

COMMENTARIES

Parallel sequencing lives, or what makes large sequencing projects successful

Javier Quilez^{1,2,*}, Enrique Vidal^{1,2}, François Le Dily^{1,2}, François Serra^{1,2,3}, Yasmina Cuartero^{1,2,3}, Ralph Stadhouders^{1,2}, Thomas Graf^{1,2}, Marc A. Marti-Renom^{1,2,3,4}, Miguel Beato^{1,2} and Guillaume Filion^{1,2}

¹Gene Regulation, Stem Cells and Cancer Program, Centre for Genomic Regulation (CRG), The Barcelona Institute of Science and Technology (BIST), Dr. Aiguader 88, 08003 Barcelona, Spain, ²Universitat Pompeu Fabra (UPF), Barcelona, Spain, ³CNAG-CRG, Centre for Genomic Regulation (CRG), The Barcelona Institute of Science and Technology (BIST), Dr. Aiguader 88, 08003 Barcelona, Spain and ⁴ICREA, Pg. Lluís Companys 23, 08010 Barcelona, Spain

*Correspondence address. Javier Quilez, Gene Gene Regulation, Stem Cells and Cancer Programme, CRG-Centre for Genomic Regulation, C/ Dr. Aiguader, 88, PRBB Building, 08003 Barcelona, Spain. Tel: +34 93 316 01 15; E-mail: javier.quilez@crg.eu

Abstract

T47D_rep2 and b1913e6c1.51720e9cf were 2 Hi-C samples. They were born and processed at the same time, yet their fates were very different. The life of b1913e6c1.51720e9cf was simple and fruitful, while that of T47D_rep2 was full of accidents and sorrow. At the heart of these differences lies the fact that b1913e6c1.51720e9cf was born under a lab culture of Documentation, Automation, Traceability, and Autonomy and compliance with the FAIR Principles. Their lives are a lesson for those who wish to embark on the journey of managing high-throughput sequencing data.

Keywords: high-throughput sequencing; management and analysis best practices; bioinformatics; FAIR Principles

The Beginning

Linda worked hard to produce a Hi-C sample in T47D cells. Upon submitting the sample for sequencing, she remembered the motto of the lab: “Make DATA more FAIR (Findable, Interoperable, Accessible, Reusable).” The team had established lab-wide habits of Documentation, Automation, Traceability, and Autonomy of experimenters. The experienced people insisted that human interfaces are always the weak link. “Every time a project fails, someone is typing on a keyboard. . . or does not bother to.” The metadata must be accurate, the code must be readable, the data must be tidy [1]. Technology helps, but this is mostly a matter of attitude. Not only had this attitude improved the perfor-

mance of the lab, but it also paved the way to meet international quality standards as those defined by the FAIR Principles [2].

Linda filled in the metadata on a low-key online Google Form. The lab had chosen this option among many others because experimenters found it the easiest. Filling the form was quick: they had to click on items from drop-down lists. As she pressed “Submit”, a shared Google Sheet was immediately updated and she received the name b1913e6c1.51720e9cf that uniquely identified her sample. These unnatural names had first left her skeptical, but she could now see the benefits of that system to collect the metadata and trace sequencing samples. She remembered the meetings with the bioinformaticians in an attempt to make the

Received: 1 August 2017; Revised: 12 September 2017; Accepted: 8 October 2017

© The Author 2017. Published by Oxford University Press. This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

data more FAIR [2]. “A project is as good as its metadata; you will see the benefit only after a year or two”, they kept saying.

Meanwhile in another lab, Pedro also worked hard to produce a Hi-C sample in T47D cells. He proudly wrote “T47D.rep2” on the tube and gave it to the sequencing facility. All the information he considered relevant was in his notebook.

By a strange coincidence, both Linda and Pedro soon found new positions. They left their respective institutes without finishing their project.

Life After Turnover

Simon was the bioinformatician in charge of analyzing T47D.rep2. He was not happy that Pedro left the institute, because he had questions about the sample. As he went to save the files in the shared repository, he saw that there were already four samples called “T47D.rep2” in different directories. Simon face-palmed and headed for the wet lab. Fortunately, Janet knew something about it: “Some of these are my experiments with glucose deprivation; the others are Pedro’s. Despite the modest sequencing coverage, he found interesting changes in the genome structure when treating with hormone, so he repeated the experiments to obtain higher coverage.” Looking into Pedro’s notes, Simon saw that, indeed, the number of reads of the initial sample was very low, hence the newest sample “T47D.rep2.” At long last, Simon had an idea of what “T47D.rep2” was.

Meanwhile, Paul, the bioinformatician in charge of analyzing b1913e6c1.51720e9cf, pulled the record from the database where the metadata in the Google Sheet were automatically dumped. The online spreadsheet was a convenient frontend for the experimenters, but the database offered a more programmatic access to the metadata—plus it was an additional backup layer. On his end, Paul launched the mapping pipeline and performed several downstream analyses that Chloe requested. He documented the procedure in the Jupyter electronic notebook he created for the analysis. The production code was run in Docker containers and pushed to a GitHub repository. The notebooks helped him (or anyone else) keep track of the analyses in a readable format, while Docker virtual machines allowed him (or anyone else) to run the code on different machines without the hassle of installing countless libraries. Finally, GitHub was as much a backup as a way to share his work.

Chloe examined the results in the online report she received from Paul and performed some additional analyses with an R Shiny web application to inspect the Hi-C data processed in the lab. It had taken some time to implement it, but now the benefits were clear: Paul could focus on other things than running basic analyses for all the lab members, and meanwhile, they were more autonomous. This last analysis provided further evidence supporting their hypothesis, so Chloe was ready to polish their manuscript. Each analysis performed by Paul was allocated in a directory with a traceable name and a clear content structure, and they were permanently accessible in the FTP site of the lab. Therefore, Chloe knew where to find the figures and tables that she needed, she updated the Methods section with the information written in the report, and she was even able to provide the scripts and parameter values used in the analysis as a GitHub repository—she knew that editors were getting more and more serious about reproducibility.

The Reviews

Chloe was very happy to hear their manuscript received positive comments from the reviewers. The only obstacle to publication

seemed to be Reviewer #3, who asked them to replicate the findings in an independent, larger data set that had been recently published. Tough but fair. Chloe panicked about having to analyze almost 100 samples in so little time; during the project, they had generated a smaller number of samples and analyzed them over time, so she worried that it would take too long. Paul reassured her: all she had to do was prepare the metadata for the new data set, as Linda had done for b1913e6c1.51720e9cf. Then, a simple command would execute the pipeline for the ~100 samples as effortlessly as for a single one, and all the required information would be retrieved automatically from the database of metadata. Running the pipeline could be parallelized in the multiple cores available in the computing cluster of the institute, so all samples were processed within a few days. In the meantime, he would start preparing the submission of the data to a public repository: a simple search within the structured directories allocated for the FASTQ and the contact matrix files and a selection of entries from the database of metadata would do much of the work. Lastly, Paul checked that the manuscript complied with the “Minimum Standards of Reporting Checklist” of the target journal and the FAIR Principles [2]. Findability and accessibility: the data and metadata were linked by the unique sample identifier and uploaded to GEO, the code was pushed to GitHub, and the URLs of both repositories were available in the manuscript. Interoperability: the Docker containers used to run the pipelines were pushed to Docker Hub. Reusability: the metadata was complete, and the data procedures were well documented.

Meanwhile, Simon was far from publication. Overall, the preliminary results of Pedro were not confirmed in the new high-coverage samples. He knew too well that trouble was only starting. Simon scavenged the directories looking for the code used to generate the plots he had seen, those that indicated a clear effect of hormone treatment on the genome structure. Unfortunately, the workflow of the analysis and the specific parameter values were not documented. Perhaps his predecessors had forgotten to remove polymerase chain reaction duplicates? And how did they correct for multiple testing, if at all? After guessing where to find the older raw data, Simon processed the initial data set with his analysis pipeline, but the differences between the old and new data sets remained. Suddenly, a thought froze Simon: ‘messy data organization, duplicated names, no metadata... could it be that the promising results Pedro found were mistakenly derived from Janet’s glucose-deprived samples?’ Simon face-palmed as if guessing the answer.

