

1 NextflowWorkbench: Reproducible and 2 Reusable Workflows for Beginners and 3 Experts.

4 Jason P. Kurs¹, Manuele Simi^{1,2}, and Fabien Campagne^{1,2,3,*}

5 ¹The HRH Prince Alwaleed Bin Talal Bin Abdulaziz Alsaud Institute for Computational
6 Biomedicine, Weill Cornell Medicine, New York, NY, United States of America

7 ²Clinical Translational Science Center, Weill Cornell Medicine, New York, NY, United
8 States of America

9 ³Department of Physiology and Biophysics, Weill Cornell Medicine, New York, NY,
10 United States of America

11 *To whom correspondence should be addressed: fac2003@campagnelab.org

12 ABSTRACT

13 Computational workflows and pipelines are often created to automate series of processing steps. For instance, workflows enable one to standardize analysis for large projects or core facilities, but are also useful for individual biologists who need to perform repetitive data processing. Some workflow systems, designed for beginners, offer a graphical user interface and have been very popular with biologists. In practice, these tools are infrequently used by more experienced bioinformaticians, who may require more flexibility or performance than afforded by the user interfaces, and seem to prefer developing workflows with scripting or command line tools. Here, we present a workflow system, the NextflowWorkbench, which was designed for both beginners and experts, and blends the distinction between user interface and scripting language. This system extends and reuses the popular Nextflow workflow description language and shares its advantages. In contrast to Nextflow, NextflowWorkbench offers an integrated development environment that helps complete beginners get started with workflow development. Auto-completion helps beginners who do not know the syntax of the Nextflow language. Reusable processes provide modular workflows. Programmers will benefit from unique interactive features that help users work more productively with docker containers. We illustrate this tool with a workflow to estimate RNA-Seq counts using Kallisto. The workflow can transparently run either on a laptop computer with docker or on a linux cluster. We found that beginners can be taught how to assemble this workflow in a two hours training session. In conclusion, the NextflowWorkbench simplifies the development of reproducible, implicitly parallel workflows. NextflowWorkbench is distributed under the Apache 2.0 license and available at <http://workflow.campagnelab.org>.

14 Keywords: Workflows, Pipelines, Reproducibility, Docker, Language Workbench, JetBrains MPS

15 INTRODUCTION

16 A computational workflow or pipeline is a description of a series of computational steps connected to each
17 other. Each step accepts one or more input(s) and transforms input(s) to produce one or more output(s).
18 Computational workflows are used in many engineering and scientific domains, but are particularly useful
19 in fields such as bioinformatics where analysis activities are repetitive and benefit from being automated.
20 Workflows can be represented as diagrams and their steps followed manually, but many solutions have
21 been developed to represent workflows electronically and automate their execution.

22 Automated workflow systems include a way for a user to edit the formal representation of the workflow,
23 and its component steps, as well as a runtime system to execute specific workflows. Two broad families
24 of workflow systems have been developed and are still in use today.

25 The first category is workflows with graphical user interfaces, which often represent workflows
26 as connected components on a 2D diagram. There is a long history for such tools, but Galaxy and
27 GenePattern are well known examples in Bioinformatics Giardine et al. [2005], Reich et al. [2006].
28 Workflow systems with graphical user interfaces are favored by beginners, or by educators who need

29 to teach beginners (see this thread, for instance, where these arguments have been made by others
30 <https://www.biostars.org/p/50034/>).

31 The second category of workflow systems is based on programming and scripting languages. A
32 workflow is expressed in a declarative or imperative language, or a combination of both. An example of
33 declarative language is the makefile input of the make/gmake tool available in UNIX, which was initially
34 developed to automate the compilation of programs, but has been used as well to implement workflows.
35 An example of imperative workflow is when bioinformaticians develop a collection of scripts and execute
36 these scripts to implement different pipelines.

37 While useful, many of these systems fail to fulfill needs that are common in bioinformatics. For
38 instance a makefile cannot easily be run in a distributed manner on multiple nodes of a cluster to parallelize
39 the processing of large collections of data. Scripts can be run in parallel on a cluster with tools like GNU
40 parallel or grid schedulers (e.g., Sun Grid Engine or SLURM), but a recurrent problem with scripts is that
41 most of them are not easily portable to new environments. Indeed, most scripts are written to assume a
42 specific location for the programs that they use (either the script expects a program to be in the PATH, or
43 the script contains a dependency on some location defined in the system where the script was developed).
44 When a dependency is not available in the location assumed by the script, the script fails and the user
45 needs to resolve the issue before retrying execution. The problem of managing dependencies for scripts
46 is often described as “dependency hell”, a term which many people believe adequately describes the
47 practical difficulties of getting scripts to run on other systems than where they have been developed.

48 Several groups have recently recognized these problems and have developed improved solutions
49 targeted at bioinformaticians. Recent developments include BigDataScript Cingolani et al. [2015] and
50 Nextflow Di Tommaso et al. [2014] (<http://nextflow.io>). These solutions address scalability and
51 portability problems and are useful to users with programming and/or scripting experience.

