

1 **Simplifying the development of portable, scalable, and** 2 **reproducible workflows**

3 Stephen R. Piccolo^{1,*}, Zachary E. Ence¹, Elizabeth C. Anderson¹, Jeffrey T. Chang², Andrea H. Bild³

4 1 - Department of Biology, Brigham Young University, Provo, UT, 84602, USA

5 2 - Department of Integrative Biology and Pharmacology, University of Texas Health Science Center at
6 Houston, Houston, TX, 77030, USA

7 3 - Department of Medical Oncology and Therapeutics, City of Hope Comprehensive Cancer Institute,
8 Monrovia, CA, 91016, USA

9 * - Please address correspondence to S.R.P. at stephen_piccolo@byu.edu. Phone: 801-422-7116.

10 **Abstract**

11 Command-line software plays a critical role in biology research. However, processes for installing and
12 executing software differ widely. The Common Workflow Language (CWL) is a community standard that
13 addresses this problem. Using CWL, tool developers can formally describe a tool's inputs, outputs, and
14 other execution details in a manner that fosters use of shared computational methods and reproducibility
15 of complex analyses. CWL documents can include instructions for executing tools inside software
16 containers—isolated, operating-system environments. Accordingly, CWL tools are portable—they can be
17 executed on diverse computers—including personal workstations, high-performance clusters, or the
18 cloud. This portability enables easier adoption of bioinformatics pipelines. CWL supports workflows,
19 which describe dependencies among tools and using outputs from one tool as inputs to others. To date,
20 CWL has been used primarily for batch processing of large datasets, especially in genomics. But it can

21 also be used for analytical steps of a study. This article explains key concepts about CWL and software
22 containers and provides examples for using CWL in biology research. CWL documents are text-based, so
23 they can be created manually, without computer programming. However, ensuring that these documents
24 confirm to the CWL specification may prevent some users from adopting it. To address this gap, we
25 created ToolJig, a Web application that enables researchers to create CWL documents interactively.
26 ToolJig validates information provided by the user to ensure it is complete and valid. After creating a
27 CWL tool or workflow, the user can create “input-object” files, which store values for a particular
28 invocation of a tool or workflow. In addition, ToolJig provides examples of how to execute the tool or
29 workflow via a workflow engine.

30 Keywords: workflows, command-line software, computational reproducibility, Common Workflow
31 Language, Web application, learn by example

32 **Introduction**

33 Software is fundamental to modern scientific research^{1,2}. It can accelerate the pace of research, formalize
34 algorithmic logic, and support reproducibility³. In a 2014 survey, 92% of academic scientists reported
35 using software in their research⁴. Our article focuses on command-line software, which scientists use in
36 many disciplines⁵ and which often provides advantages over point-and-click software, especially for
37 computational reproducibility^{3,5}. Command-line software includes tools that scientists have prepared for
38 general use within the research community; but it also includes custom scripts that scientists create to
39 analyze particular dataset(s). Popular languages for creating such scripts include Python, R, Perl, and
40 bash. In this article, we describe processes for using command-line tools and scripts in research; we
41 emphasize challenges associated with using them, combining them, and ensuring that others can use
42 them; and we describe how *workflow engines* and *software containers* help to address these challenges.
43 Our target audience includes bioinformaticians who wish to develop and deploy software tools. But we

44 also demonstrate how workflow engines and software containers can be used to support reproducibility of
45 research analyses, including those that use custom scripts.

46 A *workflow* is a defined series of steps for processing or analyzing research data⁶. Each step applies one
47 or more computational tools to the data with specific inputs, outputs, and configuration settings. Defining
48 these steps as a workflow provides advantages over using command-line tools alone. Most importantly, a
49 workflow provides a vocabulary for describing how the tools are interconnected. It also provides more
50 flexibility regarding the computing environments that can be used. For example, the Snakemake⁷ and
51 Nextflow⁸ workflow engines facilitate execution of workflows on local-, cluster-, or cloud-based
52 computers. This process is made easier when command-line tools are encapsulated in software containers⁹
53 because the tools are portable.

54 Different workflow-management systems use different methodologies for defining workflows and for
55 interfacing with software used within those workflows⁶. This heterogeneity motivated creation of the
56 Common Workflow Language (CWL), an open specification for describing command-line tools and
57 workflows¹⁰. CWL documents are recognized by many workflow engines, including Snakemake⁷,
58 *cwltool*, Toil¹¹, Apache Airflow¹², Tibanna¹³, and Arvados¹⁴.

59 It is possible to use workflow engines like Snakemake and Toil independently of CWL. It is also possible
60 to use software containers on their own. However, by creating CWL documents, scientists can describe
61 tools and workflows in a way that is standards-based and agnostic to the workflow engine(s) on which it
62 is executed. CWL documents are merely text-based files and thus can be created manually via a text
63 editor. Despite this simplicity, many researchers have yet to adopt CWL. Useful tutorials are available
64 online¹⁵, but as CWL specification provides flexibility, researchers face a learning curve to ensure that
65 documents are specified correctly. Some software applications are available to aid in this process. For
66 example, *Rabix Composer* is a desktop application that enables researchers to create and edit CWL
67 documents¹⁶. This application supports many features in the CWL specification and provides both text-
68 based and visual editors. However, researchers unfamiliar with nuanced details of the CWL specification

69 may find Rabix Composer too advanced for their needs, and it does not support the latest versions of
70 CWL. Alternatively, CWL plugins exist for many code editors¹⁵. In addition, developers have created
71 software packages that enable researchers to build CWL documents via application programming
72 interfaces. However, many researchers who could benefit from CWL lack the programming expertise to
73 use these resources.

74 In this article, we illustrate how to use CWL to describe command-line tools and workflows and to
75 perform reproducible research analyses. We provide 27 examples of CWL documents for completing
76 diverse types of research tasks, ranging from simple (e.g., printing custom messages to the console) to
77 advanced (identifying differentially expressed genes or calling somatic mutations in cancer genomes). In
78 addition, we introduce *ToolJig*, a Web application that enables researchers to create CWL documents.
79 *ToolJig* provides a simple, interactive interface that requires no installation and includes prompts to guide
80 the user. Via *ToolJig*, a researcher can specify details about a tool’s expected inputs and outputs,
81 operating-system environment, and auxiliary files (e.g., scripts, configuration files). Researchers can also
82 create workflows that integrate these tools. *ToolJig* checks the information provided by the user to ensure
83 it is complete and valid. After successfully describing a tool or workflow, the researcher can download
84 CWL files and use *ToolJig* to create “input-object” files, which store input values for a particular
85 invocation of a tool or workflow. In addition, *ToolJig* provides examples of how to execute the tool or
86 workflow via a workflow engine.

87 **Using containers to manage software installation and configuration**

88 First, we address software installation and configuration, which are essential steps to creating CWL tools
89 and workflows. These seemingly simple steps are fraught with challenges. Although the software may be
90 downloadable from a public website, installation instructions are sometimes vague, the process may
91 involve many steps, and these steps may differ for each operating system. Such challenges led
92 computational biologist Ian Holmes to quip, “You can download our code from the URL supplied. Good
93 luck downloading the only postdoc who can get it to run”¹⁷. Package managers, such as *bioconda* and

94 *bioconductor*, have helped to ameliorate these challenges, providing mechanisms to install software
95 dependencies and track versions. These package managers function in a way that is mostly agnostic to the
96 user's operating system^{18,19}. However, some software tools are not available via package managers,
97 package managers depend on operating-system components that cannot be installed using the package
98 managers themselves, and package managers may not guarantee that older versions of software and their
99 dependencies remain available³. Software *containers* have gained popularity among scientists in recent
100 years because they help to overcome these limitations.

101 Software containers provide a mechanism to encapsulate specific versions of software and their
102 dependencies in a fully configured, operating-system environment⁹. Here we focus on the *Docker*
103 ecosystem, which is commonly used for building, managing, and sharing software containers^{9,20}. Other
104 containerization systems are also available²¹⁻²⁵; these systems are typically compatible with Docker. In
105 academic computing environments, such as university-run cluster computers, Singularity has gained
106 popularity because containers can be executed without administrative privileges.

