

# 1 Digestiflow: from BCL to FASTQ with ease

2

3 Manuel Holtgrewe<sup>1,2</sup>, Mikko Nieminen<sup>1,3</sup>, Clemens Messerschmidt<sup>1,2</sup>,

4 Dieter Beule<sup>1,3</sup>

5

6<sup>1</sup> Berlin Institute of Health, Core Unit Bioinformatics, Charitéplatz 1, 10117 Berlin

7<sup>2</sup> Charité – Universitätsmedizin Berlin, Charitéplatz 1, 10117 Berlin

8<sup>3</sup> Max-Delbrück-Center for Molecular Medicine, Robert-Rössle-Straße 10, 13125 Berlin

9

10 Corresponding Author:

11 Dieter Beule<sup>1,3</sup>

12 Charitéplatz 1, 10117 Berlin

13 Email address: dieter.beule@bihealth.de

14

## 15 Abstract

16 Management raw sequencing data and its preprocessing (conversion into sequences and  
17 demultiplexing) remains a challenging topic for groups running sequencing devices. They face  
18 many challenges in such efforts and solutions ranging from manual management of spreadsheets  
19 to very complex and customized LIMS systems handling much more than just sequencing raw  
20 data. In this manuscript, we describe the software package DigestiFlow that focuses on the  
21 management of Illumina flow cell sample sheets and raw data. It allows for automated extraction  
22 of information from flow cell data and management of sample sheets. Furthermore, it allows for  
23 the automated and reproducible conversion of Illumina base calls to sequences and the  
24 demultiplexing thereof using bcl2fastq and Picard Tools, followed by quality control report  
25 generation.

26

## 27 1 Introduction

28 Laboratories operating modern sequencing facilities face a multitude of challenges. These  
29 include sample tracking, raw data preprocessing (conversion of raw sequencer output into  
30 sequences and demultiplexing of pooled experiments which is usually done in the same step),  
31 quality control of sequencing results, and delivery to the requesting party. While there is no clear  
32 consensus of what comprises a Laboratory Information Management System (LIMS), the term  
33 LIMS is often used to describe systems supporting these step. Simple “pure peopleware”  
34 implementations consist of spreadsheets on network shares while comprehensive commercial  
35 packages such as Illumina BaseSpace Clarity LIMS offer highly adjustable but very expensive

36 solutions. A number of academic and open solutions fall in between, offering a variable number  
37 of features and degrees of customizability.

38 The general lack of agreement of what a LIMS should cover or not cover stems from the fact that  
39 sequencing laboratories alone differ greatly. Areas of difference include the types of samples  
40 accepted (tissues/blood, DNA/RNA, final libraries/pools, or a subset thereof), and the type of  
41 data generated (raw base calls, sequences, aligned reads, or bioinformatics analytical reports). In  
42 addition, the surrounding information technology (IT) infrastructure varies greatly as does the  
43 degree of integration with such additional IT systems.

44 In this manuscript we present our approach DigestiFlow (DF) that addresses the different needs  
45 of organizations by focusing on a small, well-defined subset of tasks: management of Illumina  
46 flow cell and sample sheet information and orchestrating the step converting base calls to  
47 sequences and demultiplexing pooled sequencing runs. To the best knowledge of the authors, in  
48 this domain DF offers unparalleled functionality. Flow cells can be filled with an arbitrary  
49 combination of libraries using any combination of index and molecular barcode reads. DF also  
50 supports the barcode being part of the template sequence. DF provides extensive features for  
51 sanity checking and comparison of expected indexing reads with those actually seen in the raw  
52 base call data.

53 This is particularly important in an era where technologies such as single cell and low input  
54 sequencing require an ever-growing complexity of barcoding and indexing schemes and the  
55 amount of sequencer throughput is growing dramatically. We have encountered flow cells with  
56 more than 600 libraries and expect this to grow with increasing sequencer throughput.

57 A fundamental link to central IT is the integration with existing authentication infrastructure via  
58 directory servers, e.g., Microsoft ActiveDirectory (AD). DF supports linking accounts to central  
59 AD instances as well as using user accounts that only exist within the system. Beyond this, the  
60 system provides its functionality through a REST API (representational state transfer application  
61 programmable interface application programming interface) such that other services can be  
62 easily integrated. Instead of covering all possible functionality and sample tracking schemes, DF  
63 avoids the complexity of a monolithic system and can be integrated as a part of a modular  
64 system. However, it can also just be standalone without integration with any other system.

65

## 66 **2 Methods**

67 DigestiFlow (DF) consist of three major components. The architecture of the system is shown in  
68 Figure 1. The figure also shows interaction with a minimal set of external systems.

69

### 70 **2.1 Digestiflow Server**

71DF Server is a web app implemented with the SODAR-Core (Nieminen, Stolpe, Schumann,  
72Holtgrewe, & Beule, n.d.) and Django frameworks in the Python programming language. It uses  
73a PostgreSQL database system. It allows for the curation of flow cells and libraries together with  
74arbitrarily complex index and barcoding schemes. Barcodes can be organized in barcode sets  
75such that their sequence can be entered once and subsequently be referred to by name.  
76Furthermore, sequencing machines can be registered with their main properties (e.g., whether the  
77second barcode read needs to be reverse-complemented, depending on the paired indexing  
78workflow used). DF Server allows the visualization of barcodes detected by DF Client in the  
79BCL files (see Figure 2) and compares them to the libraries and barcodes entered by the users  
80into the flow cell sample sheet. DF Server provides a number of sanity checks for both barcodes  
81from raw sequencing data, including a barcode frequency distribution and recognizing expected  
82spike-ins such as PhiX sequence. It can also cross check between sample sheets and barcodes in  
83the raw sequences, detecting barcodes present in one but missing in the other. Furthermore, users  
84can add comments to flow cells and attach arbitrary files, which is useful for exchanging  
85spreadsheets or concentration measurement reports from the wet lab.

86

## 872.2 DF Client

88DF Client is meant to be called periodically via a cron job to monitor the storage volume where  
89sequencer(s) write output data. It reads the metadata files and registers any new flow cell with  
90DF Server (or alternately, flow cells can be pre-registered in DF Server and their properties are  
91then updated by DF Client). Once the barcodes have been sequenced completely, the DF Client  
92extracts and evaluates their sequences and posts this information to DF Server. The client also  
93detects when sequencing has succeeded (and various failure conditions) and updates the  
94information in the server. DF Client is written in the Rust programming language.

95

## 962.3 Digestiflow Demux

97DF Demux is also meant to watch the storage volume where the sequencers write their output  
98data. Once sequencing of a flow cell is complete and marked as ready in DF Server, it starts the  
99preprocessing by first obtaining the flow cell information from DF Server. Flow cells can be  
100marked for delivery as base call (BCL) files, (possibly) demultiplexed sequences, or both. If raw  
101BCL files are to be delivered, DF Server simply creates a TAR (tape archive) file for each lane  
102that contains all the information required for demultiplexing this one lane.

103For preprocessing, it first checks whether the flow cell can be processed by simply calling the  
104Illumina vendor software bcl2fastq (version 1 or 2, depending on the needs of the raw data) and  
105calls the program accordingly. Otherwise, it generates a series of calls to bcl2fastq and Picard

106Tools (<http://broadinstitute.github.io/picard/>) to perform the required preprocessing. An example  
107for this is the Agilent XT protocol where molecular barcode sequences are stored in the second  
108barcode read which bcl2fastq does not support. We refer to homogenous flow cell loads that can  
109be processed with the bcl2fastq as basic preprocessing while flexible preprocessing allows  
110arbitrary combination of library indexing and barcoding schemes.

111Once preprocessing is complete, FastQC

112(<https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>) is run on the results and the quality  
113control results are be collected with MultiQC (Ewels, Magnusson, Lundin, & Källér, 2016).

114Finally, the MultiQC report is posted as a message to the flow cell in DF Server using the REST  
115API together with the log files as attachment. This then allows the sequencing staff to review the  
116results and react accordingly. DF Demux is implemented in the Python programming language  
117with a Snakemake (Köster & Rahmann, 2012) workflow using Bioconda (Grüning et al., 2018)  
118for software package management.

119

## 1203 Results

### 1213.1 The States of a Flow Cell

122DF tracks three components of sequencing: (a) the sequencing process itself, (b) preprocessing,  
123and (c) data delivery. The possible state values differ for each of these three steps. See the  
124manual in the Supplemental Material for full details, but they can be summarized as follows.  
125Each step starts in the initial state. For preprocessing, an operator has to set the state of the  
126preprocessing step to ready which signals DF Demux to start. Once the client detects that the  
127sequencing of a flow cell has started the sequencing state changes to running. Similarly, once DF  
128Demux has started, the preprocessing state changes to running. If sequencing or demultiplexing  
129fails or succeeds, the corresponding state is updated accordingly (failed/success). For both, a  
130human operator can set a special confirmed failure/success state manually. For example, a  
131confirmed failure state will be set manually after they determine preprocessing failed due to  
132overall low sequencing quality and it is not possible to “rescue” preprocessing by fixing a sample  
133sheet. Or a confirmed success state may be set after a human operator determines that QC passes  
134visual inspection. A special success with warning state allows users to flag situations such as  
135sequencing which succeeded for all but one lane due to technical issues.