Behind the Scenes

The human factor is the greatest hurdle to reaching the standard of the FAIR Principles [2]. People change their minds, they resist change, they follow their own rules, and they plan for the short term. As an insurance against fiasco (Table 1), a scientific team must develop habits and tools for sharing data and analyses. The main idea is to limit or control human intervention by automating every step.

1. Achieving this has a cost. The most significant is the time spent in maintenance and continuous improvement of the software. As a rule of thumb, this toll should not exceed one-fifth of the productivity. We recommend choosing the hardware platform that minimizes the overall training time of the users and the maintenance time of the developers. Likewise, maintenance and training are important aspects of choosing the software for a project.

Table 1: Challenges associated with the management and analysis of high-throughput sequencing data

Challenge	Impact	Consideration
Poor sample description	<ul style="list-style-type: none"> • Prevents data processing and quality control • Incorrect analysis and results • Lack of reproducibility • Delays publication 	Metadata collection
Unsystematic sample naming	<ul style="list-style-type: none"> • Duplicated or similar names • Ambiguous identification • Precludes computational treatment • Data disclosure 	Sample identifier scheme
Untidy data organization	<ul style="list-style-type: none"> • Data cannot be found • Time consumption • Inability to automate searches 	Structured and hierarchical data organization
Yet another analysis	<ul style="list-style-type: none"> • Repeated manual execution of analyses • Inability to deconvolute analysis, producing different results • Compulsory linear execution 	Scalability, parallelization, automatic configuration, and modularity
Undocumented procedures	<ul style="list-style-type: none"> • Poor understanding of results • Irreproducibility • Hampers catching errors 	Documentation
Data overflow	<ul style="list-style-type: none"> • No access to data • Size and number of files make individual inspection inefficient 	Interactive web applications

What did go wrong with the T47D_rep2 sample? Its description and metadata were not collected, digitized, and stored in a central repository; orphan of a sample identification scheme, it received a duplicated name; why this and previous samples were generated, where their data were located, and how these were processed and with which methods were not documented. As storified with the lives of b1913e6c1.51720e9cf and T47D_rep2, managing and analyzing the growing amount of sequencing data present several challenges. This table details their impact on scientific quality and proposes considerations to address them.

- The absolute priority is metadata collection. We propose a scheme for collection and file naming (Fig. 1a and Additional file 1), but any system will do, as long as it is (i) agreed upon and understood by people using it, (ii) backed up automatically, (iii) future-aware and flexible, and (iv) there is someone responsible for maintenance and validation of the metadata.
- The second priority is to locate the data and the analyses. We propose a hierarchical organization that can evolve according to future needs (Fig. 1b). Again, any scheme with the properties above will do.
- Next, the analyses must be documented. Here a flurry of tools help the analysts keep track of and organize their work as it unfolds. The most popular are Jupyter Notebook and R Studio. Here we recommend using widely accepted tool kits as this facilitates sharing between the members of the team and the rest of the world.
- Such tools partly address the next priority, which is reproducibility. However, today we can go one step further with virtual machines. Providing the environment to run the analyses makes it easier to reproduce them and to share code. In this area, Docker has taken the lead, but the platform evolves too fast to be reliable in the long term. No executable is guaranteed to run at the 10-year horizon, so archiving the source files (including Docker files) on version controls systems such as GitHub is essential.
- Finally, experimenters should be empowered to perform basic analyses. The most efficient teams are made of specialists, so researchers should do what they are expert at (or become expert at what they do). But bioinformatics is fast becoming “common knowledge.” Building interfaces for standard analyses is a way to free bioinformaticians to focus

on the most technical parts of the project, while allowing all the members to contribute to the analyses. Many modern tools such as R Shiny can help build such interfaces. Here, the most important is that the developer be proficient with the chosen tool and that the users understand how to use the interface.

Data accumulates at a rapid pace in life sciences (Additional file 2), and stories similar to that of b1913e6c1.51720e9cf and T47D_rep2 have taken place in many research groups (Additional files 3–5). We propose that data-producing teams focus on Documentation, Automation, Traceability, and Autonomy (DATA) as main priorities, with the purpose of being “human-proof.” The scheme implemented in our own projects is shown in Figures 1 and 2, the tools are listed in Table 2, and in Additional file 6 we comment on the considerations and costs of implementing and maintaining a DATA lab culture. To illustrate our recommendations, we also provide a didactic data set (the actual sample b1913e6c1.51720e9cf) via GitHub [10].

