

# Application of BagIt-Serialized Research Object Bundles for Packaging and Re-execution of Computational Analyses

Kyle Chard  
Computation Institute  
University of Chicago  
Chicago, IL  
chard@uchicago.edu

Bertram Ludäscher  
School of Information Sciences  
University of Illinois at Urbana-Champaign  
Champaign, IL  
ludaesch@illinois.edu

Ian Taylor  
Center for Research Computing  
University of Notre Dame  
South Bend, IN  
ian.j.taylor@gmail.com

Niall Gaffney  
Texas Advanced Computing Center  
University of Texas at Austin  
Austin, TX  
ngaffney@tacc.utexas.edu

Timothy McPhillips  
School of Information Sciences  
University of Illinois at Urbana-Champaign  
Champaign, IL  
tmcphill@illinois.edu

Thomas Thelen  
NCEAS  
University of California at Santa Barbara  
Santa Barbara, CA  
thelen@nceas.ucsb.edu

Matthew B. Jones  
NCEAS  
University of California at Santa Barbara  
Santa Barbara, CA  
jones@nceas.ucsb.edu

Jarek Nabrzyski  
Center for Research Computing  
University of Notre Dame  
South Bend, IN  
jaroslaw.nabrzyski.1@nd.edu

Matthew J. Turk  
School of Information Sciences  
University of Illinois at Urbana-Champaign  
Champaign, IL  
mjturk@illinois.edu

Kacper Kowalik  
NCSA  
University of Illinois at Urbana-Champaign  
Urbana, IL  
kowalikk@illinois.edu

Victoria Stodden  
School of Information Sciences  
University of Illinois at Urbana-Champaign  
Champaign, IL  
vcs@stodden.net

Craig Willis<sup>†</sup>  
School of Information Sciences  
University of Illinois at Urbana-Champaign  
Champaign, IL  
willis8@illinois.edu

<sup>†</sup>Corresponding author

**Abstract**—In this paper we describe our experience adopting the Research Object Bundle (RO-Bundle) format with BagIt serialization (BagIt-RO) for the design and implementation of “tales” in the Whole Tale platform. A *tale* is an executable research object intended for the dissemination of computational scientific findings that captures information needed to facilitate understanding, transparency, and re-execution for review and computational reproducibility at the time of publication. We describe the Whole Tale platform and requirements that led to our adoption of BagIt-RO, specifics of our implementation, and discuss migrating to the emerging Research Object Crate (RO-Crate) standard.

**Index Terms**—Reproducibility of Results, Standards, Packaging, Interoperability, Software, Digital Preservation

## I. INTRODUCTION

Whole Tale (<http://wholetale.org>) is a web-based, open-source platform for reproducible research supporting the creation, sharing, execution, and verification of “tales” [2], [5]. Tales are executable research objects that capture the code, data, and environment along with narrative and workflow information needed to re-create computational results from scientific studies. A goal of the Whole Tale platform is to produce an archival package that is exportable, publishable, and

can be used for verification of computational reproducibility, for example as part of the peer-review process.

Since its inception, the Whole Tale platform has been designed to bring together existing open science infrastructure. Researchers can ingest data from various scientific archival repositories; launch popular analytical tools (such as Jupyter and RStudio); create and customize computational environments (using `repo2docker`<sup>1</sup>); conduct analyses; create/upload code and data; and publish the resulting package back to an archival repository. Tales are also downloadable and re-executable locally, including the ability to retrieve remotely published data.

With the May 2019 release of version 0.7 of the platform we adopted the Research Object Bundle BagIt serialization (BagIt-RO) format [17]. By combining the BagIt-RO serialization with our `repo2docker`-based execution framework and the BDBag tools [4], we were able to define and implement a standards-compliant, self-describing, portable, re-executable research object with the ability to retrieve remotely published data.

<sup>1</sup><https://repo2docker.readthedocs.io/>

In this paper we describe the Whole Tale platform and requirements that led to our adoption of the BagIt-RO format. The paper is organized as follows. In Section II, we present a motivating example of the use of the Whole Tale platform followed by a brief description of the system architecture in Section III. In Section IV we outline the requirements that led to our adoption of the BagIt-RO format. In Section V we describe our implementation in more detail followed by a discussion and conclusions.

## II. EXAMPLE SCENARIO: ANALYZING SEAL MIGRATION PATTERNS

We begin with a motivating example to illustrate the end-to-end Whole Tale workflow for creating, exporting, and publishing a tale based on existing data archived using the Research Workspace<sup>2</sup>, a DataONE member node. This example is based on tutorial material described in [14].

