# 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]

*Corresponding author


Institutional addresses:

[1]Gene 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

[2]Universitat Pompeu Fabra (UPF), Barcelona, Spain

[3]CNAG-CRG, Centre for Genomic Regulation (CRG), The Barcelona Institute of Science and Technology (BIST), Dr. Aiguader 88, 08003 Barcelona, Spain

[4]ICREA, Pg. Lluis Companys 23, 08010 Barcelona, Spain

Email addresses:

JQ: javier.quilez@crg.eu

EV: enrique.vidal@crg.eu

FD: francois.ledily@crg.eu

FS: francois.serra@cnag.crg.eu

YC: yasmina.cuartero@cnag.crg.eu

RS: ralph.stadhouders@crg.eu

TG: thomas.graf@crg.eu

MAM-R: martirenom@cnag.crg.eu

MB: miguel.beato@crg.eu

GF: guillaume.filion@crg.eu

## Abstract

T47D_rep2 and b1913e6c1_51720e9cf were two 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, 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". The team had established lab-wide habits of Documentation, Automation, Traceability and Autonomy of experimenters. The old-timers 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. Technology helps, but this is mostly a matter of attitude. Not only had such attitude improved the performance of the lab but it also paved the way to meet international quality standards as those defined by the FAIR Principles [1].

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

53    attempt to make the data more FAIR [1]. "A project is as good as its metadata; you will see the benefit

54    only after a year or two" they kept telling.

55    Meanwhile in another lab, Pedro also worked hard to produce a Hi-C sample in T47D cells. Things

56    had gone wrong in the past, but this time all the quality controls looked good. He proudly wrote

57    "T47D_rep2" on the tube and gave it to the sequencing facility. All the information *he* considered

58    relevant was in his notebook.

59    By a strange coincidence, both Linda and Pedro soon found a new position. They left their respective

60    institutes without finishing their project.

61

## 62    Life after turn-over

63    Simon was the bioinformatician in charge of analyzing T47D_rep2. He was not happy that Pedro left the

64    institute, because he had questions about the sample. As he meant to save the files in the shared

65    repository, he realized that there were already four samples called "T47D_rep2" in different directories.

66    Simon facepalmed and headed for the wet lab. Fortunately, Janet knew something about it: "Some of

67    these are my experiments; the others are Pedro's. Despite the modest sequencing coverage, he found

68    interesting changes in the genome structure when treating with hormone, so he repeated the

69    experiments to obtain higher coverage". Looking into Pedro's notes, Simon saw that indeed the

70    sequencing quality of the raw reads was very poor, hence the newest sample "T47D_rep2". At long last,

71    Simon had an idea of what "T47D_rep2" was…

72    Meanwhile, Paul, the bioinformatician in charge of analyzing b1913e6c1_51720e9cf pulled the record

73    from the database where the metadata in the Google Sheet were automatically dumped. The online

74    spreadsheet was a convenient frontend for the experimenters, but the database offered a more

75    programmatic access to the metadata — plus it was an additional backup layer. On his end, Paul

76    launched the mapping pipeline and performed several downstream analyses that Chloe requested. He

77    documented the procedure in the Jupyter electronic notebook he created for the analysis. The

78    production code was run in Docker containers and pushed to a GitHub repository. The notebooks

79    helped him (or anyone else) keep track of the analyses in a readable format, while Docker virtual

80  machines allowed him (or anyone else) to run the code on different machines without the hassle of

81  installing countless libraries. Finally, GitHub was as much a backup as a way to share his work.