52 We recently asked the question of whether we could design a hybrid between scripting and user
53 interface workflow systems. Such a system would make it possible for beginners and experts to collaborate
54 more closely by using the same platform to represent and execute workflows. This manuscript describes the
55 NextflowWorkbench, a workflow platform that addresses this question. We developed NextflowWorkbench
56 with Language Workbench Technology implemented in the JetBrains MPS system ([jetbrains.com/
57 mps](http://jetbrains.com/mps)) Campagne [2014, 2015], Simi and Campagne [2014], Benson and Campagne [2015]. This platform
58 leverages Nextflow, a workflow language developed for users familiar with the command line and
59 scripting, which offers many advantages also available in NextflowWorkbench. We report on the design
60 of NextflowWorkbench and our experience teaching this new workflow system to biologists and clinicians
61 with no prior scripting experience.

62 **METHODS**

63 We have used the MPS Language Workbench (<http://jetbrains.com/mps>), as also described
64 in Campagne [2014] and Campagne [2015]. For an introduction to Language Workbench Technology
65 (LWT) in the context of bioinformatics see Simi and Campagne [2014] and Benson and Campagne [2015],
66 Campagne and Simi [2015] in the context of data analysis. Methods for this study are similar to those
67 described in Campagne et al. [2015].

68 **Language Design**

69 JPK designed and developed the core MPS languages of NextflowWorkbench during a three month
70 summer internship in the Campagne laboratory. These core languages include functionality to represent
71 Processes, Workflow and execution Scripts. Additional developments were conducted by MS and FC to
72 extend these languages and add docker (<https://docker.com>) and GobyWeb functionality. Full
73 language development logs are available on the GitHub code repository ([https://github.com/
74 CampagneLaboratory/NextflowWorkbench](https://github.com/CampagneLaboratory/NextflowWorkbench)). Briefly, we designed abstractions to represent
75 Nextflow scripts in a modular fashion (decoupling Processes from Workflows). These abstractions were
76 implemented with the structure, editor, constraints and typesystem aspects of MPS languages (described
77 in Campagne [2014]). Importantly, the typesystem aspect makes it possible to typecheck a workflow as it
78 is being developed and provide feedback to the developer. Nextflow scripts are generated from nodes of
79 the languages using the MPS textgen aspect (Campagne [2014]).

80 Workflow Execution

81 Workflows can be executed directly from within the MPS LW. Execution is supported on the developers' machine (for Linux and Mac OS platforms only, since Nextflow does not run on Windows), or on a remote 82 Linux node, where one of the execution mechanisms supported by Nextflow must be available (e.g., 83 Sun Grid Engine, SLURM, Apache Ignite). This capability was implemented with Run Configurations 84 (see Campagne [2015], Chapter 5). 85

86 Scripts

87 To implement scripts as text with auto-completion, we used the MPS RichText plugin (developed and 88 distributed by members of the MBEDDR project Voelter et al. [2012]). The plugin implements the 89 approach described in Voelter [2013].

90 Documentation

91 Project documentation has been developed with L^AT_EX and the Editor2PDF language and plugin (<https://github.com/CampagneLaboratory/Editor2PDF>). Complete documentation is available at 92 Kurs and Campagne [2015]. 93

94 RESULTS

95 **A Workflow System for Beginners and Experts** In this study, we designed a workflow system aimed 96 at the full spectrum of workflow users, from beginners to computational experts. Figure 1 presents the 97 advantages of this new workflow system across its intended spectrum of users. The following results 98 section describes the design of this system and the innovations introduced to help with the development 99 and maintenance of reproducible and high-performance workflows. The section also addresses the 100 question of whether this system can be taught effectively to beginners with no programming or command 101 line experience.