A research team is preparing to publish a manuscript describing a computational model for estimating animal movement paths from telemetry data. The source data for their analysis, tracking data for juvenile seals in Alaska [3], has been published in Research Workspace, a DataONE network member. Using the Whole Tale platform, the researchers register the external dataset. They then create a new tale by launching an RStudio environment based on images maintained by the Rocker Project [1]. Using the interactive environment, they clone a Github repository, modify an R Markdown document, customize the environment by specifying OS and R packages via `repo2docker` configuration files, and execute their code to generate outputs. They download the package in a compressed BagIt-RO format and run locally to verify their tale. Finally, they enter descriptive metadata and publish the final package back to DataONE to archive the package and obtain a persistent identifier to include in publication.

This scenario is further illustrated in Figure 1.

## III. SYSTEM ARCHITECTURE

This section provides a brief overview of the Whole Tale system architecture illustrated in Figure 2. Whole Tale provides a scalable platform based on the Docker Swarm container orchestration system, exposing a set of core services via REST APIs and Single Page Application (SPA). Key components include:

- **Whole Tale Dashboard:** An Ember.js single page application
- **Whole Tale API:** A REST API built using the Girder<sup>3</sup> framework to expose key features including authentication, user/group management, tale lifecycle, data management, and integration with remote repositories

<sup>2</sup><https://www.researchworkspace.com>

<sup>3</sup><https://girder.readthedocs.io>

- **Whole Tale File System:** A custom filesystem based on WebDav and FUSE used to mount user and registered data into running container environments
- **Image registry:** A local Docker registry used to host images associated with tales
- **Jobs and task Management:** A task distribution and notification framework based on Girder and Celery
- **Data Management System (DMS):** A system for fetching, caching, and exposing externally published datasets

Several aspects of the Whole Tale system are related to the BagIt-RO serialization format including filesystem organization, user-defined environments, metadata as well as the export and publication functions. We describe these in more detail below.

1) *Tale workspace:* Each tale has a *workspace* (folder) that contains user-created code, data, workflow, documentation and narrative information. The workspace also contains `repo2docker`-compatible configuration files defining the tale environment, described below. This appears as the *workspace* folder mounted into the running tale environment (i.e., container)

2) *External data:* Optionally, each tale can include references to externally published data. The data is then registered with the Whole Tale system and managed by the DMS. Externally referenced data appears in the *data* folder, a sibling to the *workspace*.

3) *Environment customization:* Users can optionally customize the tale environment using `repo2docker`-compatible configuration files. Whole Tale extends `repo2docker` via the `repo2docker_wholetale`<sup>4</sup> package, which adds buildpacks to support Rocker, Spark, and OpenRefine images.

4) *Metadata:* Tales have basic descriptive metadata including creator, authors, title, description, keywords as well as information about the selected environment, licenses, and associated persistent identifiers. The tale metadata is included in the metadata directory both in the `manifest.json` and `environment.json` files. The license is included in the BagIt payload directory, but not as part of the tale workspace.

5) *Exporting tales:* Tales can be exported in a BagIt-RO serialized archive that contains the contents of the tale workspace (code, local data, narrative, workflow, `repo2docker` configuration files) as well as references to external data, tale metadata, and a script to run the tale locally. BDBag [4] is used to materialize “holey” bags by downloading files specified in the `fetch.txt` file, initially via HTTP(S) and eventually via DOI, Globus, Agave schemes.<sup>5</sup> The script to run locally (`run-local.sh`) is stored at the root of the exported BagIt archive.

Table I describes the contents of an exported tale in the BagIt-RO format. A complete example is available at <https://doi.org/10.5281/zenodo.2641314>.

<sup>4</sup>[https://github.com/whole-tale/repo2docker\\_wholetale](https://github.com/whole-tale/repo2docker_wholetale)

<sup>5</sup>BagIt-Profile-Info example available at <https://raw.githubusercontent.com/fair-research/bdbag/master/profiles/bdbag-ro-profile.json>