107 Steps for configuring the operating-system environment and installing software within a container are
108 documented in a "Dockerfile". Using such a file, researchers can build a container *image*, a layered set of
109 operating-system components. Commonly, the base layer is a minimal implementation of a Linux
110 distribution (for example, Debian 10.3 or Ubuntu 18.04). Subsequent layers consist of software
111 dependencies, configuration files, environment variables, etc. Once a container image has been created, it
112 is portable and immutable. This is advantageous for computational reproducibility because one researcher
113 can share an image with another researcher and know that its contents have remained static. A software
114 *container* is an actively executing instance of a particular container image. Multiple containers of the
115 same image can be executing simultaneously on the same computer (or different computers). Docker
116 containers are always Linux-based; this is convenient for biology research because bioinformatics
117 software is predominantly designed for Linux. Even though a container is Linux-based, it can be executed

118 on non-Linux operating systems, such as Windows 10 or Mac OS 10, via a container *engine*. Container
119 engines use virtual machines to facilitate this interaction³.

120 In addition to packaging scientific software, container images can package analysis code. For example,
121 upon analyzing a given dataset, a researcher may wish to share the code with the research community.
122 Many researchers post analysis code on websites such as GitHub (<https://github.com>) or Open Science
123 Framework²⁶. This practice can enable others to verify and reuse the code; it also benefits the original
124 researcher who otherwise might lose track of analysis details³. But even when analysis code resides in the
125 public domain, third-party researchers may experience difficulty executing it. Scripting languages like
126 Python²⁷ and R²⁸ require interpreters. Analysis scripts may depend on specific versions of interpreters, but
127 the third-party researcher may have a different version on their computer. In addition, most analyses rely
128 on ancillary software packages. Such packages provide logic for parsing a certain type of file, performing
129 statistical tests, creating graphics, etc. Versioning is critical: older (or newer) versions of software
130 packages may be incompatible with the analysis code. Researchers can facilitate reproducibility by
131 providing container images that include specific versions of any script interpreters or software packages
132 that are necessary to execute an analysis.

133 *Binder* facilitates the containerization process for analysis code stored in GitHub repositories²⁹. To use
134 Binder, a researcher creates a configuration file that indicates which software is needed in the container
135 image. For Python and R analyses, these configuration files indicate packages that must be installed, as
136 well as version information. In other cases, a Dockerfile can be used to configure the environment more
137 flexibly. After the researcher places the configuration file (and analysis code) in a GitHub repository,
138 other researchers can visit the Binder website and re-execute the analysis. Behind the scenes, Binder
139 provisions a cloud-based computer and executes the code within a container. This solution is effective for
140 relatively short-running analyses that require modest computational resources, that use small datasets, and
141 that are ready to be released publicly. However, many analyses do not meet these criteria.

142 Longer, more data-intensive analyses can be executed via workflow-management systems.

143 A *workflow* is a defined series of steps for processing or analyzing research data⁶. Each step applies one
144 or more computational tools to the data with specific inputs, outputs, and configuration settings. Defining
145 these steps as a workflow provides advantages over using command-line scripts alone. Most importantly,
146 a workflow provides a vocabulary for describing how the tools are executed and interconnected. It also
147 provides more flexibility regarding the computing environments that can be used. For example, the
148 Snakemake⁷ and Nextflow⁸ workflow engines facilitate execution of workflows on local-, cluster-, or
149 cloud-based computers. This process is made easier when command-line tools are encapsulated in
150 container images because the tools are portable.

151 **Basic elements of a Common Workflow Language tool description**

152 To describe execution of a command-line tool, a researcher creates a text-based file according to the
153 CWL Command Line Tool Description specification³⁰. CWL files can be structured using either the
154 YAML or JSON data-serialization formats^{31,32}. Here we provide an overview of key components of CWL
155 tool descriptions.

156 A CWL tool describes inputs that will be used when the command-line tool is executed. A data type (or
157 schema) should be defined for each input, indicating whether it represents a string, number, Boolean
158 value, file, directory, or array (a data structure with multiple values). In practice, these inputs are
159 generally data files and configuration settings for the software. These definitions help users of the tool
160 understand the nature of each input and make it easier for inputs to be validated. For example, if a
161 command-line tool expects a particular input value to be an integer (e.g., number of threads), a workflow
162 engine can verify that the user has specified an integer before executing the tool.

163 After inputs have been defined, a tool description must provide instructions for constructing a command
164 from the inputs. These commands can be based on a sorted ordering of the inputs. Alternatively, the
165 researcher can specify a template for the command, using placeholders for the inputs. Such templates can
166 represent either a single command or multiple commands.

167 As a tool executes, it can generate three types of outputs that might be useful to a researcher: 1) standard
168 output, 2) standard error, and 3) new files. *Standard output* often consists of informational messages
169 printed to the console; but it may also contain data to be used as input for another tool. *Standard error*
170 typically consists of errors, warnings, or diagnostic information printed to the console. Many command-
171 line tools produce new data files that have resulted from execution of the tool. A CWL tool description
172 must indicate which of these outputs are expected so that a workflow engine can “collect” them after
173 executing the tool.

174 **CWL tool examples for printing simple command-line messages**

175 In the GitHub repository that accompanies this article, we have provided example tool descriptions,
176 formatted according to the CWL specification (<https://github.com/srp33/ToolJig/tree/master/examples>).
177 The first series of examples is stored in the `hello` subdirectory.

178 The first example, `01_hello.cwl`, requires two inputs: 1) a person’s given name and 2) the person’s
179 surname. It uses the `baseCommand` field to construct a command based on ordinal positions specified
180 for the inputs; the resulting command prints a greeting for that person. In this simple example, the only
181 output is the greeting sent to standard output, which is redirected to a file called `01_output.txt`. In
182 the GitHub repository, `hello_objects.yml` is an example input-object file for this tool. It species a
183 person’s given name (“Fernanda”) and surname (“Dantas”). The user could execute the tool via the *cwl*
184 workflow engine using this command: `cwltool 01_hello.cwl hello_objects.yml`. The
185 output would be, “Hello, Fernanda Dantas”.

186 Suppose we wished to alter the greeting to include an exclamation point and to indicate the person’s age.
187 First, we would add an input for the person’s age (an integer). Second, we would need to update the
188 command that will be executed. However, CWL’s `baseCommand` field provides limited flexibility for
189 constructing commands. Instead, our example provides a command template using the

190 ShellCommandRequirement and arguments fields and uses placeholders within the template for
191 each input. As shown in 02_hello.cwl (and Figure 1A), we use the following template:

```
192 echo Hello, $(inputs.given_name) $(inputs.surname)! You are  
193 $(inputs.age) years old.
```

194 We use a \$ character and parentheses to indicate placeholders for input variables. We prefix each input
195 variable with “inputs.” to indicate that they will be specified as inputs.

196 By default, the first two examples would be executed within the same operating-system environment as
197 the workflow engine. Accordingly, these tools can only be executed on operating systems that support the
198 echo command. Many commands are only supported on particular operating systems—or their behaviors
199 differ by operating system. So in 03_hello.cwl (and Figure 1B), we use a Docker image based on the
200 “buster” release (version 10.3) of the Debian Linux operating system. The DockerRequirement field
201 is added (in this case, two lines of text). Before executing the tool, a researcher would install a container
202 engine such as Docker Desktop. Then, when executing the tool, the workflow engine would integrate
203 itself with the container engine, which would identify any input files or directories and create container
204 volumes so that the files or directories could be accessed from within the container. With this option
205 enabled, the command-line tool becomes portable—it can be executed on any computer that supports the
206 container and workflow engines.

