

# The International Journal of Digital Curation

## Volume 8, Issue 1 | 2013

### Generation of a Skeleton Corpus of Digital Objects for the Validation and Evaluation of Format Identification Tools and Signatures

Ross Spencer,

Independent Digital Preservation Researcher

#### Abstract

To preserve digital information it is vital that the format of that information can be identified, in-perpetuity. This is the major focus of research within the field of Digital Preservation. The National Archives of the UK called for the Digital Preservation and Digital Curation communities to develop a test corpus of digital objects to help further develop tools to aid this purpose. Following that call, an attempt has been made to develop the suite.

This paper initially outlines a methodology to generate a skeleton corpus using simple user-generated digital objects. It then explores the lessons learnt in the generation of a corpus using scripting language techniques from the file format signatures described in The National Archives PRONOM technical registry. It will also discuss the use of the digital signature for this purpose, the benefits of developing a test corpus using this technique. Finally, this paper will outline a methodology for future research before exploring how the community can best make use of the output of this project and how this project needs to be taken forward to completion.





## Introduction

The original ‘call-to-arms’ by The National Archives (Fetherston and Gollins, [2012](#)) described the development of a test corpus of digital objects that could be used for the validation and evaluation of format identification tools and techniques. The benefits of such a suite for the community are numerous. The suite would create a baseline for testing format identification tools against each other, create a baseline for improving performance, and generate a mechanism by which identification and validation techniques could be improved. Like the PUID in PRONOM<sup>1</sup>. It would allow developers and researchers alike to speak unambiguously about the formats they have used to create new techniques, offering greater transparency for those who want to accurately recreate and verify results and measurements.

Fetherston and Gollins highlighted the difficulty of setting up a fully-formed test corpus and the considerations anyone taking up the challenge would need to make. While this paper makes no attempt to describe the successful generation and publication of a fully-formed test corpus, it does describe the generation of a skeleton corpus of objects for testing format identification tools, which is expected to provide a sustainable and practical solution for a number of issues described in the original paper. While this paper offers an interim solution for the ‘call-to-arms’ it is expected that it will continue to compliment the existence of the fully-formed suite, should such a corpus of objects become available.

## The Skeleton Corpus

During the development of the DROID 6.1 project at The National Archives, mechanisms were required for the testing of the identification engine. This would take the form of unit tests developed in JUnit.<sup>2</sup> For example, a unit test might be written to test the identification engine’s ability to parse syntactical constructs. Such a test would involve the matching of a digital object against a digital signature. Where bugs had previously been raised, the best digital object to test against would be the original that caused the problem. Using a test-driven development approach allows the developer to confirm the existence of the bug and demonstrate it by creating a unit test that fails under the initial conditions. Once modifications have been made to the engine to correct the issue, the unit test will eventually pass.

The digital object needs to be distributed with the unit test. The package, therefore, gets distributed with the source code. This scenario faces two of the issues highlighted by Fetherston and Gollins ([2012](#)). The **usability** of the source code is impacted, as the size of the source code becomes bloated by the size of the digital objects that need to be distributed with DROID. Further, **intellectual property rights (IPR) and copyright** are not easily addressed. In this instance, The National Archives needs to hold the IPR and copyright for any problem files they might want to distribute with the source code. This isn’t always possible if files from long obsolete systems have

---

<sup>1</sup> Description of the PUID scheme in Wikipedia:

[http://en.wikipedia.org/wiki/PRONOM\\_technical\\_registry#The\\_PRONOM\\_Persistent\\_Unique\\_Identifier\\_28PUID.29\\_scheme](http://en.wikipedia.org/wiki/PRONOM_technical_registry#The_PRONOM_Persistent_Unique_Identifier_28PUID.29_scheme)

<sup>2</sup> Java unit testing framework: <http://junit.org>





(Brown, 2006). As well as being a syntax to create signatures, we can use it to deconstruct signatures into one or more byte sequences that can then be fed into a file, which will then become a skeleton file for the test corpus.

To illustrate, two basic signatures might look like this:

1. CAFEBAFE[01:0F]CAFED00D
2. BAADCAFE(0A|0D0A)BAADF00D

The syntax in the first example describes a “sequence of bytes which lies lexicographically between 01 and 0F,” that is, a run of numerical values. If we deconstruct the signature, the bytes CAFEBAFE and CAFED00D can be placed directly into a new file. The bracketed syntax needs to be parsed and a decision made in the code about outputting an integer between one and 15 between the two sequences.

