an item.[REDACTED]	
<b>3.4</b>	<b>Cyber Risk Assessment (RA) Reference</b>	
	[REDACTED]  If stated, this must be completed by the contractor before a contract can be awarded. In accordance with the <a href="#">Supplier Cyber Protection Risk Assessment (RA) Workflow</a> please complete the Cyber Risk Assessment available at <a href="https://suppliercyberprotection.service.xgov.uk/">https://suppliercyberprotection.service.xgov.uk/</a>	

<b>4.</b>	<b>Government Furnished Assets (GFA)</b>				
GFA to be Issued - Yes					
<b>GFA No.</b>	<b>Unique Identifier/ Serial No</b>	<b>Description:</b>	<b>Available Date</b>	<b>Issued by</b>	<b>Return Date or Disposal Date (T0+)</b>

GFA-1		GFI – the final deliverables from Phases 1 and 2 of the CDF work			[REDACTED]
		Datasets from [REDACTED] and other sources may be required			[REDACTED]

<b>5.</b>	<b>Proposal Evaluation criteria</b>										
<b>5.1</b>	<b>Technical Evaluation Criteria</b>										
	<p>The standard EW&amp;C scoring criteria will be used :</p> <table> <tr> <td><b>Technical – to include individual team members experience and expertise</b></td><td>45% Weighting</td></tr> <tr> <td><b>Value for Money</b></td><td>15% Weighting</td></tr> <tr> <td><b>Delivery – to include plan and resource management</b></td><td>15% Weighting</td></tr> <tr> <td><b>Quality</b></td><td>20% Weighting</td></tr> <tr> <td><b>Commercial</b></td><td>5% Weighting</td></tr> </table>	<b>Technical – to include individual team members experience and expertise</b>	45% Weighting	<b>Value for Money</b>	15% Weighting	<b>Delivery – to include plan and resource management</b>	15% Weighting	<b>Quality</b>	20% Weighting	<b>Commercial</b>	5% Weighting
<b>Technical – to include individual team members experience and expertise</b>	45% Weighting										
<b>Value for Money</b>	15% Weighting										
<b>Delivery – to include plan and resource management</b>	15% Weighting										
<b>Quality</b>	20% Weighting										
<b>Commercial</b>	5% Weighting										
<b>5.2</b>	<b>Commercial Evaluation Criteria</b>										
	As per framework evaluation criteria										