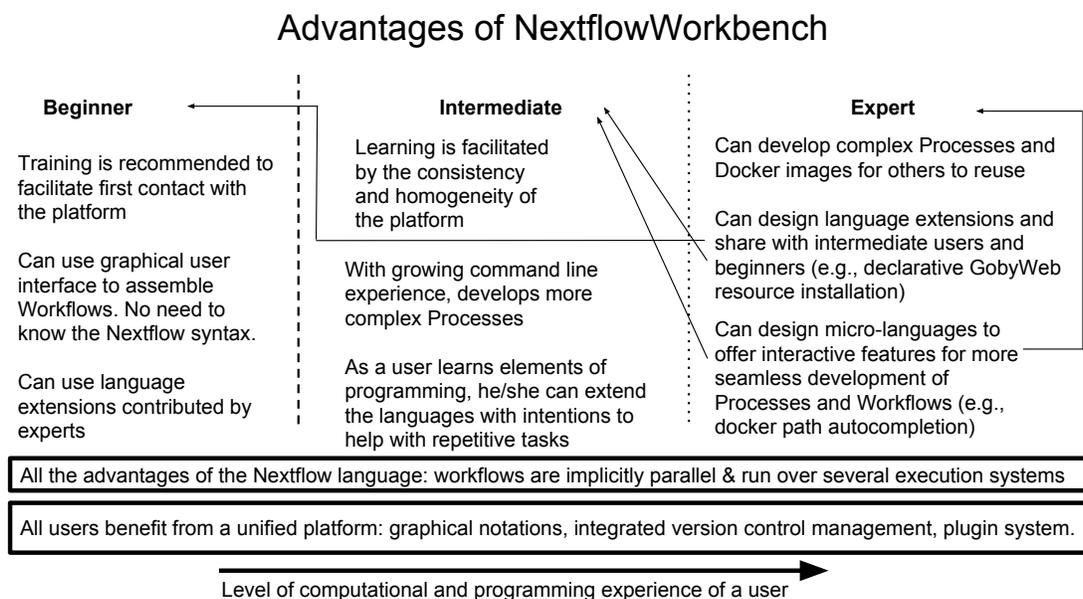


Figure 1. Advantages of the NextflowWorkbench workflow system across its intended user spectrum.

102 **Evaluation of Workflow Systems** We selected Nextflow as the target language for the NextflowWork-
103 bench platform after comparing three systems with similar goals. The comparison included BigDataScript,
104 Nextflow and the Swift language (<http://swift-lang.org/>). These systems were selected for
105 evaluation because they support parallel execution of workflows on multi-node clusters and provide a
106 reasonable level of abstraction to express workflows.

107 To compare these systems, two evaluators tried to implement a simple analysis pipeline with each of
108 them. One evaluator was an experienced software engineer with decades of programming and scripting
109 experience (MS). The second evaluator (JPK) was a sophomore undergraduate student with intermediate
110 level programming and scripting experience. Both evaluators completed the implementation of the test
111 pipeline with BigDataScript and Nextflow. The more experienced evaluator developed a partial prototype
112 workflow with Swift, but reported that locating information in the online documentation was tedious and
113 that the semantic of the language was far from intuitive.

114 This short evaluation revealed a number of characteristics, advantages and drawbacks of the three
115 systems:

116 **Swift** requires that all programs used in the workflow are already installed on each machine of a
117 cluster. A text file binds the name of a program used in the Swift script to the path of this program on the
118 machine. This requirement means that Swift assumes that dependencies of a workflow are pre-installed
119 and fixed, and makes no effort to facilitate the installation of programs and dependencies on cluster nodes.

120 **BigDataScript** was presented in detail in Cingolani et al. [2015]. During our evaluation, we
121 noticed a bug report in the forum that indicated a major error in dataflow analysis of the script by
122 the BigDataScript compiler/interpreter (see <https://groups.google.com/forum/#!topic/bigdatascript-users/r7rQ03LBYIc>). The magnitude of this error suggested to us that Big-
123 DataScript was not, at the time, a robust language for developing workflows.
124

125 **Nextflow** performed as expected and the evaluators found the documentation sometimes difficult
126 to follow, but overall fairly complete and sufficient. The language offers direct support for running
127 steps of a workflow— called Processes in the Nextflow language— inside a docker container. This
128 feature is extremely useful to develop workflows that can execute on other systems without tedious
129 dependency installation Di Tommaso et al. [2015]. Negatives were the lack of language modularity,
130 making it impossible to develop libraries of reusable processes and the difficulty of knowing at first glance
131 —at least for beginners who do not know the language well— what type of data is exchanged between
132 processes.

133 The evaluators concluded that of the three systems, Nextflow was the more promising system for
134 representing workflows as scripts.

135 **Requirements for an Improved Workflow Language** Following up this evaluation, we decided to
136 design a variant of the Nextflow language that would directly address the limitations that our evaluation
137 had identified. Specifically, we wanted:

- 138 • A modular workflow language that would make it straightforward to reuse processes developed
139 by others in new workflows. Modularity can enable experts to develop processes and share these
140 processes with beginners.
- 141 • An explicitly typed workflow language. We believe that an explicitly typed language makes it more
142 obvious to beginners what data are expected as input to a process and what data will be produced
143 as output. Coupled with a mechanism to check type compatibility (a type system) at runtime and
144 highlight type errors, explicit types make it easier for beginners to develop correct workflows. A
145 type system is also useful to experts because it highlights errors that could be missed and only
146 become apparent when trying to execute the workflow.
- 147 • A language that does not require the user to know or remember its precise syntax to use it. Such
148 languages can be built with language workbench technology to provide auto-completion that guides
149 beginners and experts. We have shown the effectiveness of these approaches to help beginners in
150 the MetaR project (Campagne et al. [2015]) and were curious to find out if they could also help
151 with workflow development.

152 **Process** A NextflowWorkbench Process is illustrated in Figure 2. A process consists of a set of inputs,
153 a set of outputs, and a script, which implements the processing of inputs to generate outputs.

154 In contrast to Nextflow, Processes in the NextflowWorkbench are created independently from a
155 Workflow script (i.e., outside the script). As standalone language constructs (implemented as MPS root
156 nodes), Processes can be developed and packaged to be shared with others (for instance as solutions
157 provided in MPS plugins).

