

A workflow reproducibility scale for automatic validation of biological interpretation results

Hiroataka Suetake ¹, Tsukasa Fukusato ², Takeo Igarashi ³ and Tazro Ohta ^{4,*}

¹Department of Creative Informatics, Graduate School of Information Science and Technology, The University of Tokyo, Tokyo, 113-0033, Japan

²Department of Computer Science, Graduate School of Information Science and Technology, The University of Tokyo, Tokyo, 113-0033, Japan

³Department of Creative Informatics, Graduate School of Information Science and Technology, The University of Tokyo, Tokyo, 113-0033, Japan

⁴Database Center for Life Science, Joint Support-Center for Data Science Research, Research Organization of Information and Systems, Shizuoka, 411-8540, Japan

*Correspondence address. Tazro Ohta, Database Center for Life Science, Joint Support-Center for Data Science Research, Research Organization of Information and Systems, Yata 1111, Mishima, Shizuoka 411-8540, Japan. E-mail: t.ohta@dbcls.rois.ac.jp

Abstract

Background: Reproducibility of data analysis workflow is a key issue in the field of bioinformatics. Recent computing technologies, such as virtualization, have made it possible to reproduce workflow execution with ease. However, the reproducibility of results is not well discussed; that is, there is no standard way to verify whether the biological interpretation of reproduced results is the same. Therefore, it still remains a challenge to automatically evaluate the reproducibility of results.

Results: We propose a new metric, a reproducibility scale of workflow execution results, to evaluate the reproducibility of results. This metric is based on the idea of evaluating the reproducibility of results using biological feature values (e.g., number of reads, mapping rate, and variant frequency) representing their biological interpretation. We also implemented a prototype system that automatically evaluates the reproducibility of results using the proposed metric. To demonstrate our approach, we conducted an experiment using workflows used by researchers in real research projects and the use cases that are frequently encountered in the field of bioinformatics.

Conclusions: Our approach enables automatic evaluation of the reproducibility of results using a fine-grained scale. By introducing our approach, it is possible to evolve from a binary view of whether the results are superficially identical or not to a more graduated view. We believe that our approach will contribute to more informed discussion on reproducibility in bioinformatics.

Keywords: workflow, provenance, reproducibility

Background

Bioinformatics is big data science and is considered the most demanding domain in terms of data acquisition, storage, distribution, and analysis [1]. Because the low cost and high throughput of measurement instruments have made it possible to generate large amounts of data, large-scale data analysis using a computer is required to extract valuable knowledge from the data [2, 3]. For each data analysis process, such as data transformation, public database referencing and merging, and statistical processing, much open-source software is developed and released [4]. Researchers typically choose appropriate software for each analysis process, build a workflow by combining the software, and execute the workflow in a computing environment [5]. However, it can be challenging to ensure the reproducibility of data analysis due to a number of factors, such as a large amount of data, the diversity of data types and software, and the complexity of the computing environment [6].

Reproducibility of research is an essential issue in the scientific community [7, 8]. However, Baker raised the alarm of a “reproducibility crisis” based on survey results that “more than 70% of researchers have tried and failed to reproduce another scientist’s experiments, and more than half have failed to reproduce their own experiments” ([9], p. 452). The key here is the requirement for research to be considered reproducible. Drummond argued that replicability and reproducibility are often confused, but they

are different concepts and need to be clearly distinguished [10]. The Association for Computing Machinery (ACM) also attempts to define the terms *repeatability*, *reproducibility*, and *replicability* (Table 1) [11]. While these definitions are in the context of computerized analysis, it should be noted that most existing studies have focused on whether the execution can be reproduced or not and have not considered the verification of the results. That is, they only state that the resulting data are exactly the same as in the original but do not adequately discuss the verification of whether the results are reproducible or not. Therefore, the reproducibility of data analysis can be divided into 2 parts: the execution of the analysis and the verification of the results. We will focus our discussion on the second part, verification.

Many workflow systems have been developed to improve the efficiency of building and executing complex data analysis [12–14]. Each system has unique characteristics, but in particular, workflow systems can have a syntax for describing the data analysis, called a workflow language. Large user communities have been formed around these workflow languages. The Common Workflow Language (CWL) [15], the Workflow Description Language (WDL) [16], Nextflow [17], and Snakemake [18] are typical examples. These systems also have execution systems that work with computational frameworks, such as job schedulers, container runtimes, and package managers. Thus, these workflow systems facilitate the execution of data analysis by different

Received: November 2, 2022. Revised: January 26, 2023. Accepted: April 28, 2023

© The Author(s) 2023. Published by Oxford University Press GigaScience. This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

Table 1: Repeatability, reproducibility, and replicability. According to the ACM [11], repeatability is defined as a researcher can reliably repeat their own computation. Reproducibility is defined as an independent group can obtain the same result using the author's own artifacts, and replicability is defined as an independent group can obtain the same result using artifacts that they develop completely independently.

Term	Team	Environment	Setup (code and data)	Result
Repeatability	Same	Same	Same	Same
Reproducibility	Different	Different	Same	Same
Replicability	Different	Different	Different	Same

teams and in different environments through the use of virtualization technology and syntax that abstracts software and computational requirements [19].

The advent and widespread use of workflow systems have facilitated data analysis reexecution. However, as mentioned earlier, to ensure reproducibility, it is necessary to verify the execution results, that is, whether the same biological interpretation is obtained or not. To address this issue, frameworks, such as Research Object Crate (RO-Crate) [20] and CWLProv [21], have been proposed to generate workflow provenance, a structured archive that packages workflow-related metadata, such as workflow descriptions, execution parameters, input and output data, tests, and documentation, in a machine-readable format. This provenance information is distributed on workflow sharing platforms, such as WorkflowHub [22], Dockstore [23], and nf-core [24]. When appropriate provenance is provided by the author, a researcher can use this information to verify new execution results, making the process reproducible.

However, the process of comparing the provenance and execution results is often incomplete and inefficient. In automatic comparison, the checksums of the output files are used; however, they do not always match. This is because these checksums may differ depending on the software version, timestamps, heuristic algorithms, and computing environments (e.g., OS and CPU architecture) [25]. However, the same biological interpretation may be obtained even when the output files do not match exactly; for example, only the timestamps in the output files may differ. Thus, a simple comparison using a checksum is incomplete in verifying results. Another method is to have humans semantically interpret the results. However, due to its inefficiency, this method is not possible when the scale of the data analysis execution is large. From the above, the verification of results using provenance remains challenging because the current procedure is limited to incomplete automatic comparison and inefficient manual checking.