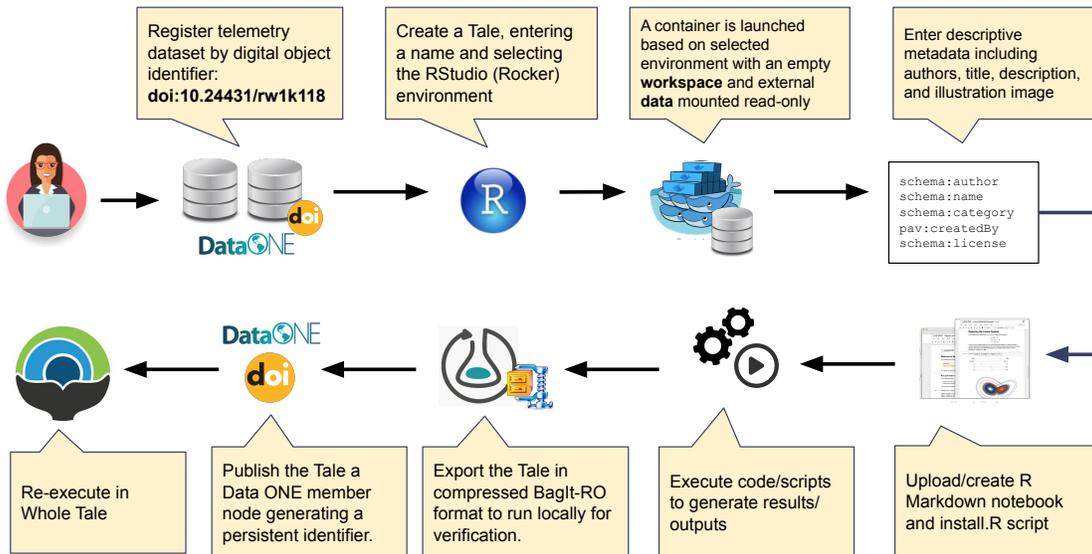


Fig. 1. Example Scenario Tale Creation and Publishing Workflow.

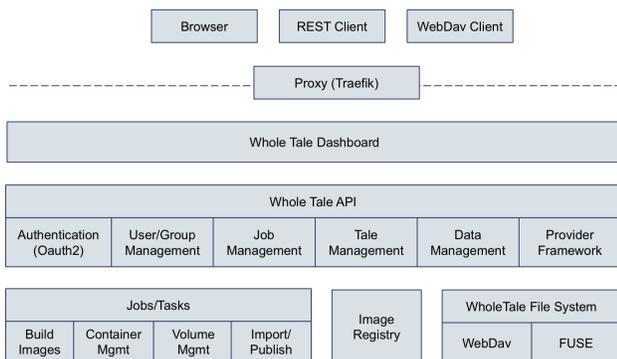


Fig. 2. Whole Tale System Architecture

#### IV. REQUIREMENTS

The scenario described in Section II highlights key requirements of the Whole Tale platform that led to our selection of the BagIt-RO serialization. These requirements include:

- 1) **Interoperability with archival repositories:** In the scenario, data is referenced from and the resulting tale published to an archival repository. In addition to DataONE network members, we are working on support for Dataverse network members, Dryad, and Zenodo. We must adopt standard formats and vocabularies to facilitate interoperability including the use of supported archival formats and identifiers (e.g., digital object identifiers).
- 2) **Interoperability with source code management (SCM):** Github is central to the workflow for the researchers in the scenario. The tale format must support publishing research objects that are based on content managed in SCM repositories.

File	Description
bag-info.txt	Bag metadata using the bdbag-ro-profile
bagit.txt	Bag declaration
data/ LICENSE workspace/ apt.txt postBuild requirements.txt wt_quickstart.ipynb	Payload directory containing tale license and workspace contents including repo2docker compatible configuration files.
fetch.txt	Fetch file
manifest-[md5, sha256].txt	Payload manifest (checksums)
metadata/ manifest.json environment.json	Tag directory containing RO manifest.json and Whole Tale environment metadata (required by repo2docker_wholetale)
tagmanifest-[md5, sha256].txt	Tag manifest (checksums)
README.md	Tale top-level readme
run-local.sh	Tale local execution script

TABLE I  
EXPORTED TALE CONTENTS

- 3) **Ability to reference external data:** The source dataset used in the scenario has been published in an archival repository. When running the tale via the Whole Tale web service or locally, externally referenced data must be resolved prior to re-execution. Whole Tale currently supports HTTP(S) resources as well as those published via Globus and in the future via the Agave Platform.
- 4) **Ability to add metadata:** The tale format must support all metadata attributes required by DataCite (<https://schema.datacite.org>) and schema.org (<https://schema.org/Dataset>) as well as attributes specific to the Whole Tale platform. In the future, we expect to also support additional metadata required by researchers

in specific domains.

- 5) **Ability to export and re-execute:** One feature of the system is the ability for users to export tales to a local machine. To re-run locally, we must be able to rebuild the environment (e.g., via Docker/repo2docker) and fetch remote data as needed.
- 6) **Simplicity and understandability:** When users view the contents of an exported or published tale, they should be able to easily understand the contents and how to explore or re-execute the tale.
- 7) **Interoperability with search engines:** Google recently unveiled Dataset Search which parses and aggregates JSON-LD embedded on dataset landing pages as an effort to lower barriers for finding datasets. Choosing JSON-LD as a representation for tale metadata provides flexibility in case we decide to expose tale information for Google. It also allows for further integration with third party publishers such as Dataverse and DataONE who may expose such metadata for Google.