82  Chloe examined the results in the online report she received from Paul and performed some

83  additional analyses with an R Shiny web application to inspect the Hi-C data processed in the lab. It had

84  taken some time to implement it, but now the benefits were clear: Paul could focus on other things than

85  running basic analyses for all the lab members and, meanwhile they were more autonomous. This last

86  analysis provided further evidence supporting their hypothesis, so Chloe was ready to polish their

87  manuscript. Each analysis performed by Paul was allocated in a directory with a traceable name, a clear

88  content structure and permanently accessible in the FTP site of the lab. Therefore, Chloe knew where to

89  find the figures and tables that she needed, updated the Methods section with the information written in

90  the report and she was even able to provide the scripts and parameter values used in the analysis as a

91  GitHub repository — she knew that editors were getting more and more serious about reproducibility.

92

## 93  The reviews

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

95  obstacle to publication seemed to be Reviewer #3, who asked to replicate the findings in an

96  independent larger dataset that had been recently published. Tough but fair. Chloe panicked about

97  having to analyze almost 100 samples in so little time; during the project they had generated a smaller

98  number of samples and analyzed them over time, so she worried that it would take too long. Paul

99  reassured her: all she had to do was prepare the metadata for the new dataset, as Linda had done for

100  b1913e6c1_51720e9cf. Then, a simple command would execute the pipeline for the ~100 samples as

101  effortlessly as for a single one, and all the required information would be retrieved automatically from the

102  database of metadata. Running the pipeline could be parallelized in the multiple cores available in the

103  computing cluster of the institute, so all samples were processed within a few days. In the meantime, he

104  would start preparing the submission of the data to a public repository: a simple search within the

105  structured directories allocated for the FASTQ and the contact matrix files as well as a selection of

106  entries from the database of metadata would do much of the work. Lastly, Paul checked that the

4

107  manuscript complied with the FAIR Principles [1]. Findability and accessibility: the data and metadata

108  were linked by the unique sample identifier and uploaded to GEO, the code was pushed to GitHub and

109  the URL to both repositories available in the manuscript. Interoperability: the Docker containers used to

110  run the pipelines were pushed to Docker Hub. Reusability: the metadata was complete and the data

111  procedures were well documented.

112  Meanwhile, Simon was far from publication. Overall, the preliminary results of Pedro were not

113  confirmed in the new high-coverage samples. Simon scavenged the directories looking for the code

114  used to generate the plots he had seen, those that indicated a clear effect of hormone treatment on the

115  genome structure. Unfortunately, the workflow of the analysis and the specific parameter values were

116  not documented. Perhaps his predecessors had forgotten to remove PCR duplicates? And how did they

117  correct for multiple testing, if at all? After guessing where to find the older raw data, Simon processed

118  the initial dataset with his analysis pipeline but the differences between the old and new datasets

119  remained. Simon facepalmed. He knew too well that trouble was only starting...

120

## Behind the scene

122  The human factor is the greatest hurdle to reaching the standard of the FAIR Principles [1]. People

123  change their mind, they resist change, they follow their own rules and they plan for the short term. As an

124  insurance against fiasco (**Table 1**), a scientific team must develop habits and tools for sharing data and

125  analyses. The main idea is to limit or control human intervention by automating every step.

126  1.  The absolute priority is metadata collection. We propose a scheme for collection and file naming

127  (**Figure 1a** and **Additional file 1**), but any system will do, as long as it is (i) agreed upon and

128  understood by people using it, (ii) backed up automatically, (iii) future-proof and (iv) there is someone

129  responsible for maintenance and validation of the metadata.

130  2.  The second priority is to locate the data and the analyses. We propose a hierarchical organization

131  that can evolve according to future needs (**Figure 1b**). Again, any scheme with the properties above

132  will do.

133   3.   Next, the analyses must be documented. Here a flurry of tools help the analysts keep track of and

134        organize their work as it unfolds. The most popular are Jupyter for Python and Rstudio for R. Here

135        we recommend using widely accepted tool kits as this facilitates sharing between the members of

136        the team and the rest of the world.

137   4.   Such tools partly address the next priority, which is reproducibility. However, today we can go one

138        step further with virtual machines. In this area, Docker has taken the lead and we recommend

139        developing ground up production scripts and exploratory analyses in Docker containers.

140   5.   Finally, experimenters should be empowered to perform basic analyses. The most efficient teams

141        are made of specialists, so researchers should do what they are expert at (or become expert at what

142        they do). But bioinformatics is fast becoming "common knowledge". Building interfaces for standard

143        analyses is a way to free bioinformaticians to focus on the most technical parts of the project, while

144        allowing all the members to contribute to the analyses. Many modern tools such as R Shiny can help

145        build such interfaces. Here, the most important is that the developer be proficient with the chosen

146        tool, and that they users understand how to use the interface.