Automation is essential for the verification of practical workflows that output many files; however, binary determination by checking checksums is not sufficient. Thus, it is necessary to introduce a fine-grained scale to determine the degree of reproducibility of the results. Automatic verification of results using this scale will make verification of workflow reproducibility practical. In this article, we propose a reproducibility scale of workflow execution results based on some discussion and experiments and a validation method using this scale. We implement a workflow execution system that generates a workflow provenance that contains metadata required for verification. This implementation is an extension of Sapporo [26] (RRID:SCR_023202), an existing workflow execution service (WES). Sapporo's extensibility makes it compatible with various workflow languages and execution systems. In addition, we adapt RO-Crate as the workflow provenance format. We also develop Tonkaz: a command-line tool that verifies the reproducibility of data analysis results by comparing the

workflow provenances. To demonstrate the effectiveness of our approach, we apply it to workflows used by researchers in real research projects. The full reproducibility of research is still an issue that has not been fully resolved. Nevertheless, we hope our approach will contribute to solving this problem by increasing the resolution of the definition of reproducibility.

Methods

Reproducibility scale of workflow execution results

A workflow is a sequence of computational steps that combine analysis tools according to their inputs and outputs. The first tool takes input data and passes its output on as input for the next tool. Thus, the result of the workflow execution is the cumulative output of each tool or the last tool in the workflow. It should be noted that the output of a tool includes not only output files but also execution logs (e.g., standard output and error) and runtime information (e.g., exit code, start time, and end time). Returning to the purpose of data analysis here, it is to obtain useful biological knowledge from the data. Therefore, it is not sufficient to consider the output files and logs as the only result of the workflow execution; the biological feature values interpreted from the output files and logs should be considered the result of the workflow execution.

The format to represent biological features obtained from data analysis is not standardized and varies depending on the analysis tool. For example, there are summarized formats (tabular and graph) and formats that express biological features themselves, such as Sequence Alignment/Map (SAM) and the Variant Call Format (VCF). To interpret and verify the results, the individual executing the workflow visually checks the output graph or uses a tool to extract a numerical feature value from the file, for example, SAMtools [27] (RRID:SCR_002105), to extract mapping statistics from the SAM format. Because these processes require domain knowledge, it is ideal that the workflow itself provides a structured summary and a way to interpret it. However, this depends on the skill and effort of the individual workflow developer, and the diversity of tools and workflows makes it challenging to provide them in a standardized way.

Several workflows provide a way to verify reproducibility using biological feature values. For example, the RNA sequencing (RNA-seq) workflow [28] distributed by the nf-core project has a test mode to verify that the workflow is working as expected. In this mode, the workflow is executed with a small test dataset, and the biological feature values are compared with the expected values. The mapping rate, which represents the percentage of reads that are mapped to the reference genome, is used as a biological feature value. If the difference between the values is within the threshold, the workflow is considered to be working as expected. As a preliminary experiment, we compared the output files with-

out using such biological feature values—that is, the checksum method was used to verify an exact match of a file. As a result, when we executed the above RNA-seq workflow twice in the same environment and compared the file output BAM files (the compressed binary version of the SAM files), we found that the checksum values were different and the file sizes differed by several bytes (see the “Results” section for details). It is ideal that the output files are exactly the same, but it is difficult to achieve this goal because the output files are generated by the analysis tools, and these tools are not designed to produce the same output all the time. Therefore, we concluded that it is not sufficient to check only the exact match of output files to verify the reproducibility of workflow execution results and that a method using biological feature values and threshold should also be introduced.

Based on the above discussion and preliminary experiments, we propose a method to verify the reproducibility of workflow execution results using biological feature values and threshold. The method consists of 2 steps: (i) extracting biological feature values from the output files and logs and (ii) comparing the extracted biological feature values with the expected values using threshold values. A detailed description of each step is provided in the “Generation of workflow provenance containing biological feature values” and “Automatic verification of reproducibility” sections. We also propose a scale to evaluate the reproducibility of workflow execution results based on the method (Table 2). This allows the reproducibility of results to be expressed at a higher resolution than a binary measure of whether the results are the same or not.

Generation of workflow provenance containing biological feature values

To verify the reproducibility of workflow execution results using biological feature values, it is necessary to package the workflow execution results as the workflow provenance in a standardized format. Because there are many workflow languages and execution engines, we first abstracted the workflow execution itself. Thus, we extended Sapporo, an existing WES implementation. Sapporo has an API scheme that satisfies the Global Alliance for Genomics and Health (GA4GH) WES standard [29], enabling the workflow execution and results acquisition in a standardized manner. In addition, due to its extensibility, it can execute workflows written in various languages, such as CWL, WDL, Nextflow, and Snakemake. Therefore, by extending Sapporo, workflow execution written in various languages can be handled in the same way.

When a workflow is executed in Sapporo, the files related to the execution are stored in the file system as workflow provenance. This provenance directory contains the workflow definition files, input files, intermediate files, output files, log files, execution parameters, runtime information, and so on. Thus, we converted Sapporo’s provenance into RO-Crate [20], a standardized format for packaging research objects expressed in JSON-LD. Because the RO-Crate use case included the packaging of workflow execution results, it was sufficient to map Sapporo’s provenance to the ontology provided by RO-Crate. However, for verification, we defined some additional terms and properties (<https://raw.githubusercontent.com/sapporo-wes/sapporo-service/main/sapporo/ro-terms.csv>). For example, we defined the property “mappedRate” to represent the mapping rate of the output file, which is a biological feature value used for verification. In addition, RO-Crate is designed to rely on the local file system for file location resolution and checksum representation. However, we prioritized the portability of being able to carry the

provenance in a single file, so we put all the information necessary for verification, such as checksums, biological feature values, and contents of files of small size, in the RO-Crate file.

Because the workflow output freely produces a large number of diverse files, it is impractical to extract biological feature values for all files. Thus, we used the file extension to determine the file type and used an appropriate tool to extract the biological feature values. We used the file types defined in the EDAM ontology [30] (RRID:SCR_006620), which are widely used to express biological interpretations (Table 3). For example, if the file type is SAM, we used SAMtools to extract the number of reads and the mapping rate, and if the file type is VCF, we used VCFtools [31] (RRID:SCR_001235) to extract the number of variants and the variant frequency. In addition, the number of lines in the file is also a biological feature value. For example, if the file type is FASTQ, 4 lines represent 1 sequence read. Fig. 1 is an example of a file entry in RO-Crate, which contains the biological feature values, such as the statistics obtained from the file, file size, and the number of lines.

For provenance enrichment and sharing, we integrated Sapporo and Yevis [32] (RRID:SRC_023204). Yevis is a system that builds a workflow registry and also acts as a client for Sapporo. The workflow metadata file used in Yevis contains not only the information required to execute the workflow but also the information for workflow availability and traceability in workflow sharing, such as author, open-source license, and documentation link. Thus, by executing the workflow in Sapporo via Yevis, the availability and traceability of the generated provenance are improved. In addition, because Yevis’s workflow sharing feature enables the attachment of generated provenance to shared workflows, the reproducibility of shared workflows can be verified by other users.