207 **CWL tool examples for performing simple data analyses**

208 The second series of examples is stored in the examples/bmi subdirectory on the GitHub site. The
209 01_bmi.cwl tool provides a simple example of a reproducible, quantitative analysis that could be
210 performed using CWL. It accepts as input a tab-separated file containing names, weights, and heights of
211 (fictional) individuals. A second input specifies the name of the column in the tab-separated file that
212 contains weight information. The third input specifies the column name containing height information.
213 The fourth input is the name of an output file that will be created. This example illustrates the use of an

214 *auxiliary file*. Under the `InitialWorkDirRequirement` field, the contents of a Python script are
215 stored. This script is used to calculate Body Mass Index (BMI) values for each person in the input file and
216 store those values in a *BMI* column in the output file. Below is the command template that we use.

```
217 python calculate_bmi.py "$(inputs.input_file.path) "  
218 "$(inputs.weight_column_name)" "$(inputs.height_column_name) "  
219 "$(inputs.output_file_name) "
```

220 The command template specifies the inputs as arguments to the Python script. When a file input is used,
221 the workflow engine stores metadata about the file in an object with multiple attributes. Thus to reference
222 the file's path within the command template, we append ".path" to the input name. As the workflow
223 engine executes the tool, it stores the auxiliary Python script inside the container, invokes the script, and
224 collects the output file that the script generates.

225 If a researcher wished to ensure that others could reproduce the BMI calculations, they would only need
226 to share the CWL file, the input-object file (`01_bmi_objects.yml`), and the data file
227 (`biometric_data.tsv`). However, many analyses use data stored in online repositories. In such
228 cases, it is convenient for a CWL tool to pull data directly from an online repository. The `02_bmi.cwl`
229 tool description and Figure 2 illustrate this approach. Similar to the previous example, it extracts names,
230 weights, and heights from a tab-separated file and adds a BMI column. However, it pulls the file from a
231 web server (in this case, our GitHub repository). The command template is similar to the previous
232 example. Again, we use a software container based on Debian Linux. But this time, we use the
233 `NetworkAccess` field to enable the container to connect to external computers. The tool emits
234 messages to both standard output and standard error; these messages are stored in files called
235 "02_output.txt" and "02_error.txt", respectively.

236 **CWL tool examples for processing transcriptomic data**

237 The third series of examples (`examples/transcriptomic` subdirectory) are wrappers around
238 existing tools for processing transcriptomic data. In both cases, we use R packages from the Bioconductor
239 suite¹⁹. Although R and Bioconductor are designed to be compatible with all major operating-systems,
240 some packages require dependencies that are operating-system specific. Furthermore, many Bioconductor
241 packages provide a large number of functions and options. Researchers can create CWL tool descriptions
242 that install dependencies (within a container image) and support a narrower range of options. The
243 researchers might then share this tool with other researchers, enabling them to apply the tool more easily
244 to their own data. Alternatively, they might use the tool as a way to support reproducibility for their own
245 analyses.

246 The Single-Channel Array Normalization (SCAN) algorithm normalizes data from gene-expression
247 microarrays, correcting for background noise and oligonucleotide-binding biases³³. The SCAN method is
248 implemented in the SCAN.UPC package in Bioconductor. It can download data directly from Gene
249 Expression Omnibus (GEO)³⁴. The `scan_normalize.cwl` example illustrates how this functionality
250 could be incorporated into a CWL tool. The base container image in this example includes the core
251 Bioconductor components; our Dockerfile extends this image by installing the SCAN.UPC package. In
252 addition, our example uses an auxiliary file containing R code that invokes the SCAN function within this
253 package to normalize a given GEO series. Upon executing, this tool produces a tab-separated output file
254 containing normalized measurements for all biological samples in the series. The tool could be
255 customized further, for example, to perform gene-level rather than probeset-level summarization or to
256 perform a quality-control analysis.

257 Commonly, researchers seek to identify genes that are differentially expressed between two conditions.
258 The DESeq2 package is popular for performing such analyses with RNA-Sequencing data³⁵. Our
259 `deseq2.cwl` example illustrates how this process could be incorporated into a CWL tool. Similar to the

260 previous example, this tool uses a container image with Bioconductor core components and installs the
261 “DESeq2” package. In addition, it installs the readr and dplyr packages^{36,37}, which we use to read and
262 parse the data before identifying differentially expressed genes. The first two inputs are URLs to data files
263 containing gene counts and phenotype information in tab-separated formats. The third input is a string
264 representing a design formula; users of the tool can customize the differential-expression calculations
265 based on the dependent variable of interest as well as any covariates. In the example input-object file
266 (`deseq2_objects.yml`), we use data from an RNA-Sequencing experiment that compared two
267 inbred mouse strains commonly used for neuroscience research³⁸; the data had been aligned to a reference
268 genome, and gene counts had been quantified previously.

269 **CWL workflow examples for performing mathematical calculations**

270 The examples so far have illustrated how to execute command-line tools in isolation. CWL workflows
271 must specify at least one input(s) and one output(s) for the entire workflow. In addition, the researcher
272 must define steps that each consist of a tool with input(s) and output(s). The researcher indicates whether
273 each step’s input should be populated by an input for the entire workflow or by the output of a previously
274 completed step. The workflow’s output(s) consist of the output(s) of one or more of the steps. As with
275 CWL tools, the researcher must create an input-object file that provides input values for a particular
276 execution of the workflow. Upon executing the workflow, the workflow engine evaluates the sequence of
277 steps that must be executed and connects inputs with outputs, as needed.

278 We provide three example workflows in the `examples/workflows/math` subdirectory of the
279 GitHub repository. The `add_sqrt_workflow.cwl` example accepts two integers as inputs, sums
280 them, calculates the square root of the sum, and then stores the square root of the sum in an output file.
281 This example illustrates the basic process of using an output from one tool as input to another. The
282 `recursive_sqrt_workflow.cwl` example reads a number from a file, calculates the square root of
283 that number, calculates the square root of the resulting number, and saves the output to a file. This

284 workflow demonstrates the ability to invoke the same tool recursively. The
285 `secondary_sqrt_workflow.cwl` example calculates the square root of a number stored in a file
286 and saves the result to an output file. It does the same for two secondary files. It then sums those values
287 and writes the sum to a file. This example demonstrates using secondary files within a workflow.
288 Secondary files are commonly used in genomics and simplify the process of working with groups of files
289 that are necessary to complete a particular task.

290 **CWL workflow example for identifying somatic variants in a cancer genome**

291 The examples in the `examples/workflows/somatic` subdirectory demonstrate a process for call
292 somatic variants from Illumina sequencing reads. We use paired-end reads from tumor and normal cells
293 for an individual from the Texas Cancer Research Biobank³⁹. (Although these data are publicly available,
294 they are subject to some data-use restrictions³⁹.) To shorten execution times, we use a subset of the data:
295 the first 10,000,000 reads from lane 2 of the sequencing run. Furthermore, our analysis is limited to
296 essential steps for preparing the data and calling variants. Additional steps like annotation and filtering
297 would improve sensitivity and specificity of the variant calls; accordingly, researchers should interpret
298 our variant calls with caution.

299 In these examples, somatic-variant calling occurs in a series of steps (Figure 3):

- 300 1. *Download and index a human reference genome (version hg38)*. We use the Linux `wget` and
301 `gunzip` utilities to download and decompress a FASTA file from the UCSC genome repository⁴⁰.
302 We also use `bwa`, `samtools`, and Picard Tools to index the FASTA file and create a sequence
303 dictionary⁴¹⁻⁴³.
- 304 2. *Preprocess reference files containing known polymorphic sites in preparation for base-quality*
305 *score recalibration (BQSR)*. We download Variant Call Format files⁴⁴ from the Genome Analysis
306 Toolkit (GATK) resource bundle⁴⁵ and use a custom Python script to adjust chromosome

- 307 identifiers that may differ across reference genomes. We also use Picard tools to sort the
308 reference files.
- 309 3. Download the raw sequencing reads (FASTQ files). The files are stored in a publicly available,
310 Open Science Framework repository²⁶.
 - 311 4. Trim adapter sequences and low-quality bases using atropos⁴⁶.
 - 312 5. Align the trimmed reads to the reference genome using *bwa mem*⁴¹. A read-group string is also
313 specified.
 - 314 6. Sort and index the resulting BAM files⁴² using sambamba⁴⁷.
 - 315 7. Mark duplicate reads and re-index the resulting BAM files using sambamba.
 - 316 8. Derive a BQSR table from each BAM file using GATK.
 - 317 9. Apply the BQSR table to the BAM files using GATK.
 - 318 10. Call somatic single-nucleotide variants and small insertions/deletions using Mutect2⁴⁸. This
319 produces a VCF file.
 - 320 11. Call somatic structural variants using Delly⁴⁹. This produces a VCF file.