147   Data accumulates at a rapid pace in life sciences (**Additional file 2**), and stories similar to that of

148   b1913e6c1_51720e9cf and T47D_rep2 have taken place in many research groups (**Additional files 3-5**).

149   We propose that data-producing teams focus on Documentation, Automation, Traceability and

150   Autonomy as main priorities, with the purpose of being "human-proof". The scheme implemented in our

151   own projects is shown in **Figures 1-2**, and the tools are listed in **Table 2**. To illustrate our

152   recommendations, we also provide a didactic data set (the actual sample b1913e6c1_51720e9cf) at the

153   following link: https://github.com/4DGenome/parallel_sequencing_lives.

154

# Abbreviations

156   3K RGP: 3,000 Rice Genomes Project; ENCODE: Encyclopedia of DNA Elements; HTS: high-throughput

157   sequencing; ID: identifier; SRA: Short Read Archive; SQL: Structured Query Language; TCGA: The

158   Cancer Genome Atlas.

159

## Declarations

**Ethics approval and consent to participate**

Not applicable.

**Consent for publication**

Not applicable.

**Availability of data and material**

The didactic dataset is available at https://github.com/4DGenome/parallel_sequencing_lives

**Competing interests**

The authors declare that they have no competing interests.

**Author's contributions**

Conceptualization: JQ, GF; Data curation: JQ; Formal analysis: JQ; Funding acquisition: TG, MAM-R, MB, GF; Methodology: JQ, EV, FD, YC, RS; Software: JQ, EV, FS; Visualisation: JQ, EV; Writing - original draft: JQ, GF; Writing - review & editing: EV, FD, FS , YC, RS, TG, MAM-R, MB. All authors read and approved the final manuscript.

185

## References

1. Wilkinson MD, Dumontier M, Aalbersberg IJJ, Appleton G, Axton M, Baak A, et al. The FAIR Guiding Principles for scientific data management and stewardship. Sci Data. 2016;3:160018. doi:10.1038/sdata.2016.18.

2. Leinonen R, Sugawara H, Shumway M, International Nucleotide Sequence Database Collaboration. The sequence read archive. Nucleic Acids Res. 2011;39 Database issue:D19-21. doi:10.1093/nar/gkq1019.

3. Cancer Genome Atlas Research Network, Weinstein JN, Collisson EA, Mills GB, Shaw KRM, Ozenberger BA, et al. The Cancer Genome Atlas Pan-Cancer analysis project. Nat Genet. 2013;45:1113–20. doi:10.1038/ng.2764.

4. 1000 Genomes Project Consortium, Abecasis GR, Auton A, Brooks LD, DePristo MA, Durbin RM, et al. An integrated map of genetic variation from 1,092 human genomes. Nature. 2012;491:56–65. doi:10.1038/nature11632.

5. ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human genome. Nature. 2012;489:57–74. doi:10.1038/nature11247.

6. Li J-Y, 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.

7. Muir P, Li S, Lou S, Wang D, Spakowicz DJ, Salichos L, 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.

8. Short Read Archive. Short Read Archive. https://trace.ncbi.nlm.nih.gov/Traces/sra/sra.cgi. Accessed 6 Apr 2017.

## Figures

209

**Figure 1. A traceable life for b1913e6c1_51720e9cf. (a)** The metadata for b1913e6c1_51720e9cf were

210

collected via an online Google Form and stored both online (Google Sheet) and in a local SQL database.