Availability of data and material

The didactic data set is available via GitHub [10].

Abbreviations

3K RGP: 3000 Rice Genomes Project; DATA: Documentation, Automation, Traceability, and Autonomy; ENCODE: Encyclopedia of DNA Elements; HTS: high-throughput sequencing; ID: identifier; SRA: Short Read Archive; SQL: Structured Query Language; TCGA: The Cancer Genome Atlas.

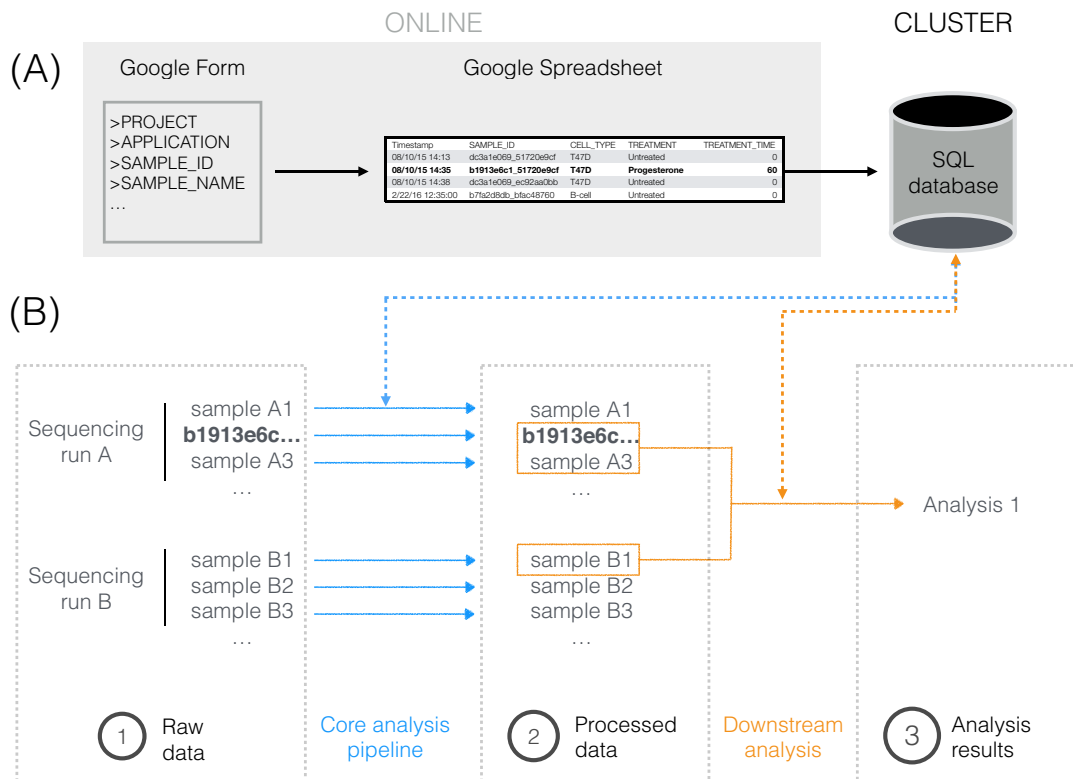


Figure 1: A traceable life for b1913e6c1.51720e9cf. (a) The metadata for b1913e6c1.51720e9cf were collected via an online Google Form and stored both online (Google Sheet) and in a local SQL database. A good metadata collection system should be (i) short and easy to complete, (ii) instantly accessible by authorized users, and (iii) easy to parse for humans and computers. (b) b1913e6c1.51720e9cf was sequenced along with other samples, whose raw sequencing data were located in a directory named after the date of the sequencing run. There one could find the FASTQ files containing the sequencing reads from b1913e6c1.51720e9cf as well as information about their quality; no modified, subsetted, or merged FASTQ file was stored to ensure that analyses started off from the very same set of reads. In a first step, the raw data of b1913e6c1.51720e9cf were processed with the Hi-C analysis pipeline, which created a “b1913e6c1.51720e9cf” directory at the same level where all processed Hi-C samples were located. “b1913e6c1.51720e9cf” had multiple subdirectories that stored the files generated in each of the steps of the pipeline, the logs of the programs, and the integrity verifications of key files. Moreover, such subdirectories accounted for variations in the analysis pipelines (e.g., genome assembly version, aligner) so that data were not overwritten. In a second step, processed data from b1913e6c1.51720e9cf and other samples were used to perform the downstream analyses Chloe asked Paul. Within the directory he allocated to her analyses, Paul created a new one called “2017-03-08.hic.validation” with the description of the analysis, along with the scripts used and the tables and figures generated.