The following requirements will be addressed in future releases of Whole Tale and relate to our selection of the BagIt-RO serialization format:

- 1) **Ability to store provenance information:** In future releases, tales will include computational and archival provenance information. We anticipate incorporating this information via standard models such as ProvONE [9].
- 2) **Verifiability:** Currently Whole Tale supports validation through re-execution of tales via a web based service or after exporting locally. Future releases will include information to allow the automatic re-execution and verification of included results/outputs and computational workflows.
- 3) **Versioning:** Since researchers iterate on their tales, share them and extend them, it is important to be able to version them over time.

In the next section, we discuss our adoption of the BagIt-RO model.

## V. ADOPTING THE BAGIT-RO MODEL

Whole Tale uses the RDF data model to encode tale information for export and exchange. We selected a JSON-LD representation for human readability, extensibility, compatibility with Whole Tale APIs, and potential interoperability with search engines and third party publishers. After developing an ad-hoc internal format, we explored emerging standards in the research object space and settled on BagIt-RO for serialization. Using the RO-Bundle specification and BagIt serialization in conjunction with the BDBag tools met many of our initial requirements. Additional tale metadata attributes which were not included in the BagIt-RO model could be added using vocabularies such as schema.org. Throughout this section, we use the `manifest.json` from the above example, with a complete listing included in Appendix A.

### A. Filesystem Artifacts

One strong point of RO-Bundle is that it treats file system artifacts as aggregates of the manifest. Doing so satisfies our requirement of being able to track where files belong, enabling us to both export and re-import tales even in the case where we must publish a hierarchical structure to a repository that can only represent a flat structure. In the case of Whole Tale, artifacts include data that were retrieved from external repositories as well as files that the user created or uploaded into the tale workspace. The tale workspace contents are included in the payload `data/workspace` directory and the external data are fetched into the payload `data/data` directory, mirroring filesystem organization on the web-based platform.

```
"aggregates": [
  {
    "uri": "../data/workspace/wt_quickstart.ipynb"
  },
  {
    "uri": "../data/workspace/apt.txt"
  }
]
```

Workspace artifacts are easily described with a single URI entry. Some files, such as the system generated `README.md` are tagged with additional metadata as shown below. In this case the additional metadata specifies the “type” of the file as a “HowTo”.

```
{
  "@type": "HowTo",
  "uri": "../README.md"
}
```

### B. External data

Whole Tale supports two types of external data: data that reside in a repository identified by persistent identifier (e.g., DOI) and data that exists at a generic HTTP(S) address. In addition to including information about external data in the `manifest.json`, the URL for each remote file, regardless of type, is included in the `fetch.txt` for retrieval using BDBag tools.

**Generic HTTP(S) Data:** For data that does not belong to a remote repository, a simple bundle is created in the aggregation section. The URI points to the HTTP(S) address where the file may be retrieved and the bundle object holds the information where the file should appear on the filesystem. This combination of information allows us to retrieve the file and place it in the correct folder (i.e., `data/data`).

**Repository Data:** For datasets that have been published to research repositories, additional metadata can be ingested when files are registered with the system. The individual files are described with a single bundle object, and linked to an additional structure that describes the dataset in more detail.

The following snippet describes a remote dataset that resides in DataONE and the aggregation recording the relationship between a file in that dataset and its ultimate location after retrieval in the payload “data” directory:

```
"dataset": [
  "@type": "Dataset",
```

```

"identifier": "doi:10.5065/D6862DM8",
"name": "Humans and Hydrology at High Latitudes...",
"@id": "doi:10.5065/D6862DM8"
],
"aggregates": [
  {
    "size": 1558016,
    "schema:isPartOf": "doi:10.5065/D6862DM8",
    "uri": "https://cn.dataone.org/cn/v2/resolve/urn
      :...",
    "bundledAs": {
      "filename": "usco2000.xls",
      "folder": "../data/data/"
    }
  }
]
]

```

### C. Describing the Computing Environment

Whole Tale uses a customized version of the Binder `repo2docker` package. In addition to including configuration files in the workspace, Whole Tale exports information about the environment including runtime information in the tale. One shortcoming of the BagIt-RO model is that there is no well-defined place for this metadata. To address this need, we define an additional tag file, `environment.json`, which encodes sufficient information about the environment so that it can be re-created. The metadata contained in this file is represented as JSON and is not yet described using standard vocabularies due as we were unable to identify a suitable convention.