211

A good metadata collection system should be (i) short and easy to complete, (ii) instantly accessible by

212

authorized users and (iii) easy to parse for humans and computers. **(b)** b1913e6c1_51720e9cf was

213

sequenced along with other samples, whose raw sequencing data were located in a directory named

214

after the date of the sequencing run. There one could find the FASTQ files containing the sequencing

215

reads from b1913e6c1_51720e9cf as well as information about their quality; no modified, subsetted or

216

merged FASTQ file was stored to ensure that analyses started off from the very same set of reads. In a

217

first step, the raw data of b1913e6c1_51720e9cf were processed with the Hi-C analysis pipeline, which

218

created a "b1913e6c1_51720e9cf" directory at the same level where all processed Hi-C samples were

219

located. "b1913e6c1_51720e9cf" had multiple subdirectories that stored the files generated in each of

220

the steps of the pipeline, the logs of the programs and the integrity verifications of key files. Moreover,

221

such subdirectories accounted for variations in the analysis pipelines (e.g. genome assembly version,

222

aligner) so that data were not overwritten. In a second step, processed data from b1913e6c1_51720e9cf

223

and other samples were used to perform the downstream analyses Chloe asked Paul. Within the

224

directory he allocated to her analyses, Paul created a new one called "2017-03-08_hic_validation" with

225

the description of the analysis along with the scripts used and the tables and figures generated.

226

227    **Figure 2. Automating the analysis and visualisation of b1913e6c1_51720e9cf data. (a)** Scalability,

228    parallelization, automatic configuration and modularity of analysis pipelines. Paul launched the Hi-C

229    pipeline for hundreds of samples with a single command (gray rectangle): the submission script

230    ("*.submit.sh") generated as many pipeline scripts as samples listed in the configuration file ("*.config").

231    The configuration file also contained the hard-coded parameters shared by all samples, such as the

232    maximum running time Paul underestimated for some samples. Processing hundreds of samples was

233    relatively fast because (i) the pipeline script for each of the samples was submitted as an independent

234    job in the computing cluster, where it was queued (orange) and eventually executed in parallel (green),

235    and (ii) the pipeline code in "*seq.sh" was adapted for running in multiple processors. For further

236    automation, each process retrieved sample-specific information (e.g. species, read length) from the

237    metadata SQL database; in addition, metrics generated by the pipeline (e.g. running time, number of

238    aligned reads) were recorded into the database. Because the pipeline code was grouped into modules,

239    Paul was able to easily re-run the "generate_matrix" module for those samples that failed in his first

240    attempt. **(b)** Interactive web application to visualise Hi-C data. b1913e6c1_51720e9cf alone generated

241    ~70 files of plots and text when passed through the Hi-C pipeline. Inspecting them might have seemed a

242    daunting task for Chloe: she did not feel comfortable navigating the cluster and lacked the skills to

243    manipulate them anyway, and even if she did, examining so many files for dozens of samples seemed

244    endless. Luckily for her, Paul had developed and interactive web application with R Shiny (**Table 2**) that

245    allowed her to visualise data and metadata and perform specific analyses in a user-friendly manner.

246

247

248 # Tables

249 **Table 1. Challenges associated to the accelerated accumulation of high throughput sequencing**

250 **data.** As storified with the lives of b1913e6c1_51720e9cf and T47D_rep2, managing and analyzing the

251 growing amount of sequencing data presents several challenges.

| Challenge | Impact | Consideration |
|---|---|---|
| Mislabelled raw sequencing data | • Underpowered analysis<br>• Erroneous results<br>• Loss of data, time and resources | Check unassigned reads and sequencing index concordance |
| Poor sample description | • Prevents data processing and quality control<br>• Incorrect analysis and results<br>• Lack of reproducibility<br>• Delays publication | Metadata collection |
| Unsystematic sample naming | • Duplicated or similar names<br>• Ambiguous identification<br>• Precludes computational treatment<br>• Data disclosure | Sample identifier scheme |
| Untidy data organisation | • Data cannot be found<br>• Time consumption<br>• Inability to automate searches | Structured and hierarchical data organisation |
| Yet another analysis | • Repeated manual execution of analyses<br>• Incapability to deconvolute analysis producing different results<br>• Compulsory linear execution | Scalability, parallelization, automatic configuration and modularity |
| Undocumented procedures | • Poor understanding of results<br>• Irreproducibility<br>• Hampers catching errors | Documentation |
| Data overflow | • No access to data<br>• Size and number of files make individual inspection inefficient | Interactive web applications |

252

253

11

254 **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/ |

255

256    ## Additional files

257    **Additional file 1. (a)** More than reads. FASTQ files may be useless if not coupled with biological,

258    technical and logistics information (metadata). Metadata are used at several stages of the high

259    throughput sequencing data. In the initial processing, for instance, the human origin of

260    b1913e6c1_51720e9cf was needed to determine hg38 as the reference genome sequence to which

261    reads would be aligned, and the restriction enzyme "DpnII" applied in the Hi-C protocol was used in the

262    mapping too. Other metadata were used for quality control (e.g. sequencing facility and/or date for

263    detecting batch effects or rescuing swapped samples using the correct index) or in the downstream

264    analysis (e.g. cell type, treatment). Furthermore, metadata is critical for data sharing and reproducibility.

265    **(b)** Choosing a name. Long before b1913e6c1_51720e9cf was generated, a scheme to name Hi-C

266    samples was envisioned. First, two sets of either biological or technical fields that unequivocally defined

267    a sequencing sample were identified. Then, for a given sample the values of the biological fields treated

268    as text are concatenated and computationally digested into a 9-mer, and the same procedure is applied

269    to the technical fields. The two 9-mers are combined to form the sample identifier (ID), as happened for

270    b1913e6c1_51720e9cf. Despite the apparent non-informativeness of this sample ID approach, it easily

271    allows identifying biological replicates and samples generated in the same batch since they will share,

272    respectively, the first and second 9-mer. While the specific fields used to generate the sample ID can

273    vary, it is important that they unambiguously define a sequencing sample (otherwise duplicated

274    identifiers can emerge) and that they are always combined in the same order to ensure reproducibility.

275    Indeed, another advantage of this naming scheme is that the integrity of the metadata can be checked,

276    as altered metadata values will lead to a different sample ID.

277    **Additional file 2. Rapid accumulation and diversity of high throughput sequencing (HTS) data.** The

278    past decade has witnessed a tremendous increase in sequencing throughput and applications, causing

279    uncontrolled accumulation of sequencing datasets. **(a)** For instance, the number of sequences deposited

280    in the Sequence Read Archive (SRA) [2], a major repository for HTS data, has skyrocketed from ~2

281    Terabases in 2009 to ~9,000 Terabases (the size of approximately 3 million human genomes) at the

282    beginning of 2017. Moreover, this is surely an underestimation of the actual amount given that only

283    sequencing experiments eventually included in a publication are deposited. Although data-intensive

284    projects like TCGA [3], 1000 Genomes Project [4], ENCODE [5] and 3K RGP [6] are top HTS data

285    generators [7], such a boost in the number of existing sequences reflects a pervasive use of HTS. **(b)** As

286    an example, while sequencing data for >90,000 studies have been submitted to the SRA, the top 10 and

287    100 contributors in terms of number of bases represent only a part of the archive (~30% and ~60%

288    respectively). **(c)** Similarly, while ~80% of SRA data derive from *Homo sapiens* and *Mus musculus*, the

289    central organisms in large sequencing projects, the remaining 20% come from a diverse number of

290    organisms (~50,000). Data were obtained from [8] and processed as described in the didactic dataset.

291    **Additional file 3. Why T47D_rep2 and b1913e6c1_51720e9cf are not singletons.**

292    **Additional file 4. Number of SRA deposited bases grouped by instrument name.** Data were obtained

293    from [8] and processed as described in the didactic dataset.
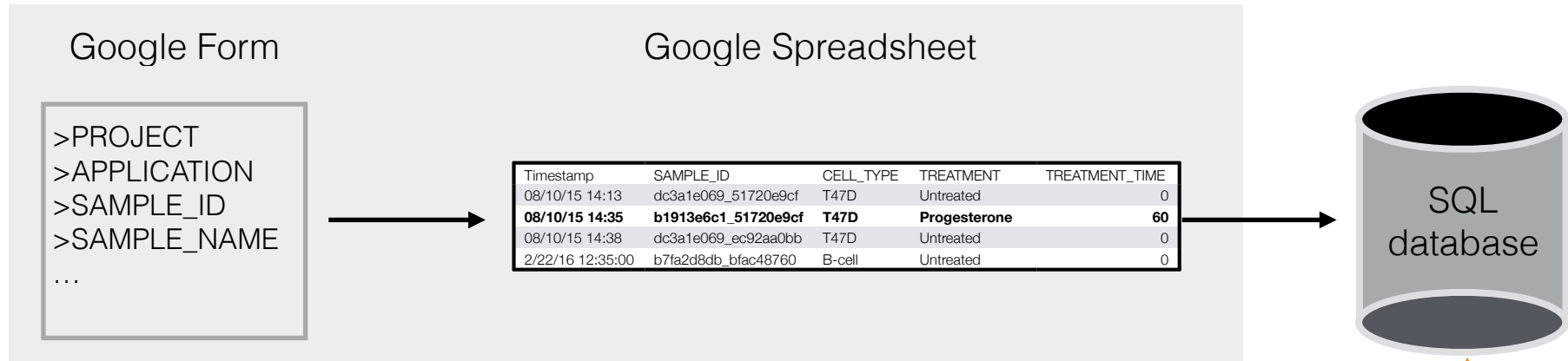
294    **Additional file 5. Number of SRA deposited bases grouped by the submitter.** For the top 25

295    contributors in terms of number of bases submitted, we searched for instances of multiple entries

296    probably referring to the same submitter (e.g. 'ncbi' and 'NCBI'). Data were obtained from [8] and

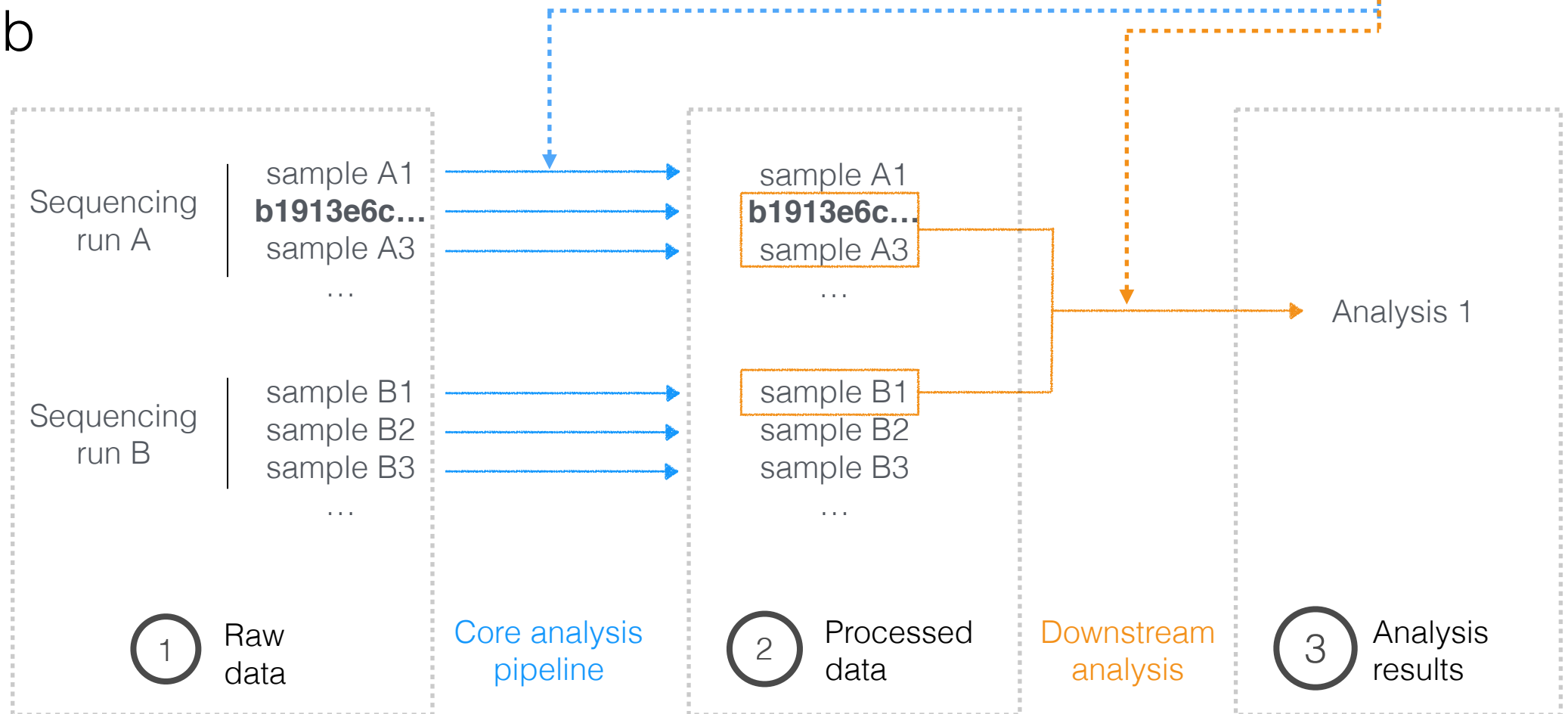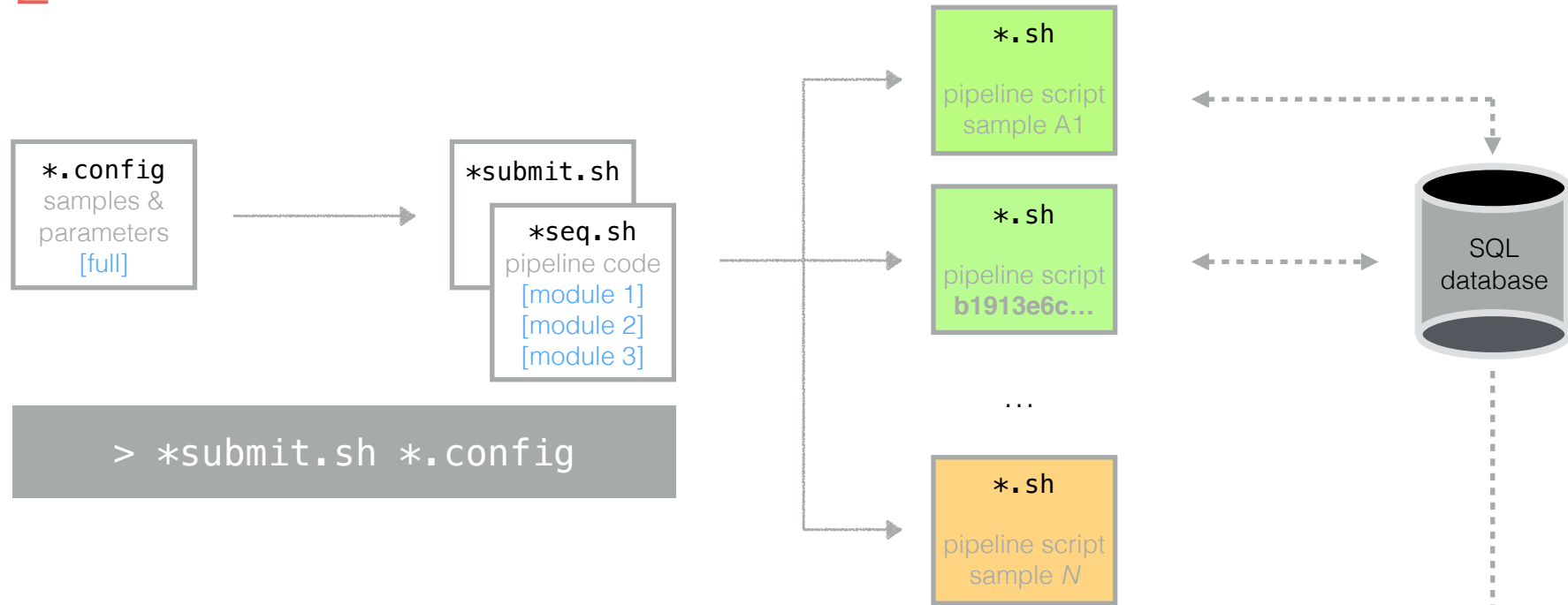297    processed as described in the didactic dataset.

14

Fig. 1

Fig. 2