

Identifying Malicious Threats to Scientific Data Integrity Using MITRE ATT&CK[®]

Emily K. Adams

(ekadams@iu.edu)

Center for Applied Cybersecurity Research, Indiana University

16 Aug, 2022 (version 1.0)

Guaranteeing the integrity of scientific workflow processing on distributed cyberinfrastructure is of paramount importance for the validity and credibility of scientific research. Data-driven research depends on the integrity of underlying scientific computational workflow and on the integrity of associated data products. Data integrity, critical in producing sound research results, refers to the properties of information that have not been altered or modified by an unauthorized person or process.

The SWIP project¹ addressed integrity checking to ensure scientific workflow computations are free from integrity errors. Informed by the SWIP project, IRIS² was initially developed to detect, diagnose, and pinpoint the root cause of unintentional (*i.e.* non-malicious) integrity errors in the scientific workflow executions on distributed cyberinfrastructure. This paper extends the IRIS *non-malicious* threat model³ to uncover *malicious* attack tactics and techniques that impact scientific data integrity within the scientific workflows by leveraging the MITRE ATT&CK[®] framework⁴ to conduct data impact analysis.

Methodology

Approach

Arguably, a number of malicious activities can threaten research activities and can have a negative influence on the availability or integrity of research processes, assets, and personnel. In this paper *malicious* tactics and techniques leveraged to *explicitly* manipulate the integrity of transient workflow data, data products, or derived metadata within scientific workflows are considered. This document leverages the MITRE ATT&CK Enterprise knowledge base of adversary tactics and techniques, based on real-world observations, as the foundation for a scoped analysis enumerating malicious attacks against data integrity within scientific workflows.

¹ <https://cacr.iu.edu/projects/swip/>

² https://www.nsf.gov/awardsearch/showAward?AWD_ID=1839900

³ <https://hdl.handle.net/2022/27980> I. Abhinit and V. Welch, “Data Integrity Threat Model using Open Science Cyber Risk Profile”

⁴ <https://attack.mitre.org/>

Open Science Cyber Risk Profile (OSCRP)

To enumerate the areas of a research workflow impacted by malicious attacks against the integrity of the data, this analysis leverages the Open Science Cyber Risk Profile (OSCRP) asset classification system used in the non-malicious threat modeling.⁵ Only three OSCRP asset categories will be used in this analysis as they are explicitly scoped to data products contained within the scientific workflow.

Data Products:⁶ Data that is taken as input by a workflow or generated by a workflow and expected to persist and be consumed by Data Consumers.

Transient Workflow Data:⁷ Transient data products created during the course of the workflow execution and that do not persist meaningfully past the completion of the workflow. E.g. intermediate data products generated by one job in the workflow and consumed by the subsequent job and then discarded.

Metadata:⁸ Data created by a workflow execution describing the workflow execution and resulting data products. This includes integrity information (e.g. hashes) about the Data Products added by SWIP project enhancements. Metadata is expected to persist in perpetuity.

For the remainder of this paper, these three data-centric OSCRP asset categories are to be collectively referenced as “scientific data”.

Classification of Attacks on Research Data Integrity

An exploited asset or a compromised process within the research workflow ecosystem may position an adversary to ultimately impact research integrity from a number of vectors and avenues of attack. However, a core assumption of this analysis is that the identified malicious techniques and tactical objectives are scoped to *explicitly* interfere with the integrity of scientific data within the research workflow.

The following attack definitions are used to assess the relevancy of ATT&CK tactics and techniques leveraged to impact scientific data integrity within scientific workflows. Only *Direct Attacks* and *Indirect Attacks* are included in this analysis; *General Attacks* are not considered in this analysis.

- **Direct Attack:** Malicious activity directly targeting scientific data assets with the intent of impacting data integrity. For example:
 - Corruption of a database of measurements gathered during research
 - Deletion of a text file containing metadata derived from

⁵ See, [APPENDIX A: OSCRP-Derived Asset Classification for Scientific Workflows](#)

⁶ We treat “Data Products” as OSCRP Public Data, <http://trustedci.github.io/OSCRP/assets/Public-Data/>.