Additional files

Additional file 1. (a) More than reads. FASTQ files may be useless if not coupled with biological, technical, and logistics information (metadata). Metadata are used at several stages of the high-throughput sequencing data. In the initial processing, for instance, the human origin of b1913e6c1.51720e9cf was needed to determine hg38 as the reference genome sequence to which reads would be aligned, and the restriction enzyme “DpnII” applied in the Hi-C protocol was used in the mapping too. Other metadata were used for quality control (e.g., sequencing facility and/or date for detecting batch effects or rescuing swapped samples using the correct index) or in the downstream analysis (e.g., cell type, treatment). Furthermore, metadata are critical for data sharing and reproducibility. (b) Choosing a name. Long before b1913e6c1.51720e9cf was generated, a scheme to name Hi-C samples was envisioned. First, 2 sets of either biological or technical fields that unequivocally defined a sequencing sample were identified. Then, for a given sample, the values of the biological fields treated as text are concatenated and computationally digested into a 9-mer, and the same procedure is applied to the technical fields. The two 9-mers are combined to form the sample identifier (ID), as happened for b1913e6c1.51720e9cf. Despite the apparent non-informativeness of this sample ID ap-

proach, it easily allows identifying biological replicates and samples generated in the same batch since they will share, respectively, the first and second 9-mer. While the specific fields used to generate the sample ID can vary, it is important that they unambiguously define a sequencing sample (otherwise duplicated identifiers can emerge) and that they are always combined in the same order to ensure reproducibility. Indeed, another advantage of this naming scheme is that the integrity of the metadata can be checked, as altered metadata values will lead to a different sample ID.

Additional file 2. Rapid accumulation and diversity of high-throughput sequencing (HTS) data. The past decade has witnessed a tremendous increase in sequencing throughput and applications, causing uncontrolled accumulation of sequencing data sets. (a) For instance, the number of sequences deposited in the Sequence Read Archive (SRA) [3], a major repository for HTS data, has skyrocketed from ~2 terabases in 2009 to ~9000 terabases (the size of approximately 3 million human genomes) at the beginning of 2017. Moreover, this is surely an underestimation of the actual amount given that only sequencing experiments eventually included in a publication are deposited. Although data-intensive projects like TCGA [4], the 1000 Genomes Project [5], ENCODE [6], and 3K RGP [7] are top HTS data generators [8], such a boost in the number of existing sequences

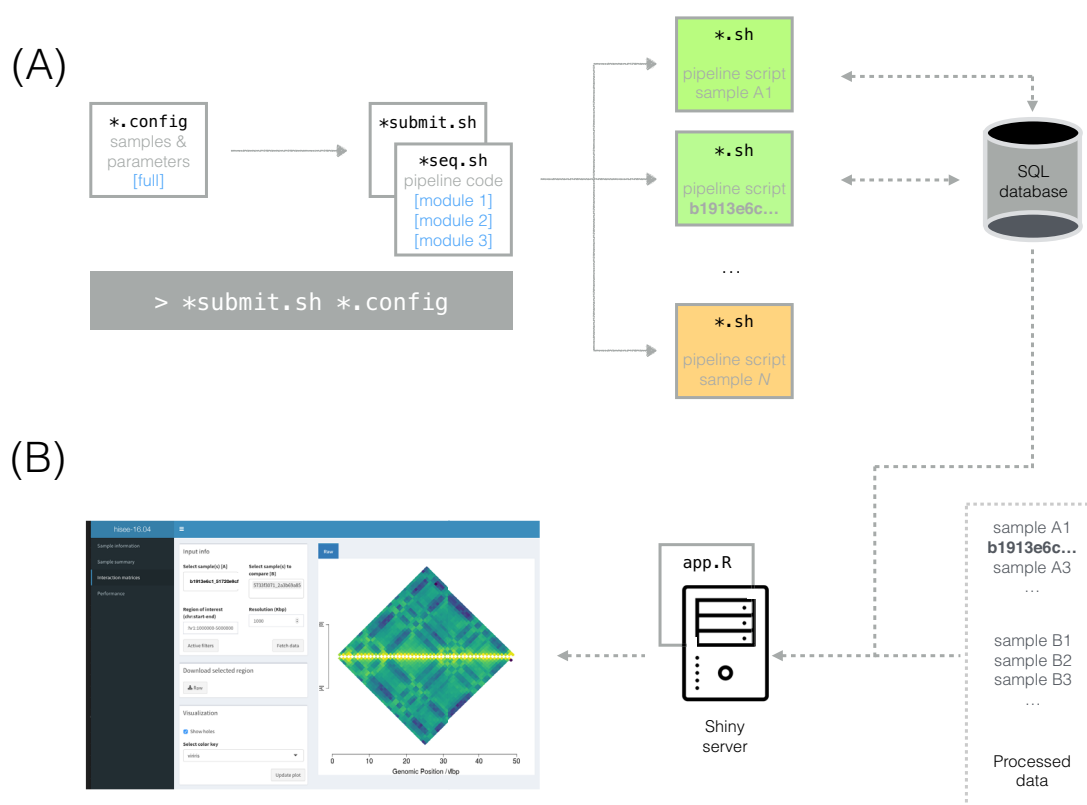


Figure 2: Automating the analysis and visualization of b1913e6c1.51720e9cf data. (a) Scalability, parallelization, automatic configuration, and modularity of analysis pipelines. Paul launched the Hi-C pipeline for hundreds of samples with a single command (gray rectangle): the submission script (“*.submit.sh”) generated as many pipeline scripts as samples listed in the configuration file (“*.config”). The configuration file also contained the hard-coded parameters shared by all samples, such as the maximum running time Paul underestimated for some samples. Processing hundreds of samples was relatively fast because (i) the pipeline script for each of the samples was submitted as an independent job in the computing cluster, where it was queued (orange) and eventually executed in parallel (green), and (ii) the pipeline code in “*.seq.sh” was adapted for running in multiple processors. For further automation, each process retrieved sample-specific information (e.g., species, read length) from the metadata SQL database; in addition, metrics generated by the pipeline (e.g., running time, number of aligned reads) were recorded into the database. Because the pipeline code was grouped into modules, Paul was able to easily re-run the “generate_matrix” module for those samples that failed in his first attempt. (b) Interactive web application to visualize Hi-C data. b1913e6c1.51720e9cf alone generated ~70 files of plots and text when passed through the Hi-C pipeline. Inspecting them might have seemed a daunting task for Chloe: she did not feel comfortable navigating the cluster and lacked the skills to manipulate them anyway, and even if she did, examining so many files for dozens of samples seemed endless. Luckily for her, Paul had developed an interactive web application with R Shiny (Table 2) that allowed her to visualize data and metadata and perform specific analyses in a user-friendly manner.

Table 2: Tools used in the story

Tool	Usage	Website
Docker	Interoperability	https://www.docker.com/
Docker Hub	Repository for Docker containers	https://hub.docker.com/
GEO	Repository for high-throughput genomics data	https://www.ncbi.nlm.nih.gov/geo/
GitHub	Version control and backup of code	https://github.com/
Google Forms and Sheets	Online collection and display of metadata	https://www.google.com/forms/about/
Jupyter Notebook	Document procedures and perform analysis	http://jupyter.org/
R Shiny	Deploy web applications	https://shiny.rstudio.com/
R Studio	Document procedures and perform analysis	https://www.rstudio.com/

Note that Jupyter Notebook and R Studio environments are not good for analyses that run for a long time and/or require heavy computational power. Therefore, we recommend them as a way to document how data are processed (even if long/heavy analyses are executed elsewhere) and to perform downstream analyses (e.g., summarizing, plotting) after the long-running ones are done.

reflects a pervasive use of HTS. (b) As an example, while sequencing data for >90 000 studies have been submitted to the SRA, the top 10 and 100 contributors in terms of number of bases represent only a part of the archive (~30% and ~60%, respectively). (c) Similarly, while ~80% of SRA data derive from *Homo sapiens* and *Mus musculus*, the central organisms in large sequencing projects, the remaining 20% come from a diverse number of organisms (~50 000). Data were obtained from the Short

Read Archive [9] and processed as described in the didactic data set [10].

Additional file 3. Why T47D_rep2 and b1913e6c1.51720e9cf are not singletons.

Additional file 4. Number of SRA deposited bases grouped by instrument name. Data were obtained from the Short Read Archive[9] and processed as described in the didactic data set [10].

Additional file 5. Number of SRA deposited bases grouped by the submitter. For the top 25 contributors in terms of number of bases submitted, we searched for instances of multiple entries probably referring to the same submitter (e.g., “ncbi” and “NCBI”). Data were obtained from the Short Read Archive [9] and processed as described in the didactic data set [10].

Additional file 6. Considerations and costs of implementing and maintaining a Documentation, Automation, Traceability, and Autonomy lab culture.

Competing interests

The authors declare that they have no competing interests.

Funding

We received funding from the European Research Council under the European Union’s Seventh Framework Programme (FP7/2007–2013)/ERC Synergy grant agreement 609 989 (4DGenome). The content of this manuscript reflects only the authors’ views, and the Union is not liable for any use that may be made of the information contained therein. We acknowledge support of the Spanish Ministry of Economy and Competitiveness, “Centro de Excelencia Severo Ochoa 2013–2017,” and Plan Nacional (SAF2016–75 006-P), as well as support of the CERCA Programme/Generalitat de Catalunya. R.S. was supported by an EMBO Long-term Fellowship (ALTF 1201–2014) and a Marie Curie Individual Fellowship (H2020-MSCA-IF-2014).

Author contributions

Conceptualization: J.Q., G.F.; data curation: J.Q.; formal analysis: J.Q.; funding acquisition: T.G., M.A.M.R., M.B., G.F.; methodology: J.Q., E.V., F.D., Y.C., R.S.; software: J.Q., E.V., F.S.; visualization: J.Q., E.V.; writing—original draft: J.Q., G.F.; writing—review and editing: E.V., F.D., F.S., Y.C., R.S., T.G., M.A.M.R., M.B. All authors read and approved the final manuscript.

Acknowledgements

We thank F. Javier Carmona and Corey T. Watson for advice on the manuscript and Henning Hermjakob and Titus C. Brown for constructive reviews.

References

1. Kenall A, Edmunds S, Goodman L et al. Better reporting for better research: a checklist for reproducibility. *Gigascience* 2015;**4**:32.
2. Wilkinson MD, Dumontier M, Aalbersberg IJ et al. The FAIR guiding principles for scientific data management and stewardship. *Sci Data* 2016;**3**:160018.
3. Leinonen R, Sugawara H, Shumway M. International Nucleotide Sequence Database Collaboration. The sequence read archive. *Nucleic Acids Res* 2011;**39**(database): D19–21.
4. Cancer Genome Atlas Research Network, Weinstein JN, Collisson EA, Mills GB et al. The cancer genome atlas pan-cancer analysis project. *Nat Genet* 2013;**45**:1113–20.
5. 1000 Genomes Project Consortium, Abecasis GR, Auton A, Brooks LD et al. An integrated map of genetic variation from 1,092 human genomes. *Nature* 2012;**491**:56–65.
6. ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human genome. *Nature* 2012;**489**:57–74.
7. Li J, Wang J, Zeigler RS. The 3,000 rice genomes project: new opportunities and challenges for future rice research. *Gigascience* 2014;**3**:8. doi:10.1186/2047-217X-3-8.
8. Muir P, Li S, Lou S et al. The real cost of sequencing: scaling computation to keep pace with data generation. *Genome Biol* 2016;**17**:53. doi:10.1186/s13059-016-0917-0.
9. Short Read Archive. Short read archive. <https://trace.ncbi.nlm.nih.gov/Traces/sra/sra.cgi>. Accessed 6 April 2017.
10. Didactic dataset, “Parallel Sequencing Lives.” GitHub 2017. <https://github.com/4DGenome/parallel.sequencing.lives>.