321 This workflow executes many of the same steps for a normal DNA sample and a tumor DNA sample.
322 These steps are independent of each other. Accordingly, when multiple computing cores are available, the
323 workflow engine may execute these steps in parallel.

324 Most of the tools in this workflow use container images from the BioContainers project, which provides
325 thousands of Docker images that encapsulate biology-related software⁵⁰. In some cases, we used a
326 BioContainers image that had been built for a specific version of a software package in the Bioconda
327 project.¹⁸ In cases where we used multiple packages for a given task, we started with a base image from
328 Biocontainers and used Bioconda to install the packages. In the case of GATK, we used container images
329 that had been created by the Broad Institute and stored on Docker Hub.

330 In most cases, we followed the recommendation that a single container image use only a single software
331 package⁵¹. However, in some cases, we determined that it was more sensible to use multiple software

332 packages in a given CWL tool. For example, when preparing index files for the reference genome, we use
333 three separate software packages. In contrast, sometimes we used the same software in multiple CWL
334 tools. For example, the BQSR steps are computationally intensive; thus, we separated them into distinct
335 CWL tools so that computational resources can be allocated at a more granular level. In this sense, each
336 CWL tool represents a practical unit for data processing, not necessarily a particular software package.

337 **ToolJig**

338 **Using ToolJig to create CWL tool descriptions**

339 In manufacturing, a “jig” is used by toolmakers to ensure that products are created in a repeatable,
340 consistent pattern. Similarly, ToolJig provides a means to generate CWL tool descriptions, workflows,
341 and input-object files in a repeatable, consistent manner. ToolJig is a Web application that uses the Vue.js
342 framework⁵². Its functionality is divided into two pages: one for creating CWL tools and one for creating
343 workflows. It is available from <https://srp33.github.io/ToolJig/>. To create a tool description in ToolJig,
344 users specify the following:

- 345 1. A unique identifier. This identifier is used in the name of the CWL file that is generated, as well
346 as for tagging the Docker image.
- 347 2. A short label that describes the tool’s purpose and function.
- 348 3. Optionally, a longer description that provides more detailed documentation about the tool.
- 349 4. Dockerfile contents. These instructions indicate the base container image that should be used and
350 any additional commands necessary to build and configure a container image for the tool. The
351 tutorial by Nüst, et al. provides helpful recommendations on authoring Dockerfiles⁵³.
- 352 5. Author information. Optionally, the tool author can specify their name and ORCID identifier⁵⁴.
353 This information helps to ensure that authors are credited for their work.

- 354 6. Software license. The tool author can select from among 7 popular licenses, thus indicating
355 conditions under which others can use the CWL document. This license may or may not be
356 identical to the license specified for the software itself.
- 357 7. Inputs. Users specify a name, type, and description for each input. Supported types are integer,
358 string, File, and “Output File.” The File type allows the user to indicate that an input file is
359 expected and asks the user to specify the EDAM format of the file⁵⁵. Additionally, input files may
360 be associated with secondary files. For instance, as our examples illustrate for somatic variant-
361 calling, BAM files must be accompanied by index files. Rather than specify these as two separate
362 inputs, we indicate that an index file is secondary to a BAM file. “Output File” is a convenience
363 type that is used when a tool author wants users to be able to specify the *name* of an output file
364 that will be generated as the tool executes. Because this requires user input, ToolJig provides it as
365 an input option. When a tool author specifies this type, ToolJig creates a *string* input for the file
366 name, along with a corresponding output element, thus simplifying this process for the user.
- 367 8. Auxiliary files. The tool author can enter the name and contents of any auxiliary files that will be
368 used. These tools are commonly programming scripts and can be written in any programming
369 language that the Docker image supports.
- 370 9. Command template. The tool author specifies a template for executing the tool at the command
371 line. Each input must be specified at least once in this template; ToolJig provides syntax
372 suggestions to the user. ToolJig uses these command templates as an alternative to the
373 `baseCommand` field. As our examples illustrate (Figure 4), command templates provide
374 flexibility in the way that commands are constructed, including the use of multiple commands.
375 They provide the additional benefit that inputs do not have an inherent order and thus can be
376 specified in the command template in any order (and multiple times, if desired).
- 377 10. Outputs. Aside from any “Output Files” that may have been specified as inputs, the user may
378 declare additional output files. For instance, the `somatic/trim_fastq.cwl` example
379 specifies that trimmed FASTQ files should have the same names as the corresponding input

380 FASTQ files. To indicate this, the user specifies a CWL-based expression:
381 `$(inputs.fastq_file.basename)`. Additionally, if the user wishes to collect standard
382 output or standard error, they may specify the names of files that will store these messages.

383 After a user has specified all required elements, ToolJig generates a YAML document that conforms to
384 the CWL specification; the user may download this document. ToolJig also generates a form in which the
385 user can indicate a value for each input. Subsequently, the user can download an input-object file in
386 YAML format.

387 **Using ToolJig to create CWL workflows**

388 When using ToolJig to create a workflow, researchers first enter metadata: a unique identifier, label,
389 description, author information, and software license. Subsequently, the researcher uploads at least one
390 tool description. The researcher then defines two or more workflow steps. For each step, they specify a
391 unique name and the tool that will be used in that step. For each of the tool's inputs, the researcher
392 indicates whether the input will be populated from the output of a preceding step or from a workflow
393 input. The user may also indicate that any of the tool's outputs will become outputs for the overall
394 workflow. As with tool descriptions, ToolJig validates the user's input and then generates a CWL
395 document and input-object file that can be downloaded.

396 **Discussion**

397 Progress in biology research is hindered when software tools are difficult to install, when inputs and
398 outputs are inadequately or inconsistently specified, and when it is difficult to combine tools into
399 workflows. The CWL specifications—in combination with package managers and software containers—
400 are helping to alleviate these longstanding challenges. Moreover, CWL tool and workflow descriptions
401 can facilitate reproducible research. Rather than simply providing analysis code alongside journal
402 manuscripts, researchers can provide CWL documents. As illustrated in our examples, CWL documents

403 provide instructions for executing analyses in software containers that encapsulate all relevant
404 dependencies (Figure 5), along with ancillary scripts and instructions for accessing data files. People who
405 read (or review) the manuscripts can then repeat the analyses, without needing to install any software
406 other than a relevant workflow engine and container engine, even if their operating system or
407 configuration differs from the authors’.

408 Multiple online repositories provide CWL documents, including the Dockstore tool registry⁵⁶ and GitHub.
409 A search on GitHub for CWL documents that use the FastQC software⁵⁷, for example, resulted in 616
410 matches (May 15, 2020). Researchers can reuse and adapt these documents. However, in cases where
411 reuse is infeasible or extensive adaptations are necessary, scientific progress may be accelerated as
412 researchers, including non-bioinformaticians, gain greater efficiency in creating CWL documents. ToolJig
413 aims to facilitate this process, enabling researchers to build CWL tools interactively, without needing to
414 gain a deep understanding of the CWL specification or YAML syntax.

415 One advantage of CWL is that it can be used with diverse workflow engines. Whether or not they support
416 CWL, most workflow engines provide custom languages or programming interfaces for creating
417 workflows. Relatively little support is available for migrating from these engine-specific solutions to
418 CWL in an automated manner. However, when these engines support execution within Docker-
419 compatible containers, researchers can migrate these tools manually using ToolJig (or other means).
420 Providing better support in existing workflow engines for exporting to CWL will be a positive step
421 toward ensuring that CWL truly becomes a common language for command-line tools and workflows.

422 The CWL specification provides considerable flexibility for describing command-line tools and
423 workflows. Our goal was to support common use cases for biology research. For the sake of simplicity
424 and to reduce barriers of entry for new creators of CWL documents, we deliberately do *not* support some
425 optional features within the CWL specification. These include input directories, dependent and exclusive
426 parameters, process requirements, hints, and output directories. To learn about augmenting CWL
427 documents with these features, researchers can consult the CWL specifications.