136For delivery, a human user has to set the state explicitly. This built-in system for keeping track of  
137the delivery state is particularly useful if more than one user is handling data delivery, especially  
138when used in conjunction with the message feature for leaving notes on flow cells.

### 1393.2 Comparison with Existing Methods

140 Existing methods include the following: openBIS ELN-LIMS (Barillari et al., 2016) which  
141 builds on top of openBIS (Bauch et al., 2011) and has a high number of features for sample  
142 submission and -tracking yet also has a large number of dependencies, MendeLIMS (Grimes &  
143 Ji, 2014) which has basic sample tracking functionality yet is bound to a rigid data processing  
144 workflow, MISO (Masella et al., 2019) which offers basic sample tracking functionality yet does  
145 not include features for preprocessing, and Parkour LIMS (Anatskiy et al., 2019) which provides  
146 extensive sample tracking and advanced lab notebook features yet also does not integrate  
147 automated preprocessing. While being out of scope of this manuscript, we note that DF could be  
148 integrated with other software packages as long as they provide an API with additional code. The  
149 integration with Parkour LIMS appears particularly appealing as it is based on the same  
150 technology as DF (Python/Django) and has few other dependencies itself. Table 1 contains a  
151 comparison of the listed tools given some important features  
152

### 153 3.3 Features for Improving Sequencing Results

154 **Sample Sheet Validation.** Based on practical experience, we greatly appreciate the automated  
155 comparison of observed adapter sequence content and sample sheet. Unexpected sequence in  
156 either set is an indication for possible errors. DF Server provides fine-grained control to  
157 acknowledge and suppress inconsistency warnings (after either fixing errors or accepting errors  
158 and then excluding corresponding data). Furthermore, common artifacts such as PhiX sequence  
159 are automatically recognized and show up as information rather than warnings or errors. Figure 2  
160 shows an example.

161 **Reproducibility, Automation, and Quality Control.** The Digestiflow Client and Demux  
162 components are available from Bioconda as Conda packages and Docker images, thus allowing  
163 for future proof installations and creating reproducible workflows. By offering REST APIs and  
164 two useful client applications, DF greatly supports sequencing and demultiplexing operators in  
165 automating their work. Further automation can be added later as the APIs are open. Automated  
166 quality control using FastQC and aggregation using MultiQC also allows users to spot problems  
167 earlier (together with the sample sheet adapter checks described above). In our experience this  
168 allows for the early detection of many common issues. For example, from time to time, it occurs  
169 that the same adapter was used for two different libraries in the same lane. This error might be  
170 hard to spot on paper or in spreadsheets but applications such as DF Server can easily detect and  
171 report such problems similar to the example shown in Figure 2.

172

## 173 References

174

175Anatskiy, E., Ryan, D. P., Grüning, B. A., Arrigoni, L., Manke, T., & Bönisch, U. (2019). Parkour LIMS:  
176high-quality sample preparation in next generation sequencing. *Bioinformatics*.  
177doi:10.1093/bioinformatics/bty820

178Barillari, C., Ottoz, D. S. M., Fuentes-Serna, J. M., Ramakrishnan, C., Rinn, B., & Rudolf, F. (2016).  
179openBIS ELN-LIMS: an open-source database for academic laboratories. *Bioinformatics*, 32(4), 638–640.  
180doi:10.1093/bioinformatics/btv606

181Bauch, A., Adamczyk, I., Buczek, P., Elmer, F.-J., Enimanev, K., Glyzowski, P., ... Rinn, B. (2011).  
182openBIS: a flexible framework for managing and analyzing complex data in biology research. *BMC*  
183*Bioinformatics*, 12(1), 468. doi:10.1186/1471-2105-12-468

184Ewels, P., Magnusson, M., Lundin, S., & Käller, M. (2016). MultiQC: summarize analysis results for  
185multiple tools and samples in a single report. *Bioinformatics*, 32(19), 3047–3048.  
186https://doi.org/10.1093/bioinformatics/btw354