From the above, by executing the workflow with Sapporo and Yevis (Fig. 2), a workflow provenance containing feature values representing a biological interpretation is generated as RO-Crate. This method also applies to workflows written in various languages and can address a wide range of use cases, such as workflow sharing. Therefore, by generating and sharing a provenance containing biological feature values, it is possible to verify the reproducibility of the workflow execution results in other users’ environments.

Automatic verification of reproducibility

We developed Tonkaz to automatically verify the reproducibility of workflow execution results by comparing the biological feature values contained in the workflow provenance. One use case of Tonkaz is to compare the expected result, which is provided by a workflow developer, and the actual result, which is generated in the user’s environment (Fig. 2). That is, Tonkaz verifies that the results are the same according to the ACM’s definition of reproducibility (Table 1). Another use case is ACM’s definition of repeatability, which is to verify that the results are the same even if a workflow is executed multiple times in the same environment, and it will not be broken by updates to dependencies. Thus, these use cases indicate that we must verify the reproducibility of the results, regardless of the differences in execution methods and environments.

We designed Tonkaz to accept as arguments 2 RO-Crates, one containing the expected provenance and the other containing the actual provenance. Tonkaz compares the biological feature values of the output files in the 2 RO-Crates and calculates the reproducibility scale for each file. Among the various types of output files, including analysis results, summary reports, or execution

Table 2: Reproducibility scale of workflow execution results. For each of the output files, this is determined by comparing the expected provenance with the provenance of the actual execution. If the file of the same name in the expected provenance and the actual execution are identical, the file is considered fully reproduced. If the file of the same name is not identical, it is determined whether its difference is acceptable or not using the feature and threshold values. If the difference is acceptable, the file is considered partially reproduced. If the file exists in the expected provenance but not in the actual execution, it is considered to be not reproduced

Reproducibility scale	Level	Description
Fully Reproduced	3	Output files are identical with the same checksum.
Acceptable Differences	2	Output files are not identical, but their biological feature values (e.g., number of reads, mapping rate, and variant frequency) are similar (within a threshold).
Unacceptable Differences	1	Output files are not identical, and their biological feature values are not similar (beyond threshold).
Not Reproduced	0	The workflow does not produce output files.

Table 3: File types and extensions defined in EDAM ontology. These file types and extensions are used to extract biological feature values from the output files.

EDAM ID	File type	Extension
format_1929	FASTA	.fa, .fasta
format_1930	FASTQ	.fq, .fastq, .fq.gz, .fastq.gz
format_1975	GFF	.gff, .gff3
format_2306	GTF	.gtf
format_2572	BAM	.bam
format_2573	SAM	.sam
format_3003	Bed	.bed
format_3004	BigBed	.bb
format_3005	Wig	.wig
format_3006	BigWig	.bw
format_3016	VCF	.vcf, .vcf.gz

```
{
  "@context": ["https://w3id.org/ro/crate/1.1/context", { ... }],
  "@graph": [
    ...
    {
      "@id": "outputs/star_salmon/RAP1_UNINDUCED_REP2.markdup.sorted.bam",
      "@type": ["File", "FormalParameter", "OutputFile"],
      "contentSize": 3279083,
      "dateModified": "2022-09-08T08:52:19.755363",
      "encodingFormat": "application/gzip",
      "format": {
        "@id": "http://edamontology.org/format_2572"
      },
      "gid": 1000,
      "mode": "-rw-r--r--",
      "sha512": "2d6c8436dd1da0e4e49f9bdfbf8d656d7740f7eae149bb2add417d6739c05aeb441aa84239041f9aac87688042c7b31bfca5c95d0dfaf742512f2e740a788979",
      "stats": {
        "@id": "#31d3ba80-21df-4ab7-93e7-558154d07161"
      },
      "uid": 1000,
      "url": "http://localhost:1122/runs/93f4d8bf-424d-4d5e-bc79-9482e1620be9/data/outputs/star_salmon/RAP1_UNINDUCED_REP2.markdup.sorted.bam"
    },
    {
      "@id": "#31d3ba80-21df-4ab7-93e7-558154d07161",
      "@type": ["FileStats"],
      "duplicateRate": 0.8027008887713803,
      "duplicateReads": 78936,
      "generatedBy": {
        "@id": "#samtools"
      },
      "mappedRate": 1.0,
      "mappedReads": 98338,
      "totalReads": 98338,
      "unmappedRate": 0.0,
      "unmappedReads": 0
    },
    ...
  ]
}
```

Figure 1: An example of a file entity in RO-Crate. This is a part of the actual workflow execution results and uses the RNA-seq workflow distributed by the nf-core project. The output file in this example is a BAM file, and its biological feature values are the file size, number of mapped reads, and mapping rate. Thus, the file entity contains properties such as “contentSize,” “stats:mappedReads,” and “stats:mappedRate.” These are defined as additional terms in Sapporo, and the values are extracted from the file using SAMtools.

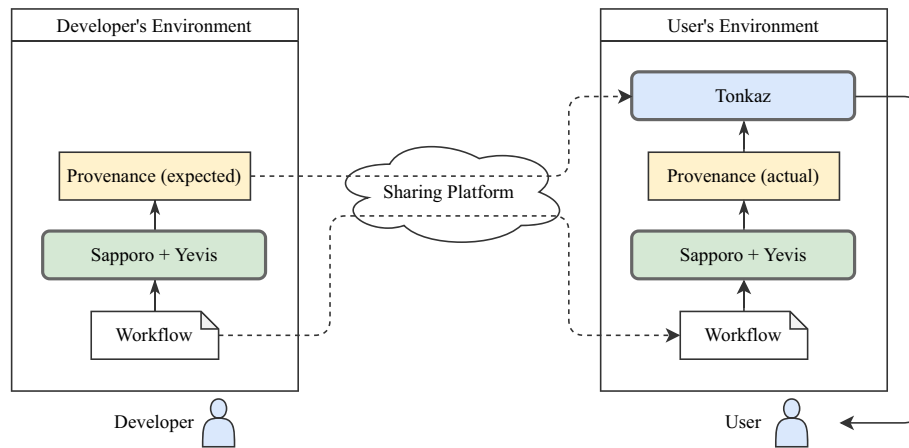


Figure 2: The flowchart representing the Tonkaz use case. The workflow built by the workflow developer is executed by WES, which is a combination of Sapporo and Yevis, and the workflow provenance, including feature values of the output files, is generated in RO-Crate format. This provenance is used as the expected value for the verification of reproducibility. Using the shared workflow, the user executes the shared workflow in their own environment using WES. Using Tonkaz, the user then compares the shared provenance with the provenance generated by the user's workflow execution and verifies the reproducibility of the workflow execution results.

