

From reproducible to reusable bioinformatics workflows

GEORGE MARCHMENT, Université Paris-Saclay, CNRS, Laboratoire Interdisciplinaire des Sciences du Numérique, France

BRYAN BRANCOTTE, Institut Pasteur, Université Paris Cité, Bioinformatics and Biostatistics Hub, France

MARIE SCHMIT, Institut Pasteur, Université Paris Cité, Bioinformatics and Biostatistics Hub, France

FRÉDÉRIC LEMOINE, Institut Pasteur, Université Paris Cité, Bioinformatics and Biostatistics Hub, France and Institut Pasteur, Université Paris Cité, CNR Virus Des Infections Respiratoires, France

SARAH COHEN-BOULAKIA, Université Paris-Saclay, CNRS, Laboratoire Interdisciplinaire des Sciences du Numérique, France

Additional Key Words and Phrases: workflows, reusability, structure, nextflow, bioinformatics

ACM Reference Format:

George Marchment, Bryan Brancotte, Marie Schmit, Frédéric Lemoine, and Sarah Cohen-Boulakia. 2018. From reproducible to reusable bioinformatics workflows. 37, 4, Article 111 (August 2018), 2 pages. <https://doi.org/XXXXXXXX.XXXXXXX>

1 ABSTRACT

Data intensive science has ushered in a new era of bioinformatics analyses, resulting in a substantial growth in the scope, complexity and variety of bioinformatics analyses. Over the past decade, workflows management systems, such as Nextflow [3], have become pivotal in the development of these analyses. They are widely used and serve as indispensable tools in the creation, execution, and sharing of complex analyses, in the form of workflows, all while enhancing reproducibility and scalability, among other benefits [4]. Workflows consist of multiple data processing steps chained together by data flow: the input of one step is connected to the output of the previous one, determining their execution order. Hence, a workflow can be represented as a directed graph (see Fig. 1), called the *workflow structure*.

If workflow management systems have solved many issues related to the reproducibility of bioinformatics analyses, their usefulness in terms of a workflow's reusability is still limited. Indeed despite the growing number of bioinformatics workflows developed, there exists a major lack of their reuse among different users. [4] showed that for all Nextflow workflows, the top-30 most reused steps are only reused in 2.42% of the total steps.

This lack of reuse may stem from the inherent complexity of workflow code, particularly challenging for biologists and bioinformaticians who lack expertise in the field [1, 5].

Authors' Contact Information: George Marchment, george.marchment@universite-paris-saclay.fr, Université Paris-Saclay, CNRS, Laboratoire Interdisciplinaire des Sciences du Numérique, Orsay, France; Bryan Brancotte, Institut Pasteur, Université Paris Cité, Bioinformatics and Biostatistics Hub, Paris, France; Marie Schmit, Institut Pasteur, Université Paris Cité, Bioinformatics and Biostatistics Hub, Paris, France; Frédéric Lemoine, frederic.lemoine@pasteur.fr, Institut Pasteur, Université Paris Cité, Bioinformatics and Biostatistics Hub, Paris, France and Institut Pasteur, Université Paris Cité, CNR Virus Des Infections Respiratoires, Paris, France; Sarah Cohen-Boulakia, sarah.cohen-boulakia@universite-paris-saclay.fr, Université Paris-Saclay, CNRS, Laboratoire Interdisciplinaire des Sciences du Numérique, Orsay, France.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM 2476-1249/2018/8-ART111

<https://doi.org/XXXXXXXX.XXXXXXX>