### D. Describing Additional Attributes

A number of properties that describe additional tale attributes (e.g., authors, keywords, description, license) are defined at the manifest root. Schema.org's vocabulary is used to describe these general metadata fields.

Attributing authorship to a tale is a requirement for tracking researcher contributions and is also used during metadata generation with publishers. The Provenance, Authoring, and Versioning (PAV) vocabulary is used instead of schema because it is already included in by RO-Bundle:

```

{
  "@id": "https://orcid.org/0000-0002-7523-5539",
  "@type": "schema:Person",
  "schema:familyName": "DeBruine",
  "schema:givenName": "Lisa"
}

```

### E. Provenance Tracking

A planned feature of Whole Tale is the ability to track executions and steps in researchers' workflows based on techniques used to capture computational provenance [6], [15], [25], [26]. The BagIt-RO model includes the ability to provide provenance information through the inclusion of the `provenance.json` file. However, this is intended to capture more archival provenance information and it is unclear whether computational provenance should be included here. Whole Tale plans to use the ProvONE model [9], an extension to W3C PROV<sup>6</sup> derived from an earlier version [16], combining retrospective provenance and workflow models.

<sup>6</sup><https://www.w3.org/TR/prov-overview/>

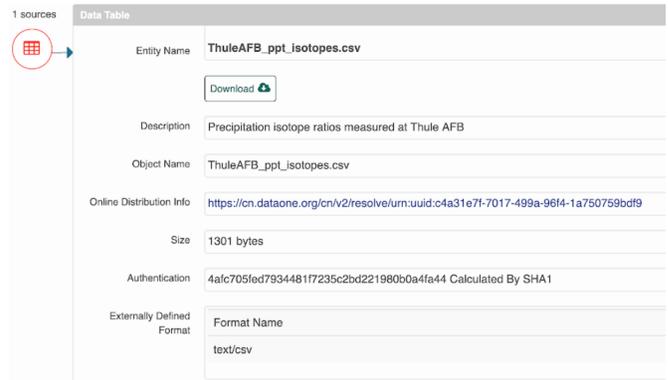


Fig. 3. Provenance rendering of a file in DataONE

The URI of each file in the manifest can be referenced inside the `provenance.json` file, enabling rich linkings of information. This information can also be transcribed to publisher-specific formats, provided that they support PROV. Figure 3 illustrates how provenance information is rendered in DataONE.

## VI. DISCUSSION

In this section, we highlight and discuss several issues related to our implementation of BagIt-RO that we hope will be of interest to workshop participants and possible input into current work on the RO-Crate specification. We discuss the importance of re-executability; the ability to reference and retrieve external data; the relationship between tales and source control repositories; and our ongoing work on computational provenance and verification workflows.

### A. Executable research objects

Tales are executable research objects. By this we mean that the research object itself may be built and re-executed for exploration, re-use, reproducibility, and verification. This is no longer a unique capability as many systems have recently been developed to support the creation of similar artifacts [6], [13], [26]. Executable research objects contain not only data, code, and documentation, but also information about the computational environment. This executability leads to additional capabilities, such as generation and comparison of computational provenance or methods of automated verification.

The FAIRDOM infrastructure initiative has made use of the Research Object framework to employ a standards based method to group its components into container platforms including BagIt [18]. We extend this approach into the Whole Tale framework and include the capability for externally referenced data and general research pipelines. Our efforts generalize those of ReproZip, which gathers and bundles dependencies for command line executions [6]. The Collective Knowledge (CK) framework gathers research objects with unique IDs and metadata in the JSON format but does not ensure re-executability [12]. Sciunits on the other hand are

self-contained bundles aimed to re-execute regardless of deployment, and targeted at scientific experiments [25], [26].

### B. External data

In the Whole Tale platform, users are presented with a fixed filesystem hierarchy that includes “workspace” and “data” directories. The workspace directory contains code, local data, and additional files (e.g., documentation) and the sibling “data” directory contains externally referenced data files (read-only).

In our v0.7 release, the BagIt payload directory of an exported tale similarly contains “workspace” and “data” directories. The `manifest.json` contains information about remotely registered datasets that is also included in the BagIt `fetch.txt`. When BDBag tools are used to fetch remote datasets, they are downloaded to the payload/data directory, matching the online filesystem organization and system capabilities. The concept of the `fetch.txt`, while primitive, is surprisingly effective when used with BDBag. We also foresee taking advantage of other BDBag capabilities, such as transferring Globus data or using DOI resolution. However, there is redundancy in tracking external information in both in the BagIt `fetch.txt` and the RO `manifest.json`.