logs, the system needs to select the files to compare. We aimed to compare the analysis results that led to a biological interpretation and to avoid the comparison of the output files that are not in a standard format. Therefore, we selected the file types to be compared as an initial set and selected the corresponding EDAM ontology terms listed in Table 3. With this selection, for example, the nf-core RNA-seq workflow produces 872 files, but only 25 files are assigned to the EDAM ontology and compared. In the process of comparing files and calculating the reproducibility scale, Tonkaz first checks whether the files are identical using a checksum (Fig. 3). If the files are identical, the reproducibility scale is “Fully Reproduced.” If the files are not identical, Tonkaz compares the biological feature values of the files using a threshold value to determine whether the differences are acceptable or not. The default threshold value used is 0.05, but this value can be changed according to the use case. This is because some workflows for medical applications or quality control of biological materials require a lower threshold value. The comparison result is finally summarized in a table, and the reproducibility scale of each file is also displayed in a table (Fig. 4). However, Tonkaz does not score the reproducibility of the entire workflow. This is because, again, the purpose of comparison may differ depending on the use case, and it is not practical to automate the final decision. Thus, we implemented an option to generate structured data in Tonkaz. In addition, we believe that workflow developers should use this option and write conditions or scripts to determine the reproducibility for each use case.

Results

To demonstrate the effectiveness of our approach, we verified the reproducibility of workflow execution results by comparing the results of public workflows used by researchers in real research projects, not simple ones for testing. We selected these workflows as they are fairly representative and mature current best practice for sequencing pipelines, implemented in different but typical workflow systems, and have similar set of genomics features that we can assess for provenance comparison. This verification was based on the following 5 practical use cases: (i) execution in the same environment, (ii) execution in a different environment, (iii) execution of different versions of the workflow, (iv) execution with

missing datasets, and (v) comparison using all output files. We used the following 3 workflows: (i) the mitochondrial short variant discovery workflow that distributed a GATK best practice workflow (hereafter referred to as GATK workflow, language: WDL) [33, 34], (ii) RNA-seq workflow distributed by nf-core (hereafter referred to as RNA-seq workflow, language: Nextflow) [24, 28], and (iii) GATK best practice-compatible germline short variant discovery workflow, which is used to process whole-genome sequencing data of the Japanese Genotype-Phenotype Archive (hereafter referred to as JGA workflow, language: CWL) [35–37]. We used the following 2 execution environments: (i) Ubuntu 20.04 LTS (CPU: Intel Xeon E5-2640 @ 2.50 GHz, RAM: 24 GB, Docker: 20.10.8) and (ii) macOS 12.5.1 (CPU: Apple M1 Max, RAM: 64 GB, Docker: 20.10.16). Table 4 shows the setting for each execution as a combination, and Table 5 summarizes the verification results based on the use cases. The methods and results of the workflow execution and verification are described in the online documentation “sapporo-wes/tonkaz - tests/README.md” [38] and are published on Zenodo [39].

The comparisons C1, C3, and C5 present execution results in the same environment. Comparison C1 was performed using the GATK workflow, and the output file types were BAM and VCF. The reproducibility scale value was level 2 (Acceptable Difference) for all files, with no differences in biological feature values expressing biological interpretation (e.g., mapping rate and variant frequency). The difference between these files was due to the fact that both the BAM and VCF files included the file paths of the original input file and timestamps in the header lines. Thus, when using the analysis tool GATK [40], it is challenging to fully reproduce the output files because of the behavior that the output files contain the file paths and timestamps. Comparison C3 was performed using the JGA workflow, and the output file types were VCF. This result also showed no differences in biological feature values, and the differences in file contents were due to the behavior of GATK. Comparison C5 was performed using the RNA-seq workflow, and the output file types were BAM, GTF, and BED. All GTF and BED files were level 3 (Fully Reproduced), and all BAM files were level 2 (Acceptable Difference). The difference between BAM files was due to the different order of the mapped reads in the BAM file. These BAM files were mapped by STAR [41] (RRID:SCR_004463) and sorted by SAMtools; however, differences occurred. These results show

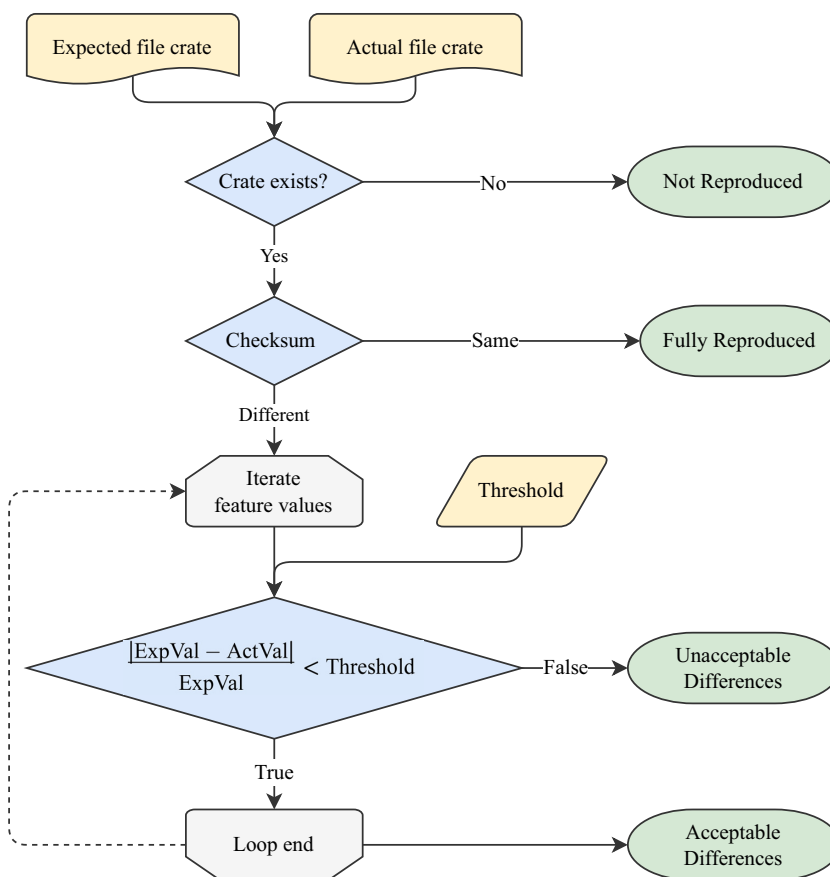


Figure 3: The process for calculating the reproducibility scale of a file. Tonkaz first checks whether the files are identical using a checksum. If the files are identical, the reproducibility scale value is “Fully Reproduced.” If the files are not identical, Tonkaz compares the biological feature values of the files using a threshold value to determine whether the differences are acceptable or not. If the differences are acceptable, the reproducibility scale value is “Acceptable Difference.” If the differences are unacceptable, the reproducibility scale value is “Unacceptable Difference.” The default threshold value used is 0.05, but this value can be changed according to the use case. If the file entity exists only in 1 of the 2 RO-Crates, the reproducibility scale value is “Not Reproduced.”

cases in which the output files were not identical, although the biological feature values were equal due to the behavior of the analysis tool.