In the second example, the syntax describes a “wildcard matching one from a list of values,” that is, a list of optional values separated by the OR ‘|’ symbol. Again, deconstructing this signature sees the sequences of BAADCAFE and BAADF00D output directly to the new file. The bracketed syntax in this case needs to be parsed and a decision made by the processing engine over which *additional* sequence to select to generate the skeleton file. We need to understand if the processing engine should output files *for each* option documented in a signature in this way.

Generating files automatically isn’t recommended for each file created for the suite, as the concept of these test files is for them to be handwritten by developers of signatures, closely mirroring the work done to generate a signature, to provide a solid human-based provenance trail for each file in the suite. The recommended use of this approach is to provide a baseline for the test suite by generating the suite from the PRONOM database as it is at present and then seeking for it to be supported by handmade files in future. The automated tool should exist to allow for bulk generation of signatures where a lot of work on signatures has already been completed but without the files being created alongside. Its use should be minimized to improve the quality of the test suite and ensure its independence from the data that is used by DROID.

### Structure of the File

The tool described above has been developed as a prototype to enable this paper to be written and to enable the challenges involved to be better understood.<sup>4</sup> Initially, it was thought to be extremely clear how skeleton files should look. Their structure should be simple, and to support the sequences in the file they should be padded by null bytes to clearly distinguish the different sections that make up a signature. Furthermore, the files created should be identifiable as their *own format*. We could give files their own magic number to identify the file type, and also create a lightweight structure inside the file to store metadata, such as the PUID for the format that the skeleton file represents.

---

<sup>4</sup> Python scripts for the Skeleton-test-suite-generator hosted on GitHub:  
<https://github.com/exponential-decay/skeleton-test-suite-generator>

Through this development, a number of issues came up which need to be considered before finalizing a suitable structure. One concern includes the multiple internal signatures that can belong to a single format in PRONOM. This means that the relationship between a PUID and a skeleton file is one-to-many (1:\*)<sup>5</sup>; this needs to be recorded in metadata if we want self-describing files. We must also consider what to do with multiple combinations of byte-sequences. In the example above a ‘Max offset’ of four means that there are four other combinations of sequence (and therefore file) that can match the signature. The purest skeleton-suite would contain all variants described by a signature, but what are the practicalities of this when we get into offsets in the hundreds or thousands of byte ranges? What should we do pragmatically that gives us the best return for our effort?

Deconstruction of signatures and generation of files proved to be more complex than anticipated. As well as dealing with multiple combinations of sequence, the code for outputting the skeleton files needs to consider what to do with wildcards, how many extra bytes to pad files with, and how it fills bytes between sequences that are statically positioned (i.e. Beginning of File (BOF) and End of File (EOF) sequences), and those which are variably positioned but bounded by two offsets.

At the time of writing, v0.2-BETA<sup>5</sup> of the skeleton-test-suite-generator outputs only the most basic skeleton-file, i.e. files matching byte-sequences described in PRONOM ‘as-is’ with minimal padding between variable positioned sequences and no additional metadata. To identify individual files we adopt the naming convention ‘*[puid]-[pronom-signature-id].[file-extension]*’. This makes it easier to debug output and keep track of permutations of file described by multiple internal signatures.

Over the two iterations of development thus far, the focus has been on the accurate output of internal signatures as skeleton bitstreams. Further development of the skeleton-test-suite-generator will help us to understand the flexibility we have in order to include greater structure, output identifying sequences for skeleton-files and additional self-describing metadata without interfering with the essence of the files themselves.

### **Proposed Methodology for Organizations Contributing Signatures**

Once a baseline suite exists, the proposed methodology will be for developers creating signatures to create matching skeleton files. When testing signatures it is important for the developer to test signatures against their own sample files, the skeleton file, and the complete skeleton test corpus. This will enable them to check for collisions in the work they are doing. New signatures should not match any existing file in the test suite. Should this happen, the developer should think about modifying their signature or look for changes that can be submitted for the existing signature and corresponding skeleton file causing the issue.

Once the testing has been completed, the DROID signature and skeleton file can then be submitted to The National Archives and to the organization/community controlling the test suite so it can be added to the PRONOM database *and* skeleton corpus respectively.

---

<sup>5</sup> Skeleton-test-suite-generator v0.2-BETA on GitHub:  
<https://github.com/exponential-decay/skeleton-test-suite-generator/tree/v0.2-BETA>

For The National Archives developing signatures on a roughly bi-monthly release cycle, every attempt should be made for a skeleton file to be uploaded to the corpus for each new signature created or each signature modified. Where large numbers of signatures are created in a single release, the automated generation of these digital objects may be appropriate. However, this approach should be made clear in the provenance trail of the test suite and every attempt should be made to generate files by hand in all other cases.