428 ToolJig has no dependencies other than a modern Web browser. Accordingly, it can be used from
429 virtually any computer with no installation process. When ToolJig is updated, the user simply needs to
430 refresh their browser. A tradeoff to this simplicity is that ToolJig does not provide a direct means of
431 testing tools or workflows. However, the `cwltest` utility provides a command-line testing framework,
432 enabling researchers to compare tool and workflow outputs with expected results. Researchers
433 implementing CWL in production systems would benefit from using such a utility for validation.

434 **Conclusions**

435 CWL documents can formalize execution of command-line tools and workflows. We have summarized
436 the key components of CWL tool descriptions and provided examples to illustrate these concepts. In
437 addition, we have described ToolJig, a Web application that enables researchers to create CWL
438 documents interactively. We hope these resources will benefit researchers from diverse backgrounds to
439 more easily create CWL documents and thus advance sharing of methods and computational
440 reproducibility.

441 **Declarations**

442 *Ethics approval and consent to participate*

443 Not applicable.

444 *Consent for publication*

445 Not applicable.

446 *Availability of data and material*

447 The source code and examples are located at <https://github.com/srp33/ToolJig>.

448 *Competing interests*

449 The authors declare that they have no competing interests.

450 *Funding*

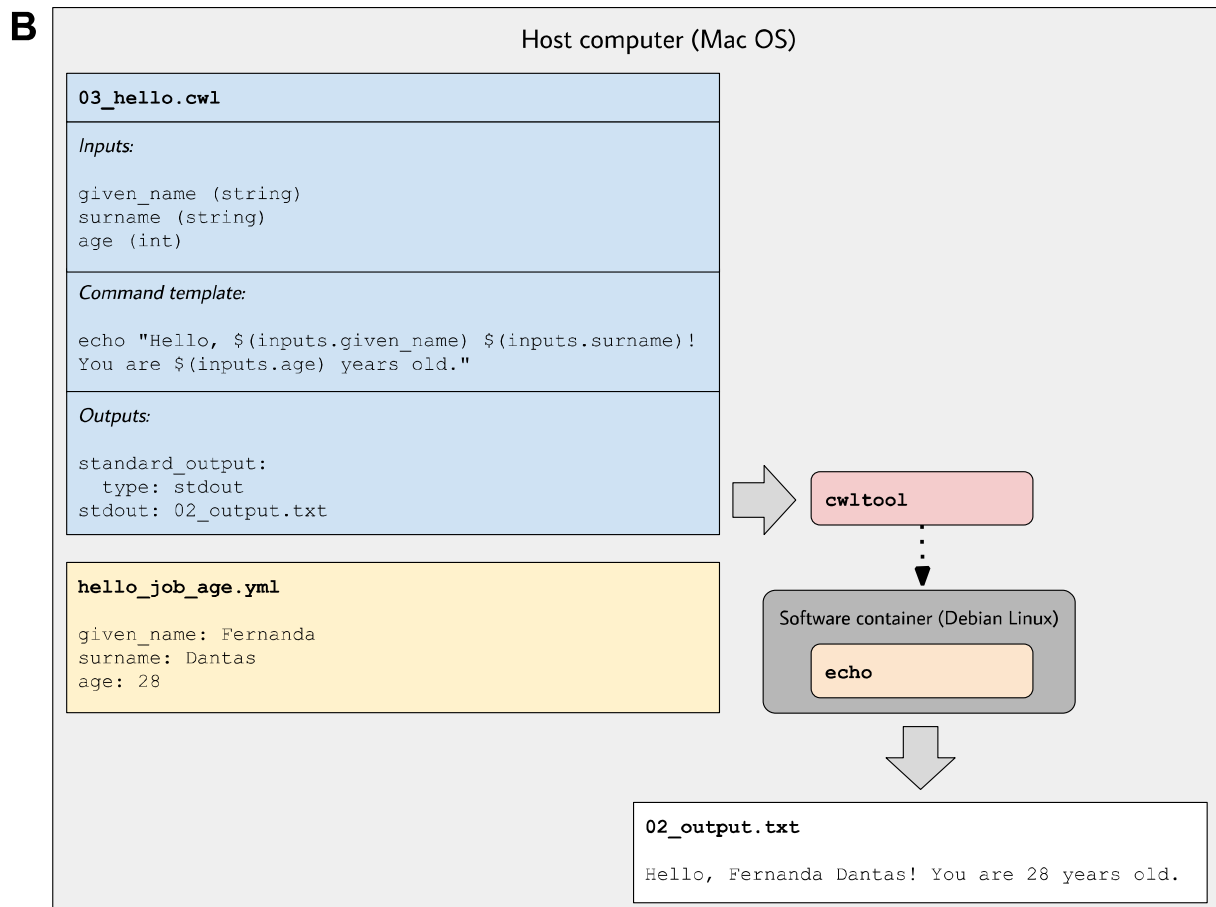
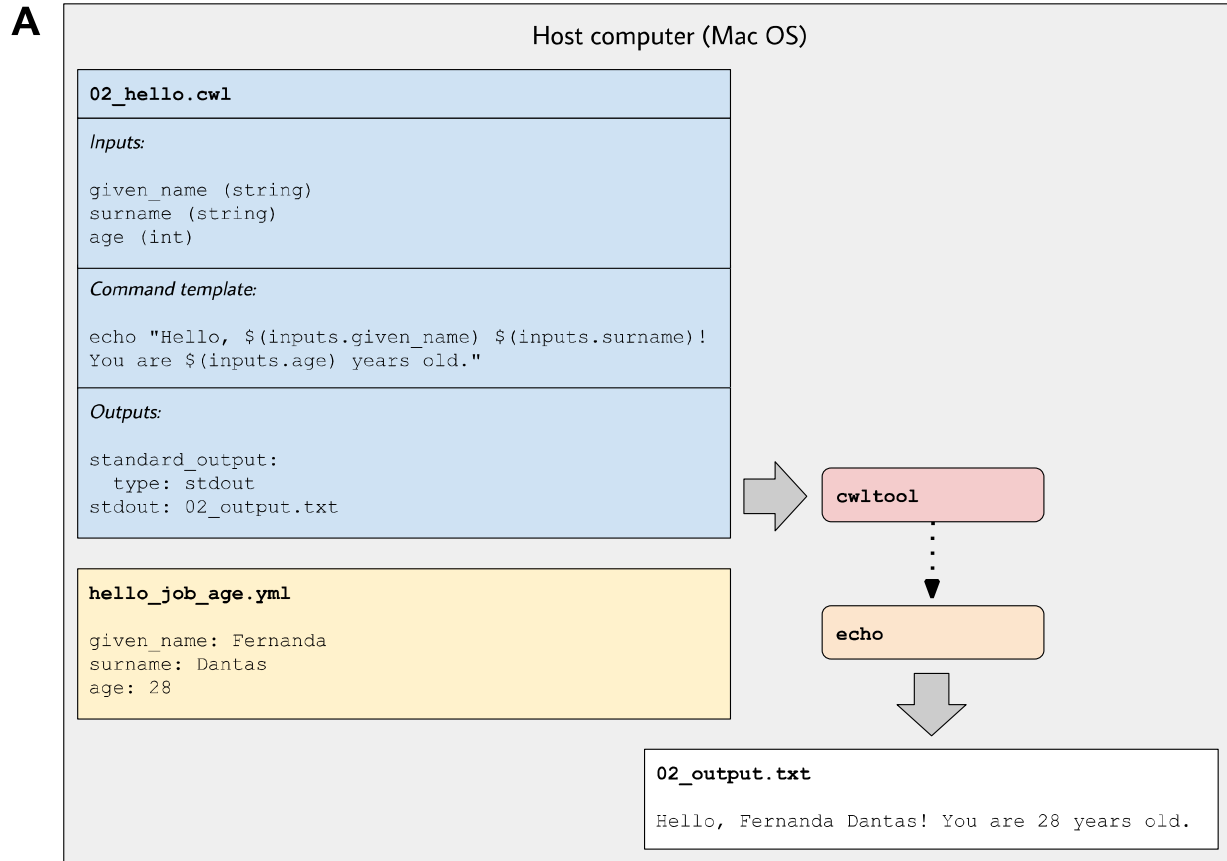
451 This work was supported by a grant from the United States National Institutes of Health (U54CA209978)
452 to SRP, AHB, JTC, and ZEE.

453 *Authors' contributions*

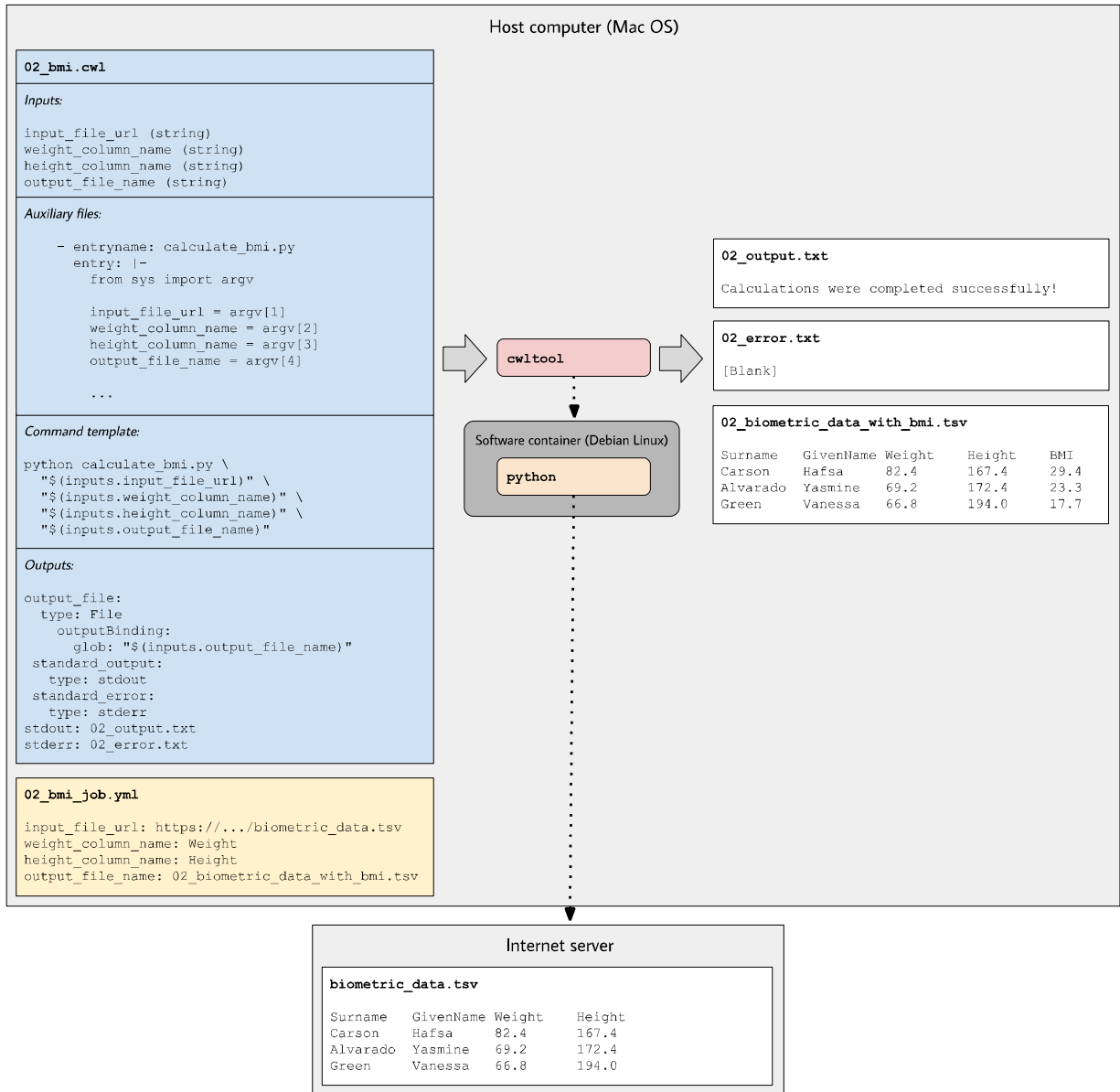
454 SRP, ECA, and ZEE conceptualized the ToolJig application and generated prototypes. SRP and ECA
455 created the application. SRP and ZEE prepared the tool examples. JTC and AHB provided critical
456 feedback on the application's functionality. SRP wrote the manuscript. All authors provided critical
457 feedback on the manuscript.

458 We acknowledge the Texas Cancer Research Biobank and Baylor College of Medicine Human Genome
459 Sequencing Center for providing openly available cancer-genome data used in our examples.

460 **Figures**



462 **Figure 1: Illustration of tool descriptions for printing simple greetings.** In the examples associated
463 with this article, we provide tool descriptions that illustrate how to print custom greetings at the command
464 line. These diagrams illustrate the `02_hello.cwl` (**A**) and `03_hello.cwl` (**B**) examples. In **A**, the
465 tool description indicates which inputs that must be specified, along with a template for executing the
466 command; it also indicates that a message will be printed to standard output and that this message should
467 be stored in a file called `02_output.txt`. The `hello_objects_age.yml` input-objects file stores
468 values for a particular invocation of the tool. In **A**, the `cwltool` workflow engine uses the host
469 computer's operating system to execute the tool; thus the `echo` command must be supported on that
470 operating system. In **B**, the tool description defines a software container environment; thus `cwltool`
471 executes the command within a container, which provides the `echo` command (packaged with the Debian
472 Linux operating system).

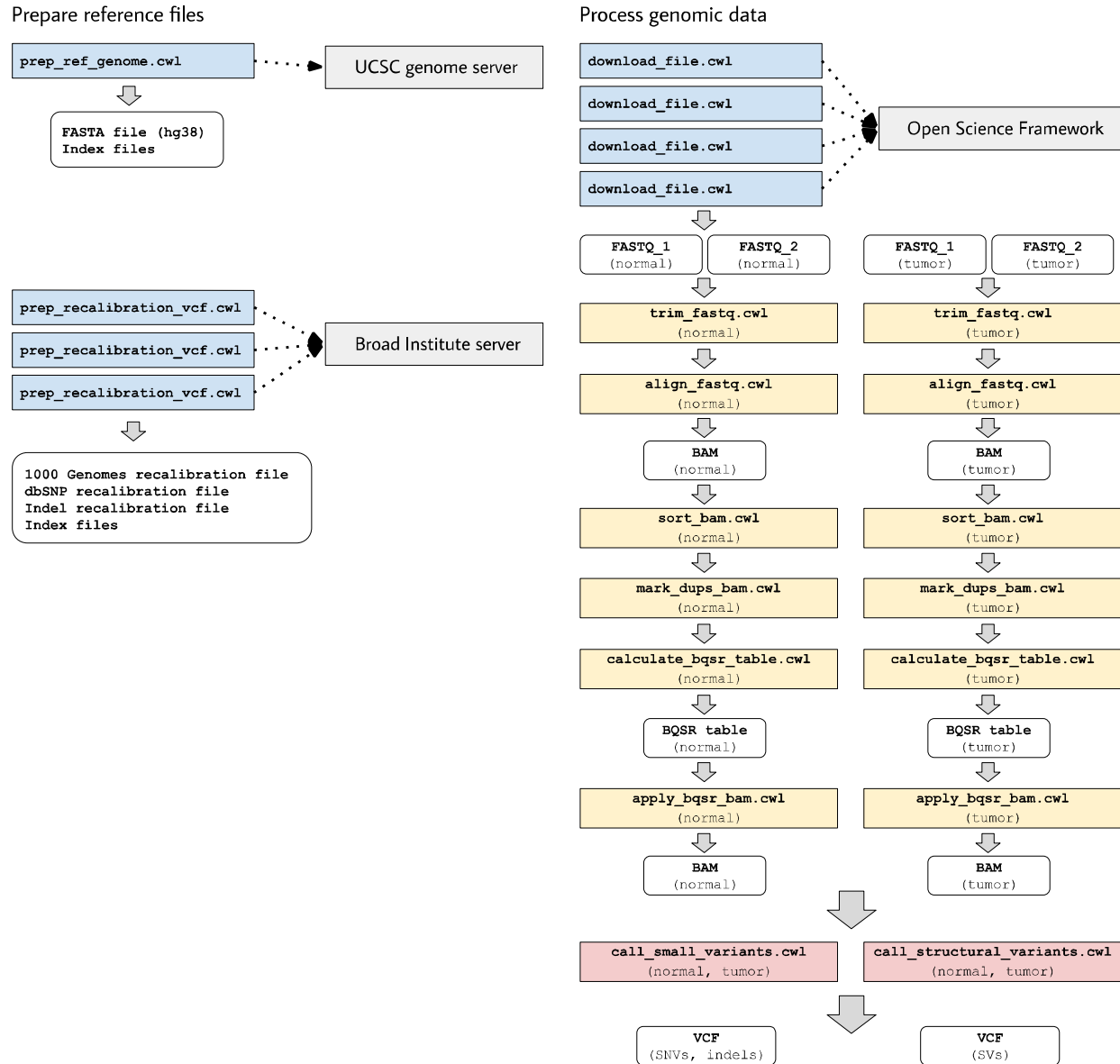


473

474 **Figure 2: Illustration of tool descriptions for calculating individuals' body mass index.** In the
 475 examples associated with this article, we provide tool descriptions that illustrate how to calculate body
 476 mass index (BMI) values based on individuals' weights and heights stored in a tab-separated-value file.
 477 This diagram illustrates the 02_bmi.cwl example. The tool description indicates the expected inputs. In
 478 this case, the URL of a data file must be provided. That file must contain a column that stores weights (in
 479 kilograms) and a column that stores heights (in centimeters). In the input-objects file
 480 (02_bmi_objects.yml), the user specifies the names of these columns. The final input is the name of

481 an output file that will be generated. This file will store the original data and a new column with the
482 calculated BMI value for each individual. As the tool executes, python (within a software container)
483 downloads the input file, performs the calculations, generates the output file, and stores standard output
484 and standard error in text files.

485



486

487 **Figure 3: Illustration of tool descriptions for calling somatic variants from a cancer genome.** In the
 488 examples associated with this article, we provide tool descriptions that illustrate how to call somatic
 489 variants from second-generation sequencing data for a cancer genome (compared against a normal
 490 genome from the same patient). This process requires execution of 11 distinct tools in a defined
 491 succession of steps (a workflow). Two tools (`prep_ref_genome.cwl` and `prep_recalibration_vcf.cwl`)
 492 prepare reference files associated with a given human reference genome. These tools download data files
 493 from public Internet servers and then create index files and standardize contig identifiers. The third tool

494 (download_file.cwl) downloads FASTQ files from an Internet server. The remaining tools process the
495 normal and tumor sequences separately before comparing the tumor genome against the normal genome
496 to identify single nucleotide variants, indels, and structural variants.

497

```
A echo "Hello, ${inputs.given_name} ${inputs.surname}! You are ${inputs.age} years old."
```

```
B Rscript run_deseq2_analysis.R "${inputs.read_counts_url}" "${inputs.phenotypes_url}" \  
    "${inputs.design_formula}" "${inputs.output_file_name}"
```

```
C bwa mem -t ${inputs.threads} ${inputs.args} \  
    -R "${inputs.read_group_string}" \  
    "${inputs.fasta_file.path}" \  
    "${inputs.fastq_file_1.path}" \  
    "${inputs.fastq_file_2.path}" | samtools view -b > "${inputs.output_file_name}"
```

```
D sambamba sort -t ${inputs.threads} -o "${inputs.output_file_name}"  
    "${inputs.bam_file.path}"  
  
sambamba index -t ${inputs.threads} "${inputs.output_file_name}"
```

```
E wget ${inputs.exclude_template_url}  
  
delly call -x `basename "${inputs.exclude_template_url}"` \  
    -o output.bcf -g "${inputs.fasta_file.path}" \  
    "${inputs.tumor_bam_file.path}" "${inputs.normal_bam_file.path}"  
  
bcftools view output.bcf > "${inputs.output_file_name}"
```

498

499 **Figure 4: Examples of command templates used in CWL tool descriptions.** These examples illustrate
500 diverse types of command templates for configuring execution of CWL tools. In each example,
501 placeholders are used for inputs. When the tools are executed, the placeholders are replaced with input-
502 objects values. **A)** A simple greeting is printed to standard output. **B)** An R script (stored as an auxiliary
503 file within the tool description) is executed; this script performs a differential-expression analysis using
504 the DESeq2 package. **C)** The *bwa* software aligns FASTQ files to a reference genome and pipes the
505 output to *samtools*; the output is then converted to BAM format. This example illustrates a scenario in
506 which two complementary software packages are used to perform a data-analysis task. Although these
507 packages could be incorporated into distinct CWL tools, we combine them because read alignment and
508 BAM conversion are typically performed jointly. **D)** The *sambamba* software sorts and then indexes a
509 BAM file. **E)** The *Delly* software identifies structural variants in a cancer genome. *Delly* can be
510 configured to exclude telomere and centromere regions as well as unplaced contigs. This example

511 downloads an exclusion file, invokes *Delly*, and converts the output to VCF format. Examples **D** and **E**
512 illustrate additional scenarios in which related tasks are executed as practical units.

513

514

```
A dockerPull: python:3.8.2-slim-buster

B dockerImageId: prep_recalibration_vcf
    dockerFile: |-
        FROM quay.io/biocontainers/picard:2.22.3--0

C dockerImageId: call_structural_variants
    dockerFile: |-
        FROM biocontainers/biocontainers:v1.0.0_cv4

        RUN conda install -c bioconda/label/cf201901 delly bcftools -y

D dockerImageId: scan_normalize
    dockerFile: |-
        FROM bioconductor/bioconductor_docker:RELEASE_3_10

        RUN R -e 'BiocManager::install("SCAN.UPC")'
```

515

516 **Figure 5: Examples of DockerRequirement specifications used in CWL tool descriptions.** These
517 examples illustrate diverse ways to configure CWL tools to be executed in software containers. In **A**, a
518 container image is pulled from Docker Hub; this image encapsulates a minimal (“slim”) version of
519 Debian Linux 10.3 (“buster”) and includes the Python 3.9 interpreter. In **B**, the contents of a Dockerfile
520 are included within the CWL description. In this case, the Dockerfile is simple—it pulls an existing image
521 from quay.io. This image is provided as part of the Biocontainers project and includes the Picard Tools
522 software. **C** uses a base image from Biocontainers and the Bioconda package manager to install the *Delly*
523 and *bcftools* software within the image. **D** uses a base image from Bioconductor and executes R code to
524 install the SCAN.UPC normalization package within the image.

525 **References**

- 526 1. Hey, T., Tansley, S., Tolle, K. & others. *The fourth paradigm: Data-intensive scientific*
527 *discovery*. vol. 1 (Microsoft research Redmond, WA, 2009).
- 528 2. Wilson, G. *et al.* Best Practices for Scientific Computing. *PLOS Biology* **12**, e1001745 (2014).
- 529 3. Piccolo, S. R. & Frampton, M. B. Tools and techniques for computational reproducibility.
530 *Gigascience* **5**, 30 (2016).
- 531 4. Hong, N. C. We are the 92%. (2014) doi:[10.6084/m9.figshare.1243288.v1](https://doi.org/10.6084/m9.figshare.1243288.v1).
- 532 5. Kumar, S. & Dudley, J. Bioinformatics software for biologists in the genomics era.
533 *Bioinformatics* **23**, 1713–1717 (2007).
- 534 6. Leipzig, J. A review of bioinformatic pipeline frameworks. *Brief Bioinform* **18**, 530–536 (2017).
- 535 7. Köster, J. & Rahmann, S. Snakemakea scalable bioinformatics workflow engine. *Bioinformatics*
536 **28**, 2520–2522 (2012).
- 537 8. Di Tommaso, P. *et al.* Nextflow enables reproducible computational workflows. *Nat. Biotechnol.*
538 **35**, 316–319 (2017).
- 539 9. Boettiger, C. An introduction to Docker for reproducible research. *ACM SIGOPS Oper. Syst. Rev.*
540 **49**, 71–79 (2015).
- 541 10. Amstutz, P. *et al.* Common workflow language, v1.0. (2016)
542 doi:[10.6084/m9.figshare.3115156.v2](https://doi.org/10.6084/m9.figshare.3115156.v2).
- 543 11. Vivian, J. *et al.* Toil enables reproducible, open source, big biomedical data analyses. *Nat.*
544 *Biotechnol.* **35**, 314–316 (2017).
- 545 12. Kotliar, M., Kartashov, A. V. & Barski, A. CWL-Airflow: A lightweight pipeline manager
546 supporting Common Workflow Language. *Gigascience* **8**, (2019).
- 547 13. Lee, S. *et al.* Tibanna: Software for scalable execution of portable pipelines on the cloud.
548 *Bioinformatics* **35**, 4424–4426 (2019).
- 549 14. Home | Arvados. <https://arvados.org>

- 550 15. Common Workflow Language User Guide. https://www.commonwl.org/user_guide/
- 551 16. Rabix: Power tools for the Common Workflow Language. *Rabix: Power tools for the Common*
552 *Workflow Language*. <https://github.com/rabix/composer>
- 553 17. Ian Holmes on Twitter: "You can download our code from the URL supplied. Good luck
554 downloading the only postdoc who can get it to run, though #overlyhonestmethods" / Twitter. *Twitter*.
555 <https://twitter.com/ianholmes/status/288689712636493824>
- 556 18. Bioconda: Sustainable and comprehensive software distribution for the life sciences | Nature
557 Methods.
- 558 19. Huber, W. *et al.* Orchestrating high-throughput genomic analysis with Bioconductor. *Nat.*
559 *Methods* **12**, 115–121 (2015).
- 560 20. Docker. *Docker*. <https://www.docker.com>
- 561 21. Kurtzer, G. M., Sochat, V. & Bauer, M. W. Singularity: Scientific containers for mobility of
562 compute. *PLOS ONE* **12**, e0177459 (2017).
- 563 22. Gomes, J. *et al.* Enabling rootless Linux Containers in multi-user environments: The udocker
564 tool. *Computer Physics Communications* **232**, 84–97 (2018).
- 565 23. CoreOS. <https://coreos.com/rkt>
- 566 24. Podman. *podman.io*.
- 567 25. Hpc/charliecloud. (2020). <https://github.com/hpc/charliecloud>
- 568 26. Foster, E. D. & Deardorff, A. Open Science Framework (OSF). *J Med Libr Assoc* **105**, 203–206
569 (2017).
- 570 27. Van Rossum, G. & others. Python Programming Language. in *USENIX Annual Technical*
571 *Conference* vol. 41 36 (2007).
- 572 28. R Core Team. *R: A language and environment for statistical computing*. (R Foundation for
573 Statistical Computing, 2020).

- 574 29. Jupyter, P. *et al.* Binder 2.0 - Reproducible, interactive, sharable environments for science at
575 scale. in *Proceedings of the 17th Python in Science Conference* (eds. Akici, F., Lippa, D., Niederhut, D.
576 & M Pacer) 113–120 (2018). doi:[10.25080/Majora-4af1f417-011](https://doi.org/10.25080/Majora-4af1f417-011).
- 577 30. Amstutz, P. & Crusoe, M. R. Common Workflow Language (CWL) Command Line Tool
578 Description, v1.1.
- 579 31. The Official YAML Web Site. <https://yaml.org>
- 580 32. JSON. <https://www.json.org/json-en.html>
- 581 33. Piccolo, S. R. *et al.* A single-sample microarray normalization method to facilitate personalized-
582 medicine workflows. *Genomics* **100**, 337–44 (2012).
- 583 34. Barrett, T. *et al.* NCBI GEO: Archive for functional genomics data sets years on. *Nucleic Acids*
584 *Res.* **39**, D1005–10 (2011).
- 585 35. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for
586 RNA-seq data with DESeq2. *Genome Biol.* **15**, 550 (2014).
- 587 36. Wickham, H., Hester, J. & Francois, R. Readr: Read Rectangular Text Data. (2018).
- 588 37. Wickham, H., François, R., Henry, L. & Müller, K. *Dplyr: A Grammar of Data Manipulation*.
589 (2018).
- 590 38. Bottomly, D. *et al.* Evaluating Gene Expression in C57BL/6J and DBA/2J Mouse Striatum Using
591 RNA-Seq and Microarrays. *PLOS ONE* **6**, e17820 (2011).
- 592 39. Becnel, L. B. *et al.* An open access pilot freely sharing cancer genomic data from participants in
593 Texas. *Sci. Data* **3**, 1–10 (2016).
- 594 40. Haeussler, M. *et al.* The UCSC Genome Browser database: 2019 update. *Nucleic Acids Res* **47**,
595 D853–D858 (2019).
- 596 41. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform.
597 *Bioinforma. Oxf. Engl.* **25**, 1754–60 (2009).
- 598 42. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinforma. Oxf. Engl.* **25**,
599 2078–9 (2009).

- 600 43. Picard Tools - By Broad Institute. <http://broadinstitute.github.io/picard/>
- 601 44. Danecek, P. *et al.* The variant call format and VCFtools. *Bioinformatics* **27**, 2156–2158 (2011).
- 602 45. Depristo, M. A. *et al.* A framework for variation discovery and genotyping using next-generation
603 DNA sequencing data. *Nat. Genet.* **43**, (2011).
- 604 46. Didion, J. P., Martin, M. & Collins, F. S. Atropos: Specific, sensitive, and speedy trimming of
605 sequencing reads. *PeerJ* **5**, e3720 (2017).
- 606 47. Tarasov, A., Vilella, A. J., Cuppen, E., Nijman, I. J. & Prins, P. Sambamba: Fast processing of
607 NGS alignment formats. *Bioinformatics* **31**, 2032–2034 (2015).
- 608 48. Benjamin, D. *et al.* Calling Somatic SNVs and Indels with Mutect2. *bioRxiv* 861054 (2019)
609 doi:10.1101/861054.
- 610 49. Rausch, T. *et al.* DELLY: Structural variant discovery by integrated paired-end and split-read
611 analysis. *Bioinformatics* **28**, i333–i339 (2012).
- 612 50. da Veiga Leprevost, F. *et al.* BioContainers: An open-source and community-driven framework
613 for software standardization. *Bioinformatics* **33**, 2580–2582 (2017).
- 614 51. Gruening, B. *et al.* Recommendations for the packaging and containerizing of bioinformatics
615 software. *F1000Research* **7**, 742 (2019).
- 616 52. Vue.js. <https://vuejs.org>
- 617 53. Nüst, D. *et al.* Ten simple rules for writing Dockerfiles for reproducible data science. *PLOS*
618 *Computational Biology* **16**, e1008316 (2020).
- 619 54. Haak, L. L., Fenner, M., Paglione, L., Pentz, E. & Ratner, H. ORCID: A system to uniquely
620 identify researchers. *Learn. Publ.* **25**, 259–264 (2012).
- 621 55. Ison, J. *et al.* EDAM: An ontology of bioinformatics operations, types of data and identifiers,
622 topics and formats. *Bioinformatics* **29**, 1325–1332 (2013).
- 623 56. O’Connor, B. D. *et al.* The Dockstore: Enabling modular, community-focused sharing of Docker-
624 based genomics tools and workflows. *F1000Res* **6**, (2017).

- 625 57. Brown, J., Pirrung, M. & McCue, L. A. FQC Dashboard: Integrates FastQC results into a web-
626 based, interactive, and extensible FASTQ quality control tool. *Bioinformatics* **33**, 3137–3139 (2017).