### C. Relationship to SCM

Many researchers use source control repositories (e.g., GitHub) to organize and collaborate on research projects. Repositories can be released and published via external tools such as Zenodo or Whole Tale. In the Whole Tale platform, the “workspace” directory can be mapped to a version controlled repository. This raises the question of whether or not the workspace (or repository) should contain everything, including information currently stored in the `manifest.json` or `environment.json`. This information is essential to the understandability and re-executability of the tale, but is currently modeled as external to the primary tale contents (as is common with descriptive metadata). During the local execution process, for technical reasons we bind mount files from the “metadata” directory into the workspace to support building the tale image. In future releases, we are considering exposing the manifest information along with computational provenance information (below) as part of the workspace instead of external to it. This means that even simple metadata would be in the workspace and easily added to version control.

### D. Reproducibility and computational provenance information

Computational provenance refers to methods of capturing provenance (“the source or origin of an object”) for computational tasks [11] and is a subset of the larger notion of reproducibility of data- and computationally-enabled results [19], [21]–[24]. We are beginning to explore methods of capturing and storing computational provenance information to enable reproducibility on computational findings in tales. In the RO-Bundle specification, provenance information is defined as “describing creators, dates, and sources” and is more concerned with the provenance of the research object itself, which we term archival provenance. Computational

provenance information is internal to the tale and could be generated by the user or the Whole Tale system directly. We view computational provenance information as a key component of transparency for evaluation and verification of tales and part of enabling reproducibility.

### E. Supporting reproducibility via verification workflows

Research communities and journals are increasingly adopting artifact review processes that include re-execution of computational analysis in support of reproducibility [20]. Examples include the workflow implemented by the Odum Institute for the American Journal for Political Science [7], the Journal of the American Statistical Association<sup>7</sup>, Biostatistics [10], and the ACM Transactions on Mathematical Software (TOMS) Replicated Computational Results<sup>8</sup> program. We see tales and related research objects being used to simplify and possibly automate aspects of the verification process. Having a standard format for the exchange of research objects that fits into these enhanced curatorial and verification workflows may significantly reduce the burden on research communities.

### F. BagIt Understandability

One drawback of the BagIt serialization is that the BagIt configuration is foregrounded and difficult to understand for the average researcher/user while the “payload” directory, which contains their work, is less apparent and confusingly named “data”. Although out of scope for the RO discussion, we are supportive of the idea of a “.bagit” directory that contains the relevant configuration information and is largely hidden from the average user.

### G. Migrating to RO-Crate

Since our adoption of the BagIt-RO model, the community has moved forward on the Research Object Crate (RO-Crate) specification<sup>9</sup>. In this section, we report the results of a preliminary analysis of changes needed to migrate to the new format. Doing so will require versioning the tale export format and we are unlikely to make changes until the community settles on a near-final version of the specification.

RO-Crate 0.2-DRAFT introduces the following changes from the RO-Bundle 1.0:

- Addition of `ro-crate-metadata.jsonld` (RO-Crate Metadata File). The relationship to the RO-Bundle `manifest.json` is unclear, since the RO-Crate Metadata File “does not necessarily list or describe all files in the package.” We have viewed the `manifest.json` as an inventory of all files in the RO (excluding those introduced by BagIt).
- The RO-Crate metadata file changes vocabulary from the set used by RO-Bundle to primarily `schema.org`, no longer using `ore:aggregates`. This also adds support for referencing external datasets, a feature not available in RO-Bundle but added in our tale format.

<sup>7</sup><https://magazine.amstat.org/blog/2016/07/01/jasa-reproducible16/>

<sup>8</sup><http://toms.acm.org/replicated-computational-results.cfm>

<sup>9</sup><https://researchobject.github.io/ro-crate/>

- The “bagged” RO-Crate structure will differ from the BagIt-RO structure as the “metadata” folder is no longer included. Our assumption is that the `ro-crate-metadata.jsonld` along with our `environment.json` will now be included in the BagIt payload. We’ve come to a similar conclusion about the tale format – that this metadata belongs in the payload not external to it.
- It is unclear whether there will be support for separate provenance metadata or whether this will need to be included in the payload.

RO-Crate promises many benefits that align with Whole Tale, namely the adoption of `schema.org` as the primary vocabulary and its ability to be used alongside a variety of serialization formats.

## VII. CONCLUSIONS

By implementing an extension to RO-Bundle with BagIt serialization and leveraging existing open science infrastructure tools including `repo2docker` and `BDBag`, we were able to effectively create an exportable, publishable, and executable research object package, in short taking a step toward the publication of “really reproducible research” [8]. While not a perfect fit, BagIt-RO met many of our platform requirements. We expect to continue work in this area as we add support for computational provenance information and automated verification and hope to contribute to the use cases and discussions that inform the development of a broader community standard.

## ACKNOWLEDGMENT

This work is supported by National Science Foundation Award OAC-1541450.

## REFERENCES

- [1] C. Boettiger and D. Eddelbuettel. An introduction to rocker: Docker containers for R. *CoRR*, abs/1710.03675, 2017.
- [2] A. Brinckman, K. Chard, N. Gaffney, M. Hategan, M. B. Jones, K. Kowalik, S. Kulasekaran, B. Ludäscher, B. D. Mecum, J. Nabrzyski, et al. Computing environments for reproducibility: Capturing the “whole tale”. *Future Generation Computer Systems*, 94:854–867, 2019.
- [3] M. Cameron, J. London, K. Frost, A. Whiting, and P. Boveng. Satellite Telemetry Dataset (Raw): Juvenile Bearded and Spotted Seals, 2004–2006. Kotzebue, Alaska, 2018.
- [4] K. Chard, M. D’Arcy, B. Heavner, I. Foster, C. Kesselman, R. Madduri, A. Rodriguez, S. Soiland-Reyes, C. Goble, K. Clark, E. W. Deusch, I. Dinov, N. Price, and A. Toga. I’ll take that to go: Big data bags and minimal identifiers for exchange of large, complex datasets. In *2016 IEEE International Conference on Big Data (Big Data)*, pages 319–328, Dec 2016.
- [5] K. Chard, N. Gaffney, M. B. Jones, K. Kowalik, B. Ludäscher, J. Nabrzyski, V. Stodden, I. Taylor, M. J. Turk, and C. Willis. Implementing computational reproducibility in the whole tale environment. In *Proceedings of the 2nd International Workshop on Practical Reproducible Evaluation of Computer Systems*, P-RECS ’19, pages 17–22, New York, NY, USA, 2019. ACM.
- [6] F. Chirigati, R. Rampin, D. Shasha, and J. Freire. Reprozip: Computational reproducibility with ease. In *Proceedings of the 2016 International Conference on Management of Data*, SIGMOD ’16, pages 2085–2088, New York, NY, USA, 2016. ACM.
- [7] T.-M. Christian, S. Lafferty-Hess, W. G. Jacoby, and T. Carsey. Operationalizing the replication standard. *IJDC*, 13(1):114–124, 2018.
- [8] J. F. Claerbout and M. Karrenbach. Electronic documents give reproducible research a new meaning. In *SEG Technical Program Expanded Abstracts 1992*, pages 601–604. Society of Exploration Geophysicists, 1992.
- [9] V. Cuevas-Vicentín, B. Ludäscher, P. Missier, K. Belhajjame, F. Chirigati, Y. Wei, S. Dey, P. Kianmajid, D. Koop, S. Bowers, I. Altintas, C. Jones, M. B. Jones, L. Walker, P. Slaughter, B. Leinfelder, and Y. Cao. Provone: A prov extension data model for scientific workflow provenance. <https://purl.dataone.org/provone-v1-dev>, May 2016.
- [10] D. L. Donoho. An invitation to reproducible computational research. *Biostatistics*, 11(3):385–388, 07 2010.
- [11] J. Freire, D. Koop, E. Santos, and C. T. Silva. Provenance for computational tasks: A survey. *Computing in Science & Engineering*, 10(3):11–21, May 2008.
- [12] G. Fursin, A. Lokhmotov, D. Savenko, and E. Upton. A collective knowledge workflow for collaborative research into multi-objective autotuning and machine learning techniques. *CoRR*, abs/1801.08024, 2018.
- [13] Jupyter-Project. Binder 2.0 - reproducible, interactive, sharable environments for science at scale. 17th Python in Science Conference, 2018.
- [14] J. M. London and D. S. Johnson. Alaska bearded and spotted seal example dataset and analysis. <https://github.com/jmlondon/crwexamplekbs>, 2019.
- [15] T. McPhillips, S. Bowers, K. Belhajjame, and B. Ludäscher. Retrospective provenance without a runtime provenance recorder. In *7th Intl. Workshop on Theory and Practice of Provenance (TaPP)*, 2015.
- [16] P. Missier, S. Dey, K. Belhajjame, V. Cuevas-Vicentín, and B. Ludäscher. D-PROV: Extending the PROV Provenance Model with Workflow Structure. In *Proc. 5th Workshop on the Theory and Practice of Provenance (TaPP)*, 2013.
- [17] S. Soiland-Reyes, M. Gamble, and R. Haines. Research object bundle 1.0, `researchobject.org` recommendation. <https://w3id.org/bundle/2014-11-05/>, 2014.
- [18] N. Stanford, F. Bacall, M. Golebiewski, O. Krebs, R. Kuzyakiv, Q. Nguyen, S. Owen, S. Soiland-Reyes, J. Straszewski, D. van Niekerk, A. Williams, K. Wolstencroft, L. Malmström, B. Rinn, J. Snoep, W. Müller, and C. Goble. FAIRDOME: Reproducible Systems Biology through FAIR Asset Management. In *Reproducibility, Standards and SOP in Bioinformatics: Combined CHARME – EMBnet and NETTAB Workshop*, 2016.
- [19] V. Stodden. Reproducible research: tools and strategies for scientific computing. *Computing in Science and Engineering*, 14:11–12, 2012.
- [20] V. Stodden, P. Guo, and Z. Ma. Toward reproducible computational research: an empirical analysis of data and code policy adoption by journals. *PLOS ONE*, 8(6):e67111, June 2013.
- [21] V. Stodden, F. Leisch, and R. D. Peng. *Implementing Reproducible Research*. CRC Press, Apr. 2014.
- [22] V. Stodden, M. McNutt, D. H. Bailey, E. Deelman, Y. Gil, B. Hanson, M. A. Heroux, J. P. Ioannidis, and M. Taufer. Enhancing reproducibility for computational methods. *Science*, 354(6317):1240–1241, 2016.
- [23] V. Stodden and S. Miguez. Best practices for computational science: Software infrastructure and environments for reproducible and extensible research. *Journal of Open Research Software*, 2, 2014.
- [24] V. Stodden, S. Miguez, and J. Seiler. Researchcompendia.org: Cyberinfrastructure for reproducibility and collaboration in computational science. *Computing in Science and Engineering*, 17(1):12–19, 2015.
- [25] D. H. T. That, G. Fils, Z. Yuan, and T. Malik. Sciunits: Reusable research objects. *CoRR*, abs/1707.05731, 2017.
- [26] Z. Yuan, D. H. T. That, S. Kothari, G. Fils, and T. Malik. Utilizing provenance in reusable research objects. *Informatics*, 5:14, 2018.