### Licensing the Skeleton Suite

The issues surrounding IPR, copyright and access rights expressed by Fetherston and Gollins (2012) are perhaps the biggest concerns for a developer hoping to benefit from a fully-featured corpus of digital objects. These concerns stem from the perceived lack of digital preservation researchers with legal experience. It is clear that for a full test suite to succeed, the committee or organization maintaining it needs a clearly defined **written** policy which states who owns the suite, what can be done with the suite, and how it can be reused. It must also clearly describe the process for contributing to the suite and the transfer of ownership of the files to the suite holders.

The skeleton suite being much simpler in nature makes the legal issues surrounding it easier to comprehend. In adopting the open source community as a model, we are provided with licenses that will be appropriate for this work.

As the skeleton suite should have minimal restrictions placed upon it, the proposed license under which the suite should be placed is ‘Creative Commons Share Alike’ (CC BY-SA)<sup>6</sup>. This should be supported by a simple waiver that is accepted by anyone contributing simple skeleton files to the suite, allowing the resource to then be licensed in this manner.

Summarizing various licenses under which research data can be released, the Digital Curation Centre (DCC) provides some useful guidance that allows us to be confident of the selection of CC BY-SA (Ball, 2012). Firstly, CC BY-SA “allows others to copy, distribute, display and perform the work as long as the creator is given due credit.” Secondly, Share Alike “inserts a strong copyleft clause into the license. This means that all derivative works must be released under the same license as the original work”. Finally, the DCC note that: “a quite general problem is that the licenses are aimed at homogeneous works.”

The attribution clause in CC BY-SA means that the provenance of the suite should remain strong through potentially numerous generations of the suite by many different organizations. The ‘Share Alike’ clause ensures that organizations are free to share and add to the suite. They are also able to release it for themselves, as long as they adopt the same license. This clause is expected to benefit the community, as it allows anyone to adopt and expand the skeleton suite for other research purposes, and to re-release it without worrying about doing so. The rest of the community can then benefit from the improved collection, provided the license is followed correctly. Further still, as we are creating large volumes of homogeneous objects – objects that do not conform to a single specification individually, but all conform to the

---

<sup>6</sup> Creative Commons CC BY-SA license text: <http://creativecommons.org/licenses/by-sa/3.0/>



With GitHub we can create a centralized host for this work containing a clear description about what the suite does, including licensing information. We are able to manage binary objects in the master branch, and for each release of the suite (per DROID signature file release) we are able to tag the test suite, package it, and enable it to be downloaded by the community. Provenance for the files in the suite is provided by the commit log and commit messages. A clear log must be written for each file committed; e.g. '*[new file] fmt/445 submitted by The National Archives for Signature file V62.*' This approach can also be used to track changes and updates to files. Further, the issues tracker will enable users of the suite to report errors in their findings and monitor the work being done. Supported by a Google Group for the test suite, a community could be built around it, promoting an environment for collaborative work as described by Fetherston and Gollins (2012).

## Recommendations: What next?

The next steps will be to complete the tool, to create a full baseline test suite and host the results on GitHub. The code will need to be thoroughly tested and unit tests created to provide assurances that the code outputs files accurately reflecting the PRONOM signatures. The testing of digital objects will also be important, and the end result we will be expecting is a one-to-one relationship between a file and an identification in DROID. Any collisions found will require investigation to determine whether the file is inaccurate or whether the signature as it is stored in PRONOM needs additional work because it is returning false-positives.

Beyond the tool itself, the community *and* The National Archives will need to decide the best way forward to continue to use the results of this work. Incorporation of the skeleton test corpus into the DROID unit tests will be an important first step, as this will give developers extra freedom and confidence to modify code, thus giving developers and DROID users assurances that identification capabilities remain consistent. It will be important for DROID to distinguish the difference between its own license (BSD 3-clause) and that of the data used for the unit tests (i.e. the skeleton corpus), but we do not anticipate any conflicts between the two sets of work.

Once the suite has been perfected for DROID, other identification tools can then be explored and we can begin to understand whether this approach will bring similar advantages to users and developers of those tools as well.

Container formats will be difficult to create in the skeleton suite initially. While creating zip-based analogues and mirroring the structure of those files for identification by DROID will be straightforward due to the accessibility of zip software, the second container type identified by DROID, OLE2, will take additional research to understand. It is conceivable that this cannot be done without the aid of a specialist tool to create the OLE2 structure. This must be understood but should not detract from the other work we are capable of doing.

Finally, during the process of prototyping the automated tool to output byte sequences in PRONOM as skeleton files, it was noticed that small improvements could be made to the quality of the information in the PRONOM database – specifically the names assigned to 'Internal Signatures' as these are different from the