158 As usual when developing a language with LWT, most parts of a Process can be extended by a user of
159 NextflowWorkbench using language composition. The next section describes an application of language

```
Process  Download_1M_Reads  docker inutano/sra-toolkit:latest
  <<no process option overrides>>

input:
  string id

output:
  tuple [ file '*_1.fastq.gz', file '*_2.fastq.gz' ]

script:
  fastq-dump -X 1000000 --split-files $id
  gzip $id_*.fastq
```

Figure 2. NextflowWorkbench Process. A Process defines inputs, outputs, an optional docker container, and a script. In the example shown, the process accepts an input called ‘id’ of type string. The string is used to query the Short Read Archive with the sra-toolkit and retrieve paired FASTQ files. Inputs must be available before the script can execute. Similarly to Nextflow, a process only executes successfully if the outputs it declared have been produced by the script execution. Notice how the Process incorporates graphical elements and colors to clearly mark different roles of the language elements. When a docker container is specified, as shown in this figure, the commands shown in the script will be executed inside the container. This semantic is implemented by the Nextflow execution runtime. NextflowWorkbench provides autocompletion for input arguments inside the script. This mechanism reduces the risk of typos in variable names and provides instant refactoring of variable names across the script when the workflow programmer renames an input variable.

160 composition where we extended the Script part of a Process with the ability to automatically install data
161 resources needed by the script.

162 **Processes with Variable Data Resources** Docker containers are useful to isolate the process from
163 the machine where the Process executed, but they have limitations. As we developed bioinformatics
164 workflows with the NextflowWorkbench, we found that docker is only a partial solution when a process
165 requires variable data resources.

166 As an example, consider the process shown in Figure 3. This Process estimates counts for RNA-seq
167 reads against a transcriptome. Different species have different transcriptomes, and different reference
168 transcriptomes or build versions exist for the same species. With the mechanisms provided by docker,
169 one would create different images for different combinations of species and reference build number. This
170 is not really practical because there are a large number of these combinations and only a few may be
171 of practical interest. If a core facility wanted to provide workflows to support multiple species, core
172 personnel would have to anticipate the needs of its user base and configure a large number of possible
173 images to support any species that the core users could need.

174 Rather than creating static images with all the data packaged in a container and for all possible choices
175 of interest, it can be more efficient to provide a mechanism to declare what specific resource is needed, and
176 use software already packaged inside the docker image to assemble data resources on demand. In this case,
177 assembling the transcriptome resource consists in downloading the appropriate transcriptome reference
178 sequence from Ensembl and indexing this reference with the Kallisto program. We have developed
179 mechanisms to support this on-demand strategy, and make it easy to construct specific data resources. The
180 use of these mechanisms is shown in Figure 3, where a simple `requires` block declares a dependency
181 on the `KALLISTO_INDEX` data resource. The same approach also supports choosing and installing
182 specific versions of software resources.

183 Resource installation scripts are built with the method previously developed for GobyWeb Dorff et al.
184 [2013]. Examples of resources configurations are distributed on GitHub at

185 [https://github.com/CampagneLaboratory/gobyweb2-plugins/tree/master/](https://github.com/CampagneLaboratory/gobyweb2-plugins/tree/master/plugins/resources)
186 `plugins/resources`

```
Process  KallistoCountsWithTuples  artifacts/kallisto-homo-sapiens :1.0.0
  <<no process option overrides>>

input:
  tuple [ file read1, file read2 ]

output:
  file 'counts-*.tsv'

script:
  requires { resource KALLISTO_INDEX 0.42.3 resolved as: KALLISTO_INDEX -> 0.42.3 }
            INDEX.organism = Homo_sapiens
            INDEX.reference-build = GRCH38
            INDEX.ensembl-version-number = 82
  #!/bin/bash (with automatic GobyWeb artifact installation)
  echo "Processing: " $read1
  TRANSCRIPT_INDEX=${artifact path KALLISTO_INDEX.INDEX}/transcripts_index
  echo ${TRANSCRIPT_INDEX}
  basename=`basename $read1`
  echo "Basename= ${basename}"

  mkdir output
  ${artifact path KALLISTO_BINARIES}/bin/kallisto quant --index=${TRANSCRIPT_INDEX} $read1 $read2
  --output-dir=./output
  #touch output/abundance.tsv
  ls -ltrR .
  cp output/abundance.tsv counts-${basename}.tsv
  exit 0
```

Figure 3. Process with GobyWeb Artifact Installation. This process uses a special type of script which declares dependencies on GobyWeb resources. GobyWeb resources can automatically install variable data resources, such as a specific transcriptome index identifier by species, reference build and Ensembl version number (as shown). In this example, the script requests installation of the KALLISTO_INDEX resource version 0.42.3. This resource was configured to retrieve the human transcriptome corresponding to GRCH38, in Ensembl version 82. Notice that rather than writing the complicated steps to download and index this transcriptome, the workflow developer can express the data dependency declaratively.