187Grimes, S. M., & Ji, H. P. (2014). MendelIMS: A web-based laboratory information management system  
188for clinical genome sequencing. *BMC Bioinformatics*. doi:10.1186/1471-2105-15-290

189Grüning, B., Dale, R., Sjödin, A., Chapman, B. A., Rowe, J., Tomkins-Tinch, C. H., ... Köster, J. (2018).  
190Bioconda: sustainable and comprehensive software distribution for the life sciences. *Nature Methods*,  
19115(7), 475–476. doi:10.1038/s41592-018-0046-7

192Köster, J., & Rahmann, S. (2012). Snakemake--a scalable bioinformatics workflow engine. *Bioinformatics*,  
19328(19), 2520–2522. doi:10.1093/bioinformatics/bts480

194Masella, A., Cooke, D., Armstrong, H., Davey, R., DeBat, T., Bian, X., ... Leipzig, J. (2019). miso-  
195lims/miso-lims: v0.2.183. doi:https://10.5281/ZENODO.3341739

196Nieminen, M., Stolpe, O., Holtgrewe, M., Beule, D. SODAR Core: a Django-based framework for scientific  
197data management and analysis web apps. under review

198

199**Tables**

200

201**Table 1** Comparison important properties and features in commercial and free software for the  
 202management of Illumina flow cells information popular in the sequencing community based.

203

Metric	DF	BSCL	OBLE	ML	MISO	PL
License	MIT	com.	f.f.n.	f.f.n.	GPL	GPL
Self-Hosted	✓	–	✓	✓	✓	✓
LDAP Auth	✓	✓	✓	✓	✓	–
(REST) API	✓	✓	✓	–	✓	✓
Sample Tracking	min.	adv.	basic	basic	basic	adv.
Basic Preproc.	✓	✓	✓	✓	–	✓
Flexible Preproc.	✓	–	–	–	–	–
Sheet Checks	✓	–	–	–	–	–
BCL Checks	✓	–	–	–	–	–

204

205Abbreviations used in the table: DF (DigestiFlow), BSCL (BaseSpace Clarity LIMS), OBLE (OpenBIS LIMS-  
 206ELN), ML (MendelLIMS), PL (Parkour LIMS), com. (commercial), f.f.n. (free for non-commercial), min. (minimal),  
 207adv. (advanced), preproc. (preprocessing).

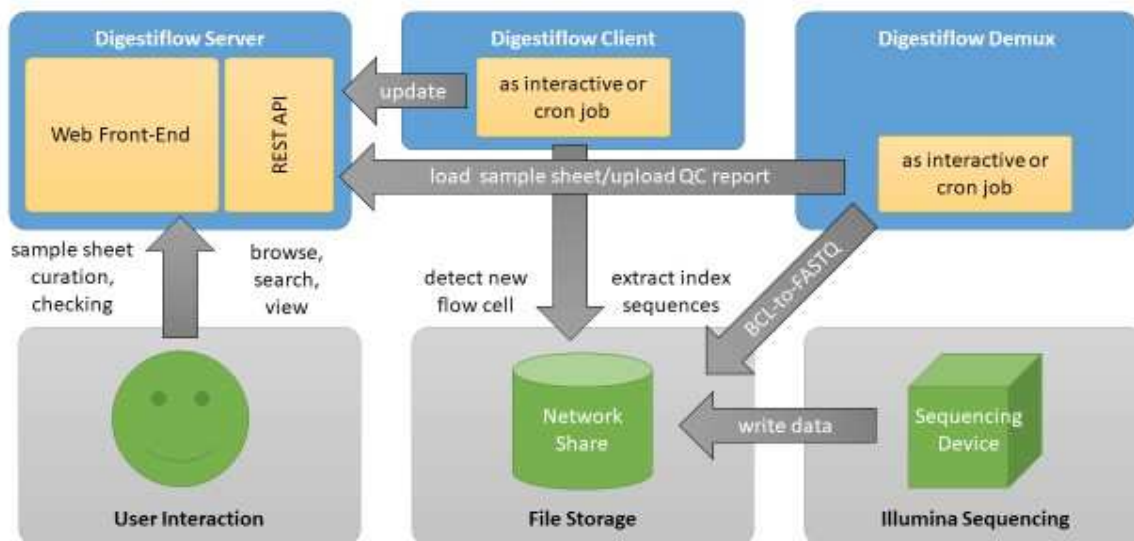


## 208 Figures

209

210 **Figure 1** Architectural overview. Sequencing instruments write data to a specified file system  
 211 storage. A periodically running DF Client detects new flow cells and registers them with the DF  
 212 Server. Once sequencing is complete and sample sheet information has been approved by the  
 213 operator, DF Demux performs the conversion to FASTQ files and creates all QC reports. Users  
 214 can not only browse and view, but also manage and curate flow cells and their sample sheets  
 215 through DF Server.