The comparisons C2, C4, and C6 present the execution results in different environments. All of these comparison results were level 0 (Not Reproduced) because all execution in the Mac environment either failed or never finished. All workflows used in this experiment were Docker containerized and designed to be very reproducible in the execution context; however, runtime errors occurred due to the Arm processor architecture of the Mac environment. Thus, even a very well-considered workflow may not be reproducible in a different environment. In such cases, it is essential to increase the debuggability of the cause of the irreproducibility of the execution results. Therefore, the importance of this debuggability indicates that it is helpful to include information about the execution environment in the workflow provenance; our approach and RO-Crate address them.

Comparison C7 presents the execution results in different versions. Workflow developers often check for workflow breakage when updating versions of analysis tools included in the workflow. In the RNA-seq workflow used in this comparison, the dependent analysis tools STAR, SAMtools, and StringTie [42] (RRID: SCR_016323) were updated with the workflow update from v3.6 to v3.7. As a result of the comparison (C7), the number of files with level 2 increased compared to C5, a comparison involving

the same version. The file types that became level 2 were GTF, BED, and BAM; the GTF and BED files were newly changed from level 3 to level 2 when compared to C5. The differences between the GTF files were due to differences in the FPKM field values and the version of StringTie included in the header line. The BED files had a different number of lines, and the BAM files had a different number of mapped reads; however, those differences were within the threshold value. This result indicates that verification using biological feature values and threshold is effective because apparent differences occur in output due to version updates and other reasons, and it is necessary to determine whether these differences are acceptable or not.

Comparison C8 presents the execution results in a case where the input dataset was partially missing. The dataset used in RNA-seq_1st contains 6 sequence read files (FASTQ) [43], while the dataset used in RNA-seq_small contained 4 sequence read files [44]. As a result of the comparison, the output files related to the sample with half the number of reads were level 1 (Unacceptable Difference), while the sample with zero reads was level 0. In this case, setting the threshold used for verification to, for example, 0.5 instead of 0.05 (default value) will verify that the workflow is functioning as expected. That is, this suggests that the threshold value and final decision may vary depending on the objectives of developers and users.

```

...

=== Level3 ★★★★★ (Same Checksum, 13/25 files)

- star_salmon/rseqc/junction_annotation/bed/RAP1_IAA_30M_REP1.junction.Interact.bed
- star_salmon/rseqc/junction_annotation/bed/RAP1_IAA_30M_REP1.junction.bed
- star_salmon/rseqc/junction_annotation/bed/RAP1_UNINDUCED_REP1.junction.Interact.bed
- star_salmon/rseqc/junction_annotation/bed/RAP1_UNINDUCED_REP1.junction.bed
- star_salmon/rseqc/junction_annotation/bed/RAP1_UNINDUCED_REP2.junction.Interact.bed
- star_salmon/rseqc/junction_annotation/bed/RAP1_UNINDUCED_REP2.junction.bed
- star_salmon/rseqc/junction_annotation/bed/WT_REP2.junction.Interact.bed
- star_salmon/rseqc/junction_annotation/bed/WT_REP2.junction.bed
- star_salmon/stringtie/RAP1_IAA_30M_REP1.coverage.gtf
- star_salmon/stringtie/RAP1_UNINDUCED_REP1.coverage.gtf
- star_salmon/stringtie/RAP1_UNINDUCED_REP2.coverage.gtf
- star_salmon/stringtie/WT_REP1.coverage.gtf
- star_salmon/stringtie/WT_REP2.coverage.gtf

=== Level2 ★★★★★ (Similar Features, 12/25 files)

- star_salmon/WT_REP1.markdup.sorted.bam
  -----
  |                               | in Crate1 | in Crate2 |
  |-----|-----|-----|
  | File Size | 6.86 MB (7196616) | 6.86 MB (7192272) |
  | Total Reads | 188243 | 188241 |
  | # Mapped | 188243 (100.00%) | 188241 (100.00%) |
  | # Duplicate | 38470 (20.44%) | 38470 (20.44%) |
  |-----|-----|-----|

- star_salmon/rseqc/junction_annotation/bed/WT_REP1.junction.Interact.bed
  -----
  |                               | in Crate1 | in Crate2 |
  |-----|-----|-----|
  | File Size | 22.46 KB (22995) | 22.32 KB (22856) |
  | Line Count | 162 | 161 |
  |-----|-----|-----|

- star_salmon/stringtie/WT_REP2.transcripts.gtf
  -----
  |                               | in Crate1 | in Crate2 |
  |-----|-----|-----|
  | File Size | 37.82 KB (38729) | 37.82 KB (38729) |
  | Line Count | 259 | 259 |
  |-----|-----|-----|

...

=== Level1 ★ (Different Features, 0/25 files)

=== Level0 (Not Found, Crate1: 0 files, Crate2: 0 files)

Summarize compare result:

  -----
  | Reproducibility | Level | Definition | File # |
  |-----|-----|-----|-----|
  | Fully Reproduced | ★★★★★ | Same Checksum | 13 files |
  | Acceptable Differences | ★★★★★ | Similar Features | 12 files |
  | Unacceptable Differences | ★★ | Different Features | 0 files |
  | Not Reproduced | | Not Found | 0 files |
  |-----|-----|-----|-----|

```

Figure 4: Example of the Tonkaz output. Tonkaz displays a table for each file and a final summary table. The user checks those summary tables to determine the reproducibility of the entire workflow and the differences between the expected and actual files (e.g., by using the diff command).

Table 4: Combination table of workflow execution and execution settings. The first column is the definition of the execution name. In the second column and below are the workflow execution settings. The blank cells in the third and fourth columns indicate that there are no differences in the execution settings.