187 **Interactive Docker Features** Developing scripts that run inside a docker container can be challenging
 188 because the programmer has to know and remember what programs and data are available inside the
 189 container and their precise location. The traditional way to build this understanding is to use interactive
 190 console sessions manually started inside a container. The shell can then be used to inspect the files and
 191 programs available in the running container and the user has to copy and paste these locations in the
 192 script. In NextflowWorkbench, we developed an auto-completion feature that shows container directories
 193 interactively when writing the Process script. With this method, a developer can associate a docker
 194 image to a Process, start an interactive container, and use auto-completion in the script to locate files
 195 or directories of interest (see Figure 4 for an illustration) without even leaving the workbench. Writing
 196 correct paths becomes seamless and no longer requires switching between Process editor and console.
 197 Similar capabilities are available to auto-complete commands in PATH as well as finding the exact location
 198 of data or programs installed as GobyWeb resources (an example is shown in Figure 3 where `$artifact`
 199 paths can be assembled using auto-completion).

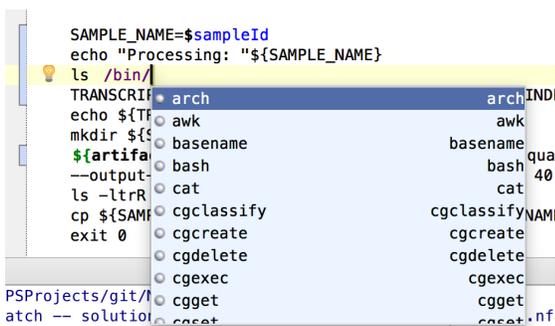


Figure 4. Interactive Path Auto-Completion. This figure illustrates interactive auto-completion of paths inside a docker container. After starting an interactive docker container, Process developers can auto-complete paths inside the docker container in the editor. Similar auto-completion functionality is offered for the GobyWeb data and program resources. These interactive features are implemented with standard features of the MPS language workbench.

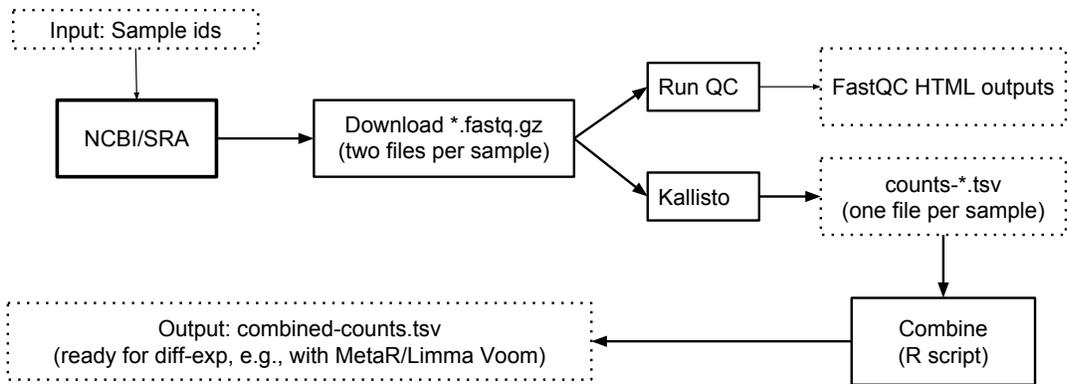


Figure 5. Diagram of Workflow. This diagram is a schematic representation of the analysis workflow shown in Figure 6.

200 **Workflow** A NextflowWorkbench Workflow consists of a set of inputs, references to processes and a
 201 list of optional report clauses. Assume that a user wishes to program a workflow to automate the analysis
 202 shown in Figure 5. Figure 6 illustrates how such a workflow can be expressed with the NextflowWorkbench
 203 language. In contrast to Nextflow, which define Processes inside a workflow script, a NextflowWorkbench
 204 Workflow accesses to processes by reference. Process references make it possible to name the Process'
 205 inputs and outputs in the context of the workflow. Such names are used to establish connections between
 206 process invocations. For instance, in Figure 6, the reference to `KallistoCountsWithTuples`

207 associates the name `B` to the input of the Process, and associates the name `result` to its output. When
208 the name `result` is defined in this way, the user becomes able to bind the name `result` in the input
209 role of another process reference.

210 In order to prevent cyclic dependencies, output names cannot be set on an input of the same process,
211 and can be set on one input at most. When an output needs to be consumed by several downstream
212 processes (e.g., to implement the fork after Download in Figure 5), the language offers an intention to
213 duplicate a name (an intention is a context dependent menu, see Simi and Campagne [2014]). In Figure 6,
214 the symbols $\rightarrow [A, B]$ indicates that the output of the `Download_1M_Reads` process is duplicated
215 and available through the names `A` and `B`. Note that the editor supports adding additional names between
216 the brackets to replicate the name as many times as needed.

Download reads from SRA, run FastQC, and estimate counts against the human transcriptome with Kallisto. Produce a
combine counts matrix.

Workflow  FastqKallis toCounts

```
with input:
  [ "SRR1514132", "SRR1514133", "SRR1514134", "SRR1514135", "SRR1514136", "SRR1514137", "SRR1514138", "SRR1514139",
    "SRR1514140", "SRR1514141" ] → [ IDstoDownload, IDstoCombine ]

do:
  IDstoDownload << ... >> → id ▶ Download_1M_Reads ▶ ['*_1.fastq.gz', '*_2.fastq.gz'] → [ A, B ] will run n times
  <<options>>;

  A << ... >> → [read1, read2] ▶ QC ▶ '*.zip' → zip will run n times <<options>>;

  B << ... >> → [read1, read2] ▶ KallistoCountsWithTuple ▶ 'counts-*.tsv' → result will run n times <<options>>;

  result.toList() → tsvs ▶ CombineCounts ▶ counts.tsv → combined will run 1 time <<options>>;
  IDstoCombine.toList() → ids

and report:
  << ... >>
```

Figure 6. Workflow Example. In this example, a set of SRA identifiers is defined in the input section of the workflow. The list is then duplicated and fed to two processes: `Download_1M_Reads` and `CombineCounts`. The `Download_1M_Reads` will run once for each identifier. When this process terminates, the output, which consists of the tuple of files `['*_1.fastq.gz', '*_2.fastq.gz']`, is duplicated to be fed to the `QC` and `KallistoCountsWithTuple` processes. The counts obtained with Kallisto are fed to `CombineCounts` along with the list of ids to produce a combined matrix of counts for all the samples analyzed.

217 **Utility to Non-Programmers** An important question is whether the user interface provided by Nextflow-
218 Workbench provides sufficient assistance to help non programmers with a biology or clinical background
219 develop workflows.

220 To address this question we developed training material and started teaching how to develop the work-
221 flow shown in Figure 6. In these training sessions, biologists and clinicians develop the `Download_1M_`
222 `Reads` and `QC` processes and reuse the `KallistoCountsWithTuple` and `CombineCounts` pro-
223 cesses from a library. Trainees are emailed instructions (<http://campagnelab.org/software/nextflow-workbench/instructions-for-workflow-tutorial/>) to install the software
224 on their laptop prior to the training session.
225

226 The main challenges we have encountered in these training sessions are related to the installation of
227 the software on the trainees machines. About 30-45 minutes of the training sessions are spent verifying
228 installations and performing some installation steps that the trainees have missed. In few cases, instructors
229 are unable to complete an installation because the trainee computers did not meet minimal specifications
230 (e.g., outdated operating system version, minimum required is Mac OS 10.8.3, or memory requirements
231 not met, minimum needed is 8GB to run Kallisto inside docker on a Mac laptop). Reducing the number
232 of installation steps and provisioning a cluster for remote execution of workflows would go a long way to
233 make this training more accessible and we are actively developing solutions to this end.

234 We found that we could teach trainees whose laptop met requirements how to assemble the workflow
235 shown in Figure 6 in less than two hours (including time to troubleshoot installations).

```
Build Image images: Kallisto_Image
Dockerfile Kallisto {
  FROM docker artifacts/software-gcc4.8:1.3.0 Pull
  MAINTAINER Campagnelab "manuele.simi@campagnelab.org"
  install gobyweb artifacts { resource KALLISTO_INDEX 0.42.3 resolved as: KALLISTO_INDEX -> 0.42.3 }
                                INDEX.organism = Homo_sapiens
                                INDEX.reference-build = GRCH38
                                INDEX.ensembl-version-number = 82
}
```

Figure 7. Support for Building Docker Images. NextflowWorkbench offers a composable language to help workflow programmers construct Docker images. In this example, we show a DockerFile root node with a special instruction type called `install gobyweb artifacts`. This special instruction generates RUN commands that install GobyWeb resources in the docker image. Pressing the Build Image button assembles the image. Built images can be used in Processes. Such instructions are useful to create frozen docker images that contain a specific data resource, for instance for clinical analysis workflows which must be frozen as required by regulations.

236 **Advanced Docker Features** The NextflowWorkbench can be used as an interactive development
237 environment for developing docker images. We have developed a composable Dockerfile language.
238 Figure 7 illustrates how docker build files can be written in the workbench to assemble a docker image. In
239 this figure, the first two instructions (`FROM` and `MAINTAINER`) will be familiar to docker programmers
240 who have written or read Dockerfiles. The last instruction, however, is not part of the standard Dockerfile
241 language. This instruction was added with language composition to make it easier to install GobyWeb
242 software or data resources inside the image.

243 DISCUSSION

244 Tools to support the development and execution of workflows have been popular among scientists who
245 need to automate repetitive execution of programs to process different inputs. Galaxy and GenePattern
246 are examples of such tools, which we call graphical workflow systems, designed to help biologists who
247 are not programmers take advantage of bioinformatics programs and automate analyses. The tools offer a
248 graphical metaphor for a workflow where boxes represent tools and lines connect the tools where data
249 feeds from one tool to the other. Graphical workflow systems in wide use today were introduced about 10
250 years ago Giardine et al. [2005], Reich et al. [2006], Hull et al. [2006]. While popularity for these tools
251 grew among computational beginners, experts have yet to widely adopt these systems. While it is unclear
252 what set of reasons explain this lack of interest across the community, understanding these reasons could
253 help design improved systems that both beginners and experts would want to use.

254 Recent work by others, including Di Tommaso et al. [2014], Cingolani et al. [2015], have developed
255 languages to express workflows with scripting or programming languages in an effort to make these tools
256 more useful to computational experts, including bioinformaticians, or biologists with a programming
257 background. In contrast to graphical workflow systems, scripting workflow systems can be installed very
258 quickly and provide strong support for high-performance computation (including support for implicit or
259 explicit parallelization). The focus of these systems is on helping expert bioinformaticians build high
260 performance parallel data analysis workflows. In addition to the workflow language, these systems offer
261 a runtime system to execute the workflow using a range of high-performance grid schedulers (such as
262 Sun Grid Engine, SLURM or Apache Ignite). Because these systems require expressing workflows in
263 text source code for a custom language, learning how to use them requires becoming familiar with the
264 syntax of a new computational language. This is often convenient only for users who have some prior
265 programming and scripting experience. In contrast to these systems, NextflowWorkbench can be thought
266 of as an integrated development environment that offers convenient graphical user interfaces and full
267 language support with auto-completion to guide new users while they learn the language.

268 Kronos is a recent workflow system presented in a pre-print Taghiyar et al. [2016]. Kronos supports
269 writing workflow as structured text files that get compiled into python scripts for execution on a variety of
270 computational platforms (local execution, cluster and cloud). The structured file format makes it possible
271 to define tasks (analogous to Nextflow processes) and subsequently connect these tasks using input output
272 I/O connections. The Nextflow language supported by NextflowWorkbench seems more flexible than

273 simple I/O connections because it is possible to transform data produced by processes with a sequence
274 of functions (e.g., see the use of the `toList()` functions in Figure 6) before the data is provided to
275 a process. A large number of pre-defined functions are supported by Nextflow as well as user-defined
276 closures that can be used to process data in custom ways. This capability is not clearly apparent in the
277 version of the Kronos preprint available as of this writing. Kronos also aims to provide a system that
278 both beginners and experts can use, but does not provide a user interface and integrated development
279 environment. Despite its simplicity, the Kronos language represent a new text-based language whose
280 syntax must be learned by beginners who will find out about errors when running the compiler. In contrast,
281 NextflowWorkbench provides an interactive graphical user interface and advanced interactive features,
282 such as real time error detection, on top of an expressive workflow language: Nextflow.

283 In a broader context, languages to express workflows are essential elements of infrastructures where
284 data analyses are moved to the data. Such infrastructures are becoming necessary in situations where
285 volumes of data are very large (i.e., a tera-byte and more). When several groups need to make computations
286 on the same set of data, transferring the data to the analysis code becomes a bottleneck. In these cases,
287 it is more efficient to move the code to the data than to do the opposite. In the USA, the National
288 Cancer Institute at the National Institutes of Health has started a series of pilots to evaluate this type of
289 infrastructure. The Broad institute, who leads one of these pilot infrastructure projects, has developed
290 the Workflow Description Language and supports executing workflows expressed in WDL on the Broad
291 infrastructure. Another, similar, but incompatible file format to express workflows is the Common
292 Workflow Language (CWL), developed by the Seven Bridges Cancer Genomics Cloud in another NIH/NCI
293 cloud pilot. CWL aims to become a widely used standard to express workflows. To this end, the project
294 organizers are trying to engage a large community of people in the design of CWL. NextflowWorkbench
295 differs from these efforts in several important ways. First, the focus is on user experience to enable
296 beginners to develop and use workflows and to make experts more productive. Neither WDL or CWL
297 address this need. Second, because we used LWT to develop NextflowWorkbench, others can develop
298 extensions of the languages that become immediately integrated with NextflowWorkbench (through
299 language composition and micro-language design, as illustrated in Campagne and Simi [2015]). In our
300 experience, in most cases, there is no need for coordination with our group to develop simple extensions
301 and share them with others. This differs strikingly from a standard development effort, which requires
302 numerous formal communications and coordination before any change can be made to the specification of
303 the “standard”.

304 CONCLUSION

305 In this study, we presented the design and implementation of a workflow system meant for a broad spectrum
306 of potential users, ranging from computational beginners to expert bioinformaticians. We applied language
307 workbench technology to develop such a system, using Nextflow as underlying workflow execution system.
308 We found that we could successfully teach this new workflow system to non-programmers who are able
309 to develop and reuse the simple workflow presented in this manuscript in a short training session of two
310 hours.

311 ACKNOWLEDGMENTS

312 The authors thank Paolo Di Tommaso for assistance integrating changes needed for the Nextflow-
313 Workbench into Nextflow distributions. We thank the members of Cédric Notredame’s Laboratory
314 (http://www.crg.eu/cedric_notredame) for developing and maintaining Nextflow as an
315 open-source project. This investigation was supported by the National Institutes of Health NIAID
316 award 5R01AI107762-02 to Fabien Campagne and by grant UL1 RR024996 (National Institutes of Health
317 (NIH)/National Center for Research Resources) of the Clinical and Translation Science Center at Weill
318 Cornell Medical College. We thank the training session participants who have provided feedback on earlier
319 versions of NextflowWorkbench and help make this new workflow system relevant to computational
320 beginners.

321 REFERENCES

322 V. M. Benson and F. Campagne. Language workbench user interfaces for data analysis. *PeerJ*, 3:e800,
323 2015.

- 324 F. Campagne. *The MPS Language Workbench*, volume I. Fabien Campagne, 2014.
- 325 F. Campagne. *The MPS Language Workbench*, volume II. Fabien Campagne, 2015.
- 326 F. Campagne and M. Simi. *MetaR Documentation Booklet*. Fabien Campagne, 2015.
- 327 F. Campagne, W. E. Digan, and M. Simi. MetaR: simple, high-level languages for data analysis with the
328 R ecosystem. *bioRxiv*, page 030254, 2015. doi: 10.1101/030254. URL [http://biorxiv.org/
329 lookup/doi/10.1101/030254](http://biorxiv.org/lookup/doi/10.1101/030254).
- 330 P. Cingolani, R. Sladek, and M. Blanchette. Bigdatascript: a scripting language for data pipelines.
331 *Bioinformatics*, 31(1):10–16, 2015.
- 332 P. Di Tommaso, M. Chatzou, P. P. Baraja, and C. Notredame. A novel tool for highly scalable computa-
333 tional pipelines. 2014. URL <http://dx.doi.org/10.6084/m9.figshare.1254958>.
- 334 P. Di Tommaso, E. Palumbo, M. Chatzou, P. Prieto, M. L. Heuer, and C. Notredame. The impact of
335 docker containers on the performance of genomic pipelines. *PeerJ*, 3:e1273, 2015.
- 336 K. C. Dorff, N. Chambwe, Z. Zeno, M. Simi, R. Shaknovich, and F. Campagne. GobyWeb: Simplified
337 Management and Analysis of Gene Expression and DNA Methylation Sequencing Data. *PLoS ONE*, 8
338 (7), 2013. ISSN 19326203. doi: 10.1371/journal.pone.0069666.
- 339 B. Giardine, C. Riemer, R. C. Hardison, R. Burhans, L. Elnitski, P. Shah, Y. Zhang, D. Blankenberg,
340 I. Albert, J. Taylor, et al. Galaxy: a platform for interactive large-scale genome analysis. *Genome
341 research*, 15(10):1451–1455, 2005.
- 342 D. Hull, K. Wolstencroft, R. Stevens, C. Goble, M. R. Pocock, P. Li, and T. Oinn. Taverna: a tool for
343 building and running workflows of services. *Nucleic acids research*, 34(suppl 2):W729–W732, 2006.
- 344 S. M. Kurs, Jason P. and F. Campagne. *NextflowWorkbench Documentation Booklet*. Fabien
345 Campagne, 2015. URL [https://play.google.com/store/books/details/Jason_
346 P_Kurs_Nextflow_Workbench_Documentation_Book?id=VQhVCgAAQBAJ](https://play.google.com/store/books/details/Jason_P_Kurs_Nextflow_Workbench_Documentation_Book?id=VQhVCgAAQBAJ).
- 347 M. Reich, T. Liefeld, J. Gould, J. Lerner, P. Tamayo, and J. P. Mesirov. Genepattern 2.0. *Nature genetics*,
348 38(5):500–501, 2006.
- 349 M. Simi and F. Campagne. Composable languages for bioinformatics: the nyosh experiment. *PeerJ*, 2014.
350 URL <https://peerj.com/articles/241/>.
- 351 M. J. Taghiyar, J. Rosner, D. Grewal, B. Grande, R. Aniba, J. Grewal, P. C. Buotros, R. D. Morin,
352 A. Bashashati, and S. Shah. Kronos: a workflow assembler for genome analytics and informatics.
353 Technical report, feb 2016. URL [http://biorxiv.org/content/early/2016/02/19/
354 040352.abstract](http://biorxiv.org/content/early/2016/02/19/040352.abstract).
- 355 M. Voelter. Integrating prose as first-class citizens with models and code. In *MPM@ MoDELS*, pages
356 17–26. Citeseer, 2013.
- 357 M. Voelter, D. Ratiu, B. Schaetz, and B. Kolb. mbeddr: an extensible c-based programming language and
358 ide for embedded systems. In *Proceedings of the 3rd annual conference on Systems, programming, and
359 applications: software for humanity*, pages 121–140. ACM, 2012.