216



217

218



219

**220Figure 2**

221When adding the sample sheet (not shown), the operator made a small mistake. The adapter P37  
 222is given twice for the same lane in the sample sheet while the adapter sequence “AAGACCGT”  
 223occurs in the raw base calls but not in the sample sheet. This information can then be used for  
 224debugging sample sheet information. This is highlighted in the sample sheet (a) and the display  
 225of the adapters read from the raw base call data (b).

(a)

Flow Cell "190118\_ST-K00107\_0120\_A\_H2WKLBBXY\_CHARITE"

#	Name	Reference	Barcode Set #1	Barcode #1	Barcode Set #2	Barcode #2	Lane
1	CH_LB_AL_001	unknown	NEBNext Multiplex Oligos f...	TTACCGAC (P1)	-	-	1-8
2	CH_LB_AL_002	unknown	NBRNext Multiplex Oligos f...	-	-	-	1-8
3	CH_LB_AL_003	unknown	NEBNext Multiplex Oligos f...	-	-	-	1-8
4	CH_LB_AL_004	unknown	NEBNext Multiplex Oligos f...	TACGGTCT (P37)	-	-	1-8
5	CH_LB_AL_005	unknown	NEBNext Multiplex Oligos f...	TACGGTCT (P37)	-	-	1-8
6	CH_LB_AL_006	unknown	NEBNext Multiplex Oligos f...	CAGGTTCA (P6)	-	-	1-8
7	CH_LB_AL_007	unknown	NEBNext Multiplex Oligos f...	CGCAACTA (P9)	-	-	1-8
8	CH_LB_AL_008	unknown	NEBNext Multiplex Oligos f...	CACAGACT (P2)	-	-	1-8
9	CH_LB_AL_009	unknown	NEBNext Multiplex Oligos f...	AGTGACCT (P2)	-	-	1-8
10	CH_LB_AL_010	unknown	NBRNext Multiplex Oligos f...	AGCCTATC (P14)	-	-	1-8

(b)

Flow Cell "190118\_ST-K00107\_0120\_A\_H2WKLBBXY\_CHARITE"

The index histogram statistics are computed from the flow cell's raw base calls. They can be used for sanity-checking your sample sheets.

Lane	Index Read	Frequencies
1	AGGACAC	( 1.86%)
	TCCAGTT	( 1.29%)
	CCAGTATC	( 1.12%)
	TATACCA	( 1.89%)
	ATAACCC	( 1.85%)
	TATCTCC	( 1.83%)
	GAGGTAC	( 8.99%)
	AACAGAC	( 8.94%)
	CCCTCSA	( 8.92%)
	AGACCTA	( 8.89%)
	GCATCC	( 8.87%)
	CAGTAC	( 5.81%)
	GTCGTAC	( 1.79%)
	GTCACTA	( 1.27%)
	TCCTCAT	( 1.11%)
	TACCSAC	( 1.89%)
	AGTACCT	( 1.84%)
	CTTACCA	( 1.83%)
	CTTAGAC	( 8.89%)
	CCGATTC	( 8.83%)
AGCTATC	( 8.82%)	
CTTCTAG	( 8.89%)	
CAGTAC	( 5.81%)	
GATCACC	( 8.99%)	
AGAGGAT	( 8.92%)	
GACTCAT	( 8.91%)	
GTCCTAG	( 8.89%)	
GAACGAG	( 8.87%)	
CATGACA	( 4.11%)	
CTTAGAC	( 1.32%)	
CTCAGAG	( 1.12%)	
TACGTCT	( 1.89%)	
CAGTCTT	( 1.87%)	
CAGTTCA	( 1.84%)	
TTCCSAA	( 1.89%)	
TGAGTTG	( 8.99%)	
ATCTAGC	( 8.92%)	
TATGACG	( 8.98%)	
TGDTARE	( 8.87%)	
CACAGACT	( 8.92%)	

**226 Supplemental Material**

227

228 • Digestiflow Server Documentation

229