Execution name	Workflow	Version	Dataset	Environment
GATK_1st	GATK			Linux
GATK_2nd				Linux
GATK_mac				Mac
JGA_1st	JGA			Linux
JGA_2nd				Linux
JGA_mac				Mac
RNA-seq_1st	RNA-seq	v3.7	Standard	Linux
RNA-seq_2nd		v3.7	Standard	Linux
RNA-seq_mac		v3.7	Standard	Mac
RNA-seq_v3.6		v3.6	Standard	Linux
RNA-seq_small		v3.7	Small	Linux

Table 5: Comparisons of execution and verification results. The definition of each execution is defined in Table 4. Five use cases are assigned according to the combination of executions. In the fifth column and below are the number of files for each reproducibility scale defined in Table 2: level 3 is “Fully Reproduced,” level 2 is “Acceptable Difference,” level 1 is “Unacceptable Difference,” and level 0 is “Not Reproduced.”

ID	Source execution	Target execution	Use case	Level 3	Level 2	Level 1	Level 0
C1	GATK_1st	GATK_2nd	Same environment	0	5	0	0
C2	GATK_1st	GATK_mac	Different environment	0	0	0	5
C3	JGA_1st	JGA_2nd	Same environment	0	4	0	0
C4	JGA_1st	JGA_mac	Different environment	0	0	0	4
C5	RNA-seq_1st	RNA-seq_2nd	Same environment	20	5	0	0
C6	RNA-seq_1st	RNA-seq_mac	Different environment	0	0	0	25
C7	RNA-seq_1st	RNA-seq_v3.6	Different version	13	12	0	0
C8	RNA-seq_1st	RNA-seq_small	Missing dataset	8	5	7	5
C9	RNA-seq_1st	RNA-seq_2nd	All output files	557	306	1	8

Comparison C9 presents the execution results in a case where all the output files were compared. Most of the files were level 3 or level 2; however, 16 files were not reproduced (level 0). These level 0 files had random names or timestamps in the file names, for example, `mqc_mqc_mplplot_gtnuqiebf1_1.pdf` and `execution_report_2022-09-08_06-28-19.html`. Therefore, it is not appropriate to use all files to verify the reproducibility of execution results; it is essential to focus on characteristic files, such as BAM and GTF files.

For the 5 practical use cases, we found that our approach was well suited to verify the reproducibility of the workflow execution results. In all use cases, existing methods that use checksums to verify exact file matches can produce false positives; this means that the workflow is considered not reproduced, even though it is working as expected. Therefore, it is important to introduce a reproducibility scale and verify the workflow execution results' reproducibility at higher resolutions.

Discussion

Despite its complexity, data analysis in bioinformatics is considered reproducible and is being shared. In particular, the workflows shared by nf-core and GATK best practices are well maintained and include test datasets, documentation, and open-source licenses. Ideally, all shared workflows would be like these; however, in reality, this is challenging because of the amount of work and domain knowledge required. Thus, we aim to facilitate workflow sharing by providing a workflow provenance model and a workflow provenance verification method. However, our approach is

not applicable in domains where it is difficult to verify the results and inferences using a computer. In such cases, it is first necessary to discuss an ontology or structured format for representing the research.

A related project, CODECHECK [45], aims to provide the verification of the reproducibility of data analysis by a third party in scientific publishing. CODECHECK proposes a procedure similar to a peer review system, in which the workflow associated with research articles is verified at the time of publication by a reviewer called a CODECHECKER. However, this project focuses on increasing the availability of the workflow and does not verify the execution results. As such, it is unlikely to address the case of our concern that the execution results are not exactly the same, but the conclusions of the study remain the same. Our proposed metrics, a reproducibility scale of workflow execution results, would be useful in such workflow reproducibility validation in publishing as well.

In a scientific context, automated verification is a crucial process that should be performed for various reasons. Workflow developers can use it to easily add or update code and improve development efficiency. Administrators of workflow registries can use it to perform quality control, such as checking for broken links between the analysis tools and data used internally. Users of the workflow can also use it to validate the behavior of the workflow as an acceptance test in their own environment, thereby improving the reliability of their research projects. Tonkaz aims to support these validation efforts in different use cases and promote open science.

When generating workflow provenance using a format, such as RO-Crate, it is important to consider licensing issues. The provenance includes not only the execution results but also the executed workflow, input datasets, and software used internally. These files and software may have different licenses, and combining them under a single license can cause relicensing problems. In RO-Crate, a license can be specified for each entity; however, this approach is not currently possible as Sapporo automatically generates provenance from run requests and execution results without the original license information. This limitation can be overcome if data and software are consistently able to present their licenses, but this would require a generic method to get the license information of files retrieved from the internet.

Software begins to degrade from the moment it is developed, and it is not easy to maintain the same quality over time. Cases in which an error, including a stack trace, is thrown are quite fortunate; in many cases, the software cannot be executed in the first place, the process does not finish, or the output is inaccurate without throwing an error. Dealing with such cases and improving debuggability is accomplished by packaging the expected behavior of the software at the time it is developed. In our approach, we were able to attach information, such as OS, CPU architecture, and dependent software versions, to the expected workflow provenance due to RO-Crate's extensibility. However, when analysis tools are used internally, as in a workflow, the behavior of the analysis tools tends to be a black box. Therefore, if an option to display the reproducibility of the execution for each analysis tool is provided, it will be possible to identify the cause of the irreproducibility of the workflow execution results.

In the Results section, we showed the cases where the differences are found in the outputs but the biological interpretation will be identical. However, there are cases where users find differences that affect the interpretation even when comparing the same workflow definitions. For example, the output results may change when the workflow has a tool that dynamically uses external databases, which may be regularly updated over time. Another case when the impact on the results can be observed is a comparison of runs of the workflow that does not explicitly specify the software version, nor is it properly packaged.

As the system allows users to change the reporting threshold in comparison to the outputs, users need to be aware of the acceptable differences in the outputs of the given workflow. Although the threshold needs to be low for workflows used in applications that require severe quality control, such as medical data analysis, users can set it higher for workflows that can generate different outputs per run. For example, workflows using external databases or used for environmental monitoring purposes may have outputs that vary per run. The system alerts when a change was found, but as the interpretation depends on the cases, users need to understand the reason from the workflow description.

Though Tonkaz aims to improve the reproducibility of data analysis, the system itself also has a challenge in the reproducibility of its function. The system uses the file extension to check the file type, then specifies an external tool to extract the biological features from the file to compare the workflow outputs. However, the extracted features may change by the updates in the external tools, which results in the inconsistency of the results of comparison by Tonkaz. Another issue we see in the reproducibility of the comparison is the system's dependency on Sapporo, our WES implementation. Ideally, the results, analysis summaries, logs, and so on generated by analysis tools should be in a standardized format so any system can generate comparable statistics. The bioinformatics community needs to have a consen-

sus for such outputs of data analysis. As a related project, MultiQC attempts to summarize the results of multiple analysis tools [46] (RRID:SCR_014982). The Tonkaz system may improve its future consistency by integrating with a community effort like MultiQC, which can share the effort to extract the information from the analysis tools.

We used RO-Crate to express the provenance of our study and added additional terms and properties to the "@context" declaration for verification purposes. These terms are currently located on our own GitHub repository, but we are discussing with the RO-Crate community moving them to a more authoritative location, such as <https://github.com/ResearchObject/ro-terms>. In future work, we are also considering using the Workflow Run RO-Crate profile [47], which is currently under development to capture the provenance of executing a computational workflow, instead of our custom terms.

The proposed method is currently dependent on our software implementation; however, it can be generalized by the following 3 steps: (i) extract the statistics of biological features from the output, (ii) represent the statistics in a standardized format, and (iii) compare the statistics and report on the reproducibility scale. Although we implemented steps (i) and (ii) in Sapporo, it is ideal to let workflow execution systems have those 2 steps rather than a WES implementation. Once steps (i) and (ii) become common, step (iii) can be implemented in many kinds of data analysis platforms, while Tonkaz only provides a CLI interface. However, the bioinformatics community needs to have a consensus on the standardized scale for reproducibility.

In the article regarding a "reproducibility crisis," Baker quoted a Johns Hopkins microbiologist as stating, "The next step may be identifying what is the problem and to get a consensus" ([9], p. 452). Subsequently, the proliferation of virtualization technology and workflow systems has lowered the bar for reexecuting data analysis that an individual or others have previously built. Despite this, workflow developers are always anxious about whether their workflows are broken. In response to this anxiety, we realized that the cause is our binary view of whether the workflow could be reproduced or not. To remove this anxiety, we proposed a new approach to verify the reproducibility of workflows by providing a range of reproducibility of execution results. With the development of sharing platforms, workflow sharing has become more active. Therefore, we hope that by verifying reproducibility and sharing the results, more workflows will be reused with confidence, which, in turn, will lead to increased scientific progress.

Availability of Source Code and Requirements

- Project name: Tonkaz
 - Project homepage: <https://github.com/sapporo-wes/tonkaz>
 - DOI: 10.5281/zenodo.7559433
 - biotoolsID: tonkaz
 - RRID: SCR_023206
 - Operating system(s): Platform independent
 - Programming language: TypeScript
 - Other requirements: Deno
 - License: Apache License, Version 2.0
-
- Project name: Sapporo-service
 - Project homepage: <https://github.com/sapporo-wes/sapporo-service>
 - DOI: 10.5281/zenodo.7559425
 - biotoolsID: sapporo-service

- RRID: SCR_023202
 - Operating system(s): Platform independent
 - Programming language: Python
 - Other requirements: Docker recommended
 - License: Apache License, Version 2.0
- Project name: Yevis-cli
 - Project homepage: <https://github.com/sapporo-wes/yevis-cli>
 - DOI: 10.5281/zenodo.7546102
 - biotoolsID: yevis-cli
 - RRID: SCR_023204
 - Operating system(s): Platform independent
 - Programming language: Rust
 - Other requirements: Docker recommended
 - License: Apache License, Version 2.0

Authors' Contributions

H.S. and T.O. conceived and developed the methodology and software and conducted the investigation. H.S., T.F., and T.O. wrote the manuscript. T.F., T.O., and T.I. supervised the project. All authors read and approved the final version of the manuscript.

Abbreviations

ACM: Association for Computing Machinery; CPU: central processing unit; CWL: Common Workflow Language; GA4GH: Global Alliance for Genomics and Health; ID: identifier; OS: operating system; RNA-seq: RNA sequencing; RO-Crate: Research Object Crate; SAM: Sequence Alignment/Map; VCF: Variant Call Format; WDL: Workflow Description Language; WES: workflow execution service.

Competing Interests

The authors declare that they have no competing interests.

Data Availability

The workflow, result, and documentation related to the experiment described in the “Results” section are available on GitHub and Zenodo as follows:

- Execution method and result [38]
- Workflow definitions [48]
- Raw data of workflow execution results [39]

Funding

This study was supported by JSPS KAKENHI (grant 20J22439 to H.S.), the CREST program of the Japan Science and Technology Agency (grant JPMJCR17A1 to T.I.), and the Life Science Database Integration Project, NBDC of Japan Science and Technology Agency.

Acknowledgments

We acknowledge and thank the following scientific communities and their collaborative events where several of the authors engaged in irreplaceable discussions and development throughout the project: the Pitagora Meetup, Workflow Meetup Japan, NBDC/DBCLS BioHackathon Series, and RO-Crate community. We also acknowledge Prof. Masahiro Kasahara for his valuable com-

ments on the project. Computations were partially performed on the NIG supercomputer at ROIS National Institute of Genetics.

References

1. Stephens, ZD, Lee, SY, Faghri, F, et al. Big data: Astronomical or genomics? *PLoS Biol* 2015;**13**(7):e1002195.
2. Stein, LD. The case for cloud computing in genome informatics. *Genome Biol* 2010;**11**(5):207.
3. Goodwin, S, McPherson, JD, McCombie, RW. Coming of age: ten years of next-generation sequencing technologies. *Nat Rev Genet* 2016;**17**(6):333–51.
4. Prins, P, de Ligt, J, Tarasov, A, et al. Toward effective software solutions for big biology. *Nat Biotechnol* 2015;**33**(7):686–7.
5. Perkel, JM. Workflow systems turn raw data into scientific knowledge. *Nature* 2019;**573**(7772):149–51.
6. Bánáti, A, Kacsuk, P, Kozlovszky, M. Evaluating the reproducibility cost of the scientific workflows. In: 2016 IEEE 11th International Symposium on Applied Computational Intelligence and Informatics (SACI). 2016. p. 187–90.
7. Software with impact. *Nat Methods* 2014;**11**(3):211.
8. Rebooting review. *Nat Biotechnol* 2015;**33**(4):319.
9. Baker, M. 1,500 scientists lift the lid on reproducibility. *Nature* 2016;**533**(7604):452–4.
10. Drummond, C. Replicability is not reproducibility: nor is it good science. In: *Proceedings of the Evaluation Methods for Machine Learning Workshop at the 26th ICML*. Vol. 1. 2009.
11. Association for Computing Machinery. Artifact review and badging version 1.1. 2020. <https://www.acm.org/publications/policies/artifact-review-and-badging-current>. (accessed: 22 April 2023).
12. Leprevost, FdV, Barbosa, VC, Francisco, EL, et al. On best practices in the development of bioinformatics software. *Front Genet* 2014;**5**:199.
13. Wratten, L, Wilm, A, Göke, J. Reproducible, scalable, and shareable analysis pipelines with bioinformatics workflow managers. *Nat Methods* 2021;**18**(10):1161–8.
14. Amstutz, P, Mikheev, M, Crusoe, MR, et al. Existing workflow systems. 2021. <https://s.apache.org/existing-workflow-systems>. (accessed: 22 April 2023).
15. Crusoe, MR, Abeln, S, Iosup, A, et al. Methods included: standardizing computational reuse and portability with the common workflow language. *Communications of the ACM* 2022;**65**(6):54–63.
16. Voss, K, Gentry, J, Auwera, GVD. Full-stack genomics pipelining with GATK4 + WDL + Cromwell. *F1000Research* 2017;**6**:1381.
17. Di Tommaso, P, Chatzou, M, Floden, EW, et al. Nextflow enables reproducible computational workflows. *Nat Biotechnol* 2017;**35**(4):316–9.
18. Köster, J, Rahmann, S. Snakemake—a scalable bioinformatics workflow engine. *Bioinformatics* 2012;**28**(19):2520–2.
19. Leprevost, FdV, Grüning, BA, Alves Aflitos, S, et al. BioContainers: an open-source and community-driven framework for software standardization. *Bioinformatics* 2017;**33**(16):2580–2.
20. Soiland-Reyes, S, Sefton, P, Crosas, M, et al. Packaging research artefacts with RO-Crate. *Data Sci* 2022;**5**(2):97–138.
21. Khan, FZ, Soiland-Reyes, S, Sinnott, RO, et al. Sharing interoperable workflow provenance: a review of best practices and their practical application in CWLProv. *Gigascience* 2019;**8**(11):giz095.
22. Goble, C, Soiland-Reyes, S, Bacall, F, et al. Implementing FAIR digital objects in the EOSC-Life workflow collaboratory. *Zenodo*. <https://doi.org/10.5281/zenodo.4605654>, 2021.

23. O'Connor, BD, Yuen, D, Chung, V, et al. The Dockstore: enabling modular, community-focused sharing of Docker-based genomics tools and workflows. *F1000Research* 2017;**6**:52.
24. Ewels, PA, Peltzer, A, Fillinger, S, et al. The nf-core framework for community-curated bioinformatics pipelines. *Nat Biotechnol* 2020;**38**(3):276–8.
25. Ivie, P, Thain, D. Reproducibility in scientific computing. *ACM Computing Surveys* 2019;**51**(3):1–36.
26. Suetake, H, Tanjo, T, Ishii, M, et al. Sapporo: a workflow execution service that encourages the reuse of workflows in various languages in bioinformatics. *F1000Research* 2022;**11**:889.
27. Danecek, P, Bonfield, JK, Liddle, J, et al. Twelve years of SAMtools and BCFtools. *Gigascience* 2021;**10**(2):giab008.
28. Patel, H, Ewels, P, Peltzer, A, et al. nf-core/rnaseq: nf-core/rnaseq v3.7. *Zenodo*. <https://doi.org/10.5281/zenodo.6513815>, 2022.
29. Rehm, HL, Page, AJH, Smith, L, et al. GA4GH: International policies and standards for data sharing across genomic research and healthcare. *Cell Genom* 2021;**1**(2):100029.
30. Ison, J, Kalas, M, Jonassen, I, et al. EDAM: an ontology of bioinformatics operations, types of data and identifiers, topics and formats. *Bioinformatics* 2013;**29**(10):1325–32.
31. Danecek, P, Auton, A, Abecasis, G, et al. The variant call format and VCFtools. *Bioinformatics* 2011;**27**(15):2156–8.
32. Suetake, H, Fukusato, T, Igarashi, T, et al. Workflow sharing with automated metadata validation and test execution to improve the reusability of published workflows. *Gigascience* 2022;**12**:giad006.
33. DePristo, MA, Banks, E, Poplin, R, et al. A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nat Genet* 2011;**43**(5):491–8.
34. Peterson, A, Verdier, A, Abdel Ghany, A, et al. broadinstitute/gatk - scripts/mitochondria_m2_wdl. *GitHub*. 2021. https://github.com/broadinstitute/gatk/tree/33bda5e08b6a09b40a729ee525d2e3083e0ecdf8/scripts/mitochondria_m2_wdl.
35. Kodama, Y, Mashima, J, Kosuge, T, et al. The DDBJ Japanese Genotype-phenotype Archive for genetic and phenotypic human data. *Nucleic Acids Res* 2015;**43**(D1):D18–22.
36. National Bioscience Database Center. Whole genome sequencing analysis—NBDC Human Database. 2023. <https://humandbs.biosciencedbc.jp/en/whole-genome-sequencing>. (accessed: 22 April 2023).
37. Bioinformatics and DDBJ Center. ddbj/jga-analysis. *GitHub*. 2023. <https://github.com/ddbj/jga-analysis>.
38. Suetake, H. sapporo-wes/tonkaz - tests/README.md. *GitHub*. 2022. <https://github.com/sapporo-wes/tonkaz/blob/main/tests/README.md>.
39. Suetake, H. Raw data of workflow execution results used in Tonkaz's experiments. *Zenodo*. <https://doi.org/10.5281/zenodo.7660388> 2023.
40. McKenna, A, Hanna, M, Banks, E, et al. The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res* 2010;**20**(9):1297–303.
41. Dobin, A, Davis, CA, Schlesinger, F, et al. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* 2013;**29**(1):15–21.
42. Pertea, M, Pertea, GM, Antonescu, CM, et al. StringTie enables improved reconstruction of a transcriptome from RNA-seq reads. *Nat Biotechnol* 2015;**33**(3):290–5.
43. nf-core Community. nf-core/test-datasets - rnaseq/samplesheet/v3.4/samplesheet_test.csv. *GitHub*. 2018. https://raw.githubusercontent.com/nf-core/test-datasets/rnaseq/samplesheet/v3.4/samplesheet_test.csv.
44. Suetake, H, Ohta, T. sapporo-wes/test-workflow: 1.0.2 - assets/nf-core_rnaseq_samplesheet_small_test.csv. *Zenodo*. 2022. https://raw.githubusercontent.com/sapporo-wes/test-workflow/1.0.2/assets/nf-core_rnaseq_samplesheet_small_test.csv.
45. Nüst, D, Eglen, SJ. CODECHECK: an open science initiative for the independent execution of computations underlying research articles during peer review to improve reproducibility. *F1000Research* 2021;**10**:253.
46. Ewels, P, Magnusson, M, Lundin, S, et al. MultiQC: summarize analysis results for multiple tools and samples in a single report. *Bioinformatics* 2016;**32**(19):3047–8.
47. Workflow Run RO-Crate Working Group. Workflow Run RO-Crate. 2023. <https://www.researchobject.org/workflow-run-crate/>. (accessed: 22 April 2023).
48. Suetake, H, Ohta, T. sapporo-wes/test-workflow: 1.0.2. *Zenodo*. <https://doi.org/10.5281/zenodo.7102664>, 2022.