## VIII. APPENDIX A

"@id": "https://data.wholetale.org/api/v1/tale/5  
cb4ffead9323600016c4d4c"

```
{
  "createdBy": {
    "@type": "schema:Person",
    "schema:givenName": "Craig",
    "@id": "willis8@illinois.edu",
    "schema:email": "willis8@illinois.edu",
    "schema:familyName": "Willis"
  },
  "schema:description": "Demonstration of how to use
  Whole Tale to develop custom analysis and
  visualization for data published externally via
  DataONE. See https://wholetale.readthedocs.io/en/
  stable/users_guide/quickstart.html for more
  information.",
  "@context": [
    "https://w3id.org/bundle/context",
    {
      "schema": "http://schema.org/"
    },
    {
      "Datasets": {
        "@type": "@id"
      }
    }
  ],
  "schema:author": [
    {
      "@type": "schema:Person",
      "schema:givenName": "Craig",
      "@id": "https://orcid.org/0000-0002-6148-7196",
      "schema:familyName": "Willis"
    }
  ],
  "schema:version": 7,
  "schema:identifier": "5cb4ffead9323600016c4d4c",
  "schema:image": "http://use.yt/upload/dclida723",
  "Datasets": [
    {
      "@type": "Dataset",
      "identifier": "doi:10.5065/D6862DM8",
      "name": "Humans and Hydrology at High Latitudes
      : Water Use Information",
      "@id": "doi:10.5065/D6862DM8"
    }
  ],
  "createdOn": "2019-04-15 22:04:26.970000",
  "schema:name": "Example Water Tale",
  "schema:category": "Examples",
  "aggregates": [
    {
      "uri": "../data/workspace/wt_quickstart.ipynb"
    },
    {
      "uri": "../data/workspace/apt.txt"
    },
    {
      "uri": "../data/workspace/requirements.txt"
    },
    {
      "uri": "../data/workspace/postBuild"
    },
    {
      "size": 1558016,
      "schema:isPartOf": "doi:10.5065/D6862DM8",
      "uri": "https://cn.dataone.org/cn/v2/resolve/
      urn:uuid:62e1a8c5-406b-43f9-9234-1415277674
      cb",
      "bundledAs": {
        "filename": "usco2000.xls",
        "folder": "../data/data/"
      }
    },
    {
      "schema:license": "CC-BY-4.0",
      "uri": "../data/LICENSE"
    },
    {
      "@type": "schema:HowTo",
      "uri": "../data/README.md"
    }
  ],
}
```