⁷ We treat “Transient Workflow Data” as OSCRP Internal Data, <http://trustedci.github.io/OSCRP/assets/InternalData/>

⁸ We treat “Metadata” as OSCRP Accounting Information, <http://trustedci.github.io/OSCRP/assets/AccountingInformation/>.

- **Indirect Attack:** Malicious activity targeting an asset or process that interfaces with or produces scientific data with the intent of impacting scientific data integrity. For example:
 - Altering software used to generate or process scientific data
 - Disabling devices used to collect raw measurement data
- **General Attack:** Malicious activity that impacts the security of assets, functions, and personnel supporting the research activities. For example:
 - Exfiltrating data over a compromised network infrastructure
 - Conduct a denial of service campaign against a computational system cluster

MITRE ATT&CK Tactics, Techniques, & Procedures

The selection of the ATT&CK Enterprise class Tactics⁹, Techniques¹⁰ included in this analysis focuses on elements of a malicious attack that may *directly* or *indirectly* impact the integrity of scientific data within scientific workflows. Tactics and techniques that facilitate or include malicious activity *not* explicitly related to data integrity (*i.e.* reconnaissance of research IT infrastructure, theft of account identifiers and passwords, lateral movement across networks) are not included in this analysis. As a number of malicious activities can ultimately situate a malicious actor to impact data integrity, ATT&CK Procedures (*i.e.* the “how” of implementing an ATT&CK Technique) are referenced to determine if an attack is designated *Direct Attack*, *Indirect Attack*, or *General Attack*.

ATT&CK Tactics

The ATT&CK Tactic classifications describe the “why” of an ATT&CK Techniques. Of the fourteen ATT&CK Tactics, the following two are leveraged in this analysis as they represent categories of malicious activity that can *directly* or *indirectly* impact the integrity of scientific data within scientific workflows.

Impact [TA0040]¹¹ - The adversary is trying to manipulate, interrupt, or destroy systems and data.

Execution [TA0002]¹² - The adversary is trying to run malicious code.

ATT&CK Techniques

MITRE ATT&CK Techniques represent what actions a malicious actor may perform or execute to succeed in a tactical goal. This includes malicious activity that can *directly* or *indirectly* impact the integrity of scientific data within the scientific workflows. Many ATT&CK Techniques include Sub-Techniques, which enumerate particular ways to carry out the action outlined in a Technique. For the purpose of this paper, ATT&CK Sub-Techniques are incorporated into the parent Technique unless otherwise indicated in Appendix B: Excluded ATT&CK Tactics, Techniques, & Sub-Techniques.

⁹ <https://attack.mitre.org/tactics/enterprise/>

¹⁰ <https://attack.mitre.org/techniques/enterprise/>

¹¹ <https://attack.mitre.org/tactics/TA0040>

¹² <https://attack.mitre.org/tactics/TA0002>

Analysis of Malicious Threats to Research Data Integrity

The tables contained within the two Tactic subsections below (*i.e.* Impact [TA0040], Execution [TA0002]) leverages the ATT&CK[®] adversary tactics and techniques, OSCRIP asset classification, and attack type to derive the data integrity concern. Columns 4 and 5 describe the scientific workflow impact to data, data products, or derived metadata.

Column 1: **ATT&CK Technique**

ATT&CK Technique an adversary performs to impact data integrity

Values: ATT&CK Technique ID number and title

Column 2: **OSCRIP Data Asset Class**

Class of scientific data products impacted by the attack

Values: “Transient Workflow Data”; “Data Products”; and/or “Metadata”

Column 3: **Attack Type**

Type of attack against scientific data or against asset(s) interfacing with the data

Values: “Direct” or “Indirect”

Column 4: **Data Integrity Concern**

Action(s) adversaries can perform on or with the data product following a successfully implemented attack

Column 5: **Scientific Workflow Impact**

The negative consequence(s) of a realized concern to the integrity of the scientific data

ATT&CK Tactic: Impact [TA0040]

Threat: The adversary is trying to manipulate, interrupt, or destroy your systems and data.

The Impact Tactic (TA0040) includes thirteen techniques that adversaries use to compromise integrity by manipulating business and operational processes. Per the ATT&CK framework, “Techniques used for impact can include destroying or tampering with data. In some cases, business processes can look fine, but may have been altered to benefit the adversaries’ goals.” Seven ATT&CK Techniques within the Impact Tactic are identified in Table 1 to *directly* or *indirectly* impact scientific workflow data integrity.

Table 1. ATT&CK Tactic Impact [TA0040]

ATT&CK Technique	OSCRP Data Asset Class	Attack Type	Data Integrity Concern (An adversary may...)	Scientific Workflow Impact
T1485 Data Destruction	Data Products Metadata	Direct	<ul style="list-style-type: none"> • Destroy data and files • Render stored data irrecoverable by forensic techniques 	<ul style="list-style-type: none"> • Scientific workflow cannot be initiated or executed • Data lost
T1486 Data Encrypted for Impact	Data Products Metadata	Direct	<ul style="list-style-type: none"> • Render stored data inaccessible by encrypting files or data 	<ul style="list-style-type: none"> • Data rendered unusable • Data lost
T1489 Service Stop	Data Products Transient Workflow Metadata	Direct & Indirect	<ul style="list-style-type: none"> • Stop or disable services on a system to render those services unavailable • Stop services or processes in order to destroy or manipulate data 	<ul style="list-style-type: none"> • Data rendered unusable • Data lost • Scientific Workflow cannot be initiated or executed • Scientific workflow producing incorrect/invalid results
T1495 Firmware Corruption	Data Products Transient Workflow Metadata	Indirect	<ul style="list-style-type: none"> • Overwrite or corrupt device firmware to render a system/device inoperable 	<ul style="list-style-type: none"> • Scientific Workflow cannot be initiated or executed • Loss/damage to research systems managing data
T1499 Endpoint Denial of Service ¹³	Data Products Transient Workflow Metadata	Indirect	<ul style="list-style-type: none"> • Stop or disable a device or component to render the resource unavailable 	<ul style="list-style-type: none"> • Scientific Workflow cannot be initiated or executed • Loss/damage to research systems managing data
T1561 Disk Wipe ¹⁴	Data Products Metadata	Indirect	<ul style="list-style-type: none"> • Wipe or corrupt raw disk data on a system or storage devices 	<ul style="list-style-type: none"> • Scientific workflow cannot be initiated or executed • Loss/damage to research systems managing data
T1565 Data Manipulation ¹⁵	Data Products Metadata	Direct	<ul style="list-style-type: none"> • Insert, delete, or manipulate data in order to influence external outcomes files or data 	<ul style="list-style-type: none"> • Scientific workflow producing incorrect/invalid results

ATT&CK Tactic: Execution [TA0002]

Threat: The adversary is trying to run malicious code.

Execution consists of techniques that result in adversary-controlled code running on a local or remote system. Per the ATT&CK framework, “Techniques that run malicious code are often paired with techniques from all other tactics to achieve broader goals, like exploring a network or stealing data. For example, an adversary might use a remote access tool to run a PowerShell script that does Remote System Discovery.” Seven ATT&CK Techniques within the Execution Tactic are identified here *directly* or *indirectly* impact scientific workflow data integrity.

¹³ Includes ATT&CK Sub-Techniques: *T1499.001* OS Exhaustion Flood, *T1499.002* Service Exhaustion Flood, *T1499.003* Application Exhaustion Flood, *T1499.004* Application or System Exploitation

¹⁴ Includes ATT&CK Sub-Technique: *T1561.001* Disk Content Wipe

¹⁵ Includes ATT&CK Sub-Techniques: *T1565.001* Stored Data Manipulation, *T1565.002* Transmitted Data Manipulation

Table 2. ATT&CK Tactic: Execution [TA0002]

ATT&CK Technique	OSCRP Data Asset Class	Attack Type	Data Integrity Concern (An adversary may...)	Scientific Workflow Impact
T1053 Scheduled Task/Job ¹⁶	Data Products Metadata	Indirect	<ul style="list-style-type: none"> Disrupt task scheduling of data processing activities Abuse task scheduling to potentially mask one-time execution 	<ul style="list-style-type: none"> Scientific workflow cannot be initiated or executed Scientific workflow is executed at unexpected intervals impacting downstream data
T1059 Command and Scripting Interpreter ¹⁷	Data Products Transient Workflow Metadata	Direct & Indirect	<ul style="list-style-type: none"> Abuse command and script interpreters to execute commands, scripts, or binaries to impact data integrity 	<ul style="list-style-type: none"> Scientific workflow cannot be initiated or executed Scientific workflow producing incorrect/invalid results
T1106 Native API	Data Products Metadata	Indirect	<ul style="list-style-type: none"> Calling low-level OS services within the kernel, such as those involving hardware/devices, memory, and processes. 	<ul style="list-style-type: none"> Scientific workflow cannot be initiated or executed
T1129 Shared Modules	Data Products Transient Workflow Metadata	Indirect	<ul style="list-style-type: none"> Subvert shared modules integrated into the scientific workflow 	<ul style="list-style-type: none"> Scientific workflow cannot be initiated or executed Scientific workflow producing incorrect/invalid results Subvert components that interface directly with scientific data
T1203 Exploitation for Client Execution ¹⁸	Data Products Transient Workflow Metadata	Indirect	<ul style="list-style-type: none"> Employ targeted exploitation of software vulnerabilities to take advantage of applications or processes used in the scientific workflow code execution. 	<ul style="list-style-type: none"> Subvert components that interface directly with scientific data
T1559 Inter-Process Communication	Data Products Transient Workflow Metadata	Direct & Indirect	<ul style="list-style-type: none"> Abuse inter-process communication (IPC) mechanisms for local code or command execution 	<ul style="list-style-type: none"> Scientific workflow is degraded or compromised Scientific workflow producing incorrect/invalid results
T1569 System Services ¹⁹	Data Products Metadata	Direct & Indirect	<ul style="list-style-type: none"> Abuse system services for one-time or temporary execution of malicious commands or payloads 	<ul style="list-style-type: none"> Scientific workflow is degraded or compromised Scientific workflow producing incorrect/invalid results

¹⁶ Includes ATT&CK Sub-Techniques: *T1053.002* At, *T1053.003* Cron, *T1053.005* Scheduled Task, *T1053.006* Systemd Timers

¹⁷ Includes ATT&CK Sub-Techniques: *T1059.001* PowerShell, *T1059.002* AppleScript, *T1059.003* Windows Command Shell, *T1059.004* Unix Shell, *T1059.005* Visual Basic, *T1059.006* Python, *T1059.007* JavaScript

¹⁸ Includes ATT&CK Sub-Techniques: *T1559.001* Component Object Model, *T1559.002* Dynamic Data Exchange, *T1559.003* XPC Services

¹⁹ Includes ATT&CK Sub-Techniques: *T1569.001* Launchctl, *T1569.002* Service Execution

Application & Future Work

By systematically identifying the adversarial tactics and techniques that underpin malicious threats to scientific data within the scientific workflow, researchers and data stewards will be equipped to proactively design scientific workflows to ensure data integrity, better understand how malicious attacks against data integrity are executed, and use this knowledge to build detections and defenses to protect scientific data.

The approach of the two threat models and resulting analysis developed through IRIS can be applied for not only detecting data integrity errors during the execution of scientific workflows, but also for pinpointing the root cause of those errors. Thus, results from the threat analysis can inform and enrich the Machine Learning models for data integrity root cause analysis to more fully align to assessing threats observed in real systems. While one can inject errors based on anecdotal evidence, as in previous works, injecting errors guided by findings from malicious and non-malicious threat models should be a far superior technique that is grounded in well-established threat modeling approaches.

APPENDIX A: OSCRP-Derived Asset Classification for Scientific Workflows

1. **Transient Workflow Data:** Transient data products created during the course of the workflow execution and that do not persist meaningfully past the completion of the workflow. E.g. intermediate data products generated by one job in the workflow and consumed by the subsequent job and then discarded.
 - a. We treat Transient Workflow Data as [OSCRP Internal Data](#).
2. **Data Products:** Data that is taken as input by a workflow or generated by a workflow and expected to persist and be consumed by Data Consumers.
 - a. We treat Data Products as [OSCRP Public Data](#) (To consider confidentiality, one could treat this as OSCRP Embargoed Data).
3. **Metadata:** Data created by a workflow execution describing the workflow execution and resulting data products. This includes integrity information (e.g. hashes) about the Data Products added by SWIP project enhancements. Metadata is expected to persist in perpetuity.
 - a. We treat Metadata as [OSCRP Accounting Information](#).
4. **Researcher System:** The interface used by the Researcher to craft and initiate the workflow. Specifically, it executes Workflow Managed System client, holds abstract workflow description and workflow plan, holds metadata (including integrity information) regarding workflows, and holds Researcher credentials for accessing Computational and Data Storage Systems.
 - a. We treat the Researcher system as an [OSCRP Desktop](#).
5. **Workflow Management System (WMS):** This system translates abstract workflow description from Researcher into work plan, maps workflow plan to Computational and Data Storage Systems, manages workflow execution, and orchestrates the creation and checking of data integrity information.
 - a. We treat the WMS as an [OSCRP Workflow](#).
6. **Computational Systems:** Computer systems that run computational aspects of the workflow, provide for temporary data storage during workflow (during computation and stage-in/out), provide the software stacks for workflow execution, and create and check data integrity information within the workflow.
 - a. We treat Computational Systems as [OSCRP Servers](#).
7. **Data Storage Systems:** Computer systems that provide for long-term storage of data consumed by and produced by workflows, and serve to make data available to Data Consumers.
 - a. We treat Data Storage Systems as [OSCRP File Stores](#).
8. **Network System:** The IT system that transports data between the Research Systems, Workflow Management System, Computational Systems, and Data Storage Systems involved in executing the workflow.
 - a. We treat the Network System as an [OSCRP Network](#).

APPENDIX B: Excluded ATT&CK Tactics, Techniques, & Sub-Techniques

Based on descriptions put forth by the MITRE ATT&CK framework, ATT&CK Tactics determined *not* to have a direct or indirect impact on data integrity are:

Initial Access [TA0001]	Discovery [TA0007]
Defense Evasion [TA0005]	Lateral Movement [TA0008]
Persistence [TA0004]	Collection [TA0009]
Reconnaissance [TA0043]	Exfiltration [TA0010]
Resource Development [TA0042]	Command and Control [TA0011]
Credential Access [TA0006]	

Of the ATT&CK Tactics assessed in this paper, the following lists itemize ATT&CK Techniques and Sub-Techniques *not* relevant to a *direct* or *indirect* attack on scientific data integrity.

ATT&CK Tactic: Impact [TA0040]

T1489	Service Stop
T1490	Inhibit System Recovery
T1491	Defacement
T1496	Resource Hijacking
T1498.001	Network Denial of Service, Direct Network Flood
T1498.002	Network Denial of Service, Reflection Amplification
T1529	System Shutdown/Reboot
T1531	Account Access Removal
T1561.002	Disk Wipe, Disk Structure Wipe

ATT&CK Tactic: Execution [TA0002]

T1047	Windows Management Instrumentation
T1053.004	Scheduled Task/Job, Launchd (deprecated)
T1053.007	Scheduled Task/Job, Container Orchestration Job
T1059.008	Network Device CLI
T1072	Software Deployment Tools
T1072	Software Deployment Tools
T1204.001	User Execution, Malicious Link
T1204.002	User Execution, Malicious File
T1204.003	User Execution, Malicious Image
T1609	Container Administration Command