

*Supplementary Material: A step-by-step guide to I-ATAC, validating pipeline with two case studies*

**I-ATAC: Interactive pipeline for the management and pre-processing of ATAC-seq samples**

**Zeeshan Ahmed<sup>1</sup> and Duygu Ucar<sup>2</sup>**

<sup>1</sup>University of Connecticut Health Center, 195 Farmington Ave, Farmington, CT, USA

<sup>2</sup>The Jackson Laboratory for Genomic Medicine, 10 Discovery Dr., Farmington, CT, USA

Corresponding authors: Zeeshan Ahmed<sup>1</sup> and Duygu Ucar<sup>2</sup>

Email address: [zahmed@uchc.edu](mailto:zahmed@uchc.edu) and [duygu.ucar@jax.org](mailto:duygu.ucar@jax.org)

**Software Availability:** I-ATAC is programmed in Java and built at both Mac-OS-X and Windows platform.

- Its source code and executable are freely available at:  
<https://github.com/UcarLab/I-ATAC>
- Example dataset is available at:  
<https://zenodo.org/record/46079#.WAe3l5MrK7Y>
- Supporting software and dependencies are available at:  
<https://zenodo.org/record/162023#.WAe3dJMrK7Y>
- For additional information, please refer to the project webpage:  
<https://www.jax.org/research-and-faculty/tools/i-atac>

24

## 25 **Table of Contents**

26	<b>1</b>	<b>Motivation .....</b>	<b>3</b>
27	<b>2</b>	<b>I-ATAC.....</b>	<b>4</b>
28	<b>3</b>	<b>Design Description.....</b>	<b>6</b>
29	3.1	Operational Workflow of I-ATAC.....	6
30	3.2	Applications integration, data processing pipeline and project's directory structure.....	7
31	3.3	Comments workflow, operating systems and physical data storage in data cluster.....	9
32	<b>4</b>	<b>GUI Description .....</b>	<b>10</b>
33	<b>5</b>	<b>Integrated Applications Details.....</b>	<b>13</b>
34	5.1	FASTQC: .....	13
35	5.2	Trimmomatic.....	13
36	5.3	BWA.....	13
37	5.4	SAMtools .....	13
38	5.5	Picard.....	14
39	5.6	BEDtools .....	14
40	5.7	ATAC_BAM_shifter_gappedAlign.pl .....	14
41	5.8	MACS2 .....	14
42	<b>6</b>	<b>Installation and Configuration .....</b>	<b>15</b>
43	<b>7</b>	<b>Case Studies .....</b>	<b>17</b>
44	7.1	Example Dataset.....	17
45	7.1.1	Dataset Details .....	17
46	7.1.2	Input .....	17
47	7.1.3	Output.....	20
48	7.2	Case Study 2: Using GM12878 – CD4 T- Cells.....	23
49	7.2.1	Dataset Details .....	23
50	7.2.2	Input .....	23
51	7.2.3	Output.....	25
52	<b>8</b>	<b>Conclusions .....</b>	<b>29</b>
53	<b>9</b>	<b>Acknowledgments.....</b>	<b>30</b>
54	<b>10</b>	<b>Funding: .....</b>	<b>30</b>
55	<b>11</b>	<b>Conflict of Interests: .....</b>	<b>30</b>
56	<b>12</b>	<b>Additional Requirements.....</b>	<b>30</b>
57	<b>13</b>	<b>References .....</b>	<b>30</b>

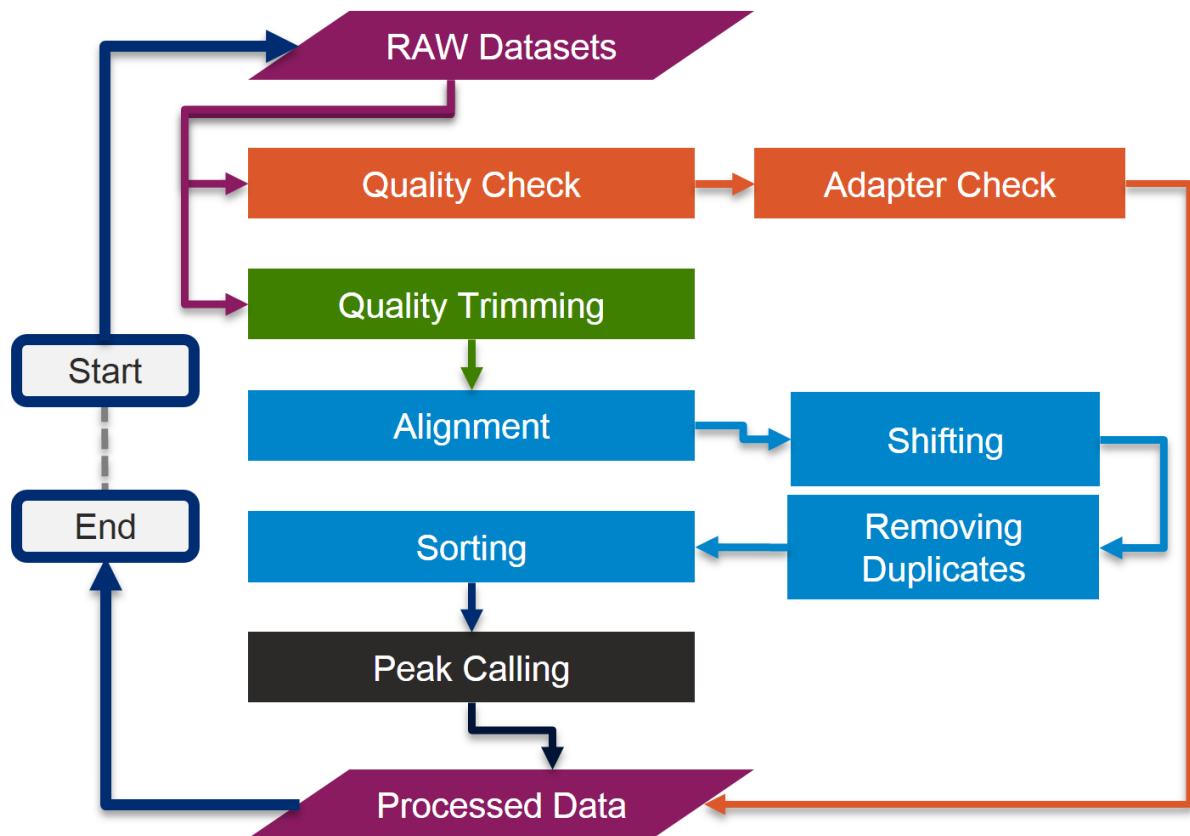
58

59

## 1 Motivation

The use of high-throughput sequencing technologies has brought an enormous increase in the amount of heterogeneous genomic data production in the last decades. The importance of genomic dataset processing in the genomic community is well known; as it plays important role in analysing the dynamics and complexities of gene regulation with modelling and implementation of different statistical methods utilizing data processing pipelines.

Traditional way of next generation sequencing (NGS) data pre-processing is complex and based on running a series of command-line applications in Unix, Linux, MAC and DOS environments, which requires good knowledge of bioinformatics tools and good programming skills. There are over 200 tools available for the genome and exome sequencing data pre-processing and analysis (Pabinger *et al.*, 2013) but most of them are non-interactive and command line based. Writing complex command line scripts and pipelines, and running non-interactive mode applications might be convenient for the scientists with good bioinformatics background but it is very hard for the biologist with no programming skills to conduct complex data analyses. The focus of our research is toward the application of a novel epigenomic profiling assay for transposase-accessible chromatin with high throughput sequencing (ATAC-seq) for integrative epigenomic analysis (Buenrostro *et al.*, 2013). ATAC-seq is a protocol to capture open chromatin sites (Buenrostro *et al.*, 2013; Buenrostro *et al.*, 2015a) by performing adaptor ligation and fragmentation of open chromatin regions (Tsompana and Buck, 2014).



S-Fig. 1. ATAC-seq data pre-processing pipeline's workflow.

ATAC-seq has been a popular chromatin profiling technology for clinical samples and has been used for the assessment of chromatin accessibility in various cells and tissues in human and model organisms e.g. (Moskowitz *et al.* 2017; Miskimen *et al.*, 2017; Bao *et al.* 2015; Buenrostro *et al.*, 2015b) etc. Due to its efficiency in requirement of biological sample and in library preparation time, many scientists are generating ATAC-seq libraries to decipher the chromatin landscape in a given cell type and condition of interest. To generate ATAC-seq libraries, a hyperactive molecule, Tn5 is used to cut the open chromatin and then short reads are sequenced typically from both ends (i.e., paired end). The next step is the processing of ATAC-seq samples. A typical ATAC-seq data processing pipeline's workflow is shown in S-Fig. 1, which starts with the quality check and adapter trimming, then alignment, shifting, removing duplicates, sorting and peak calling to find potential open chromatin sites, indicating active regulatory elements in each cell.

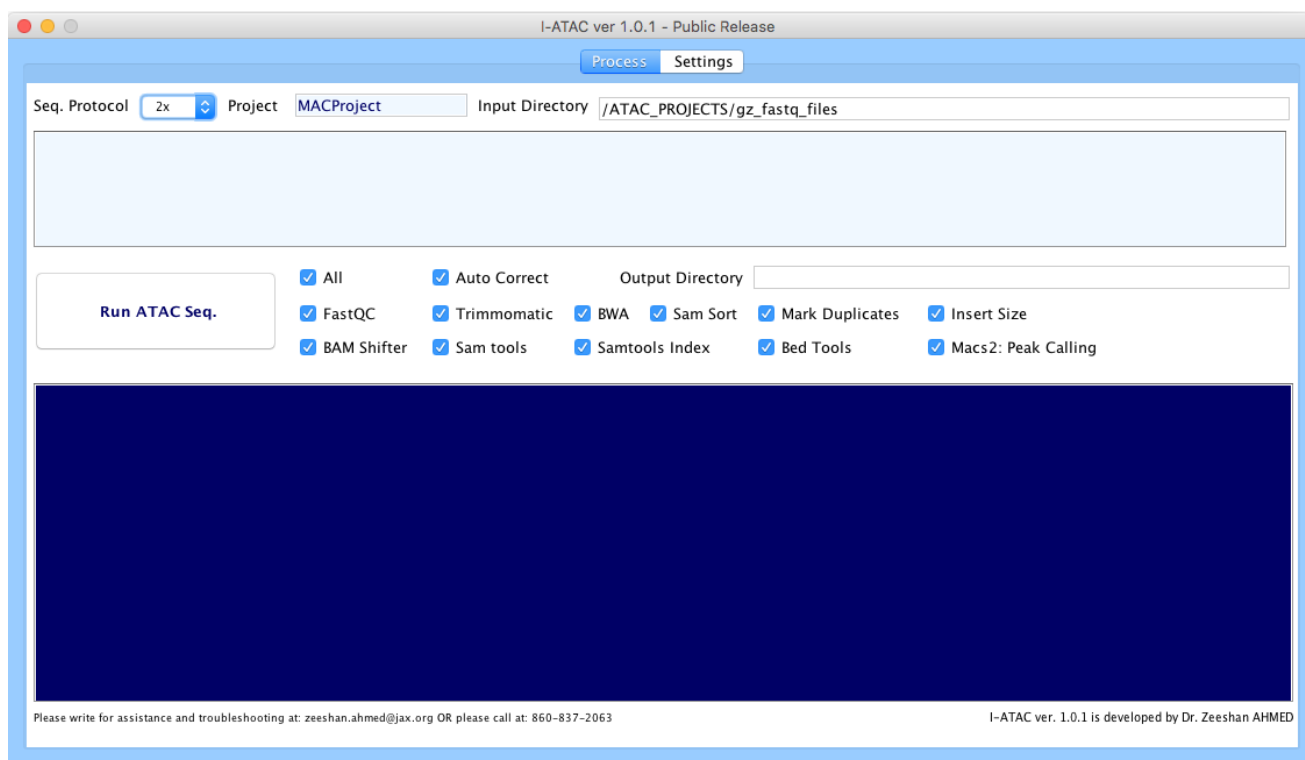
Processing and analysis of large number of ATAC-seq samples is a challenge for non-computational scientists since usually multiple tests are required to find the optimal algorithms and parameter settings. Interactive-ATAC (I-ATAC) (Ahmed and Ucar, 2017) is the first interactive, cross platform, user-friendly desktop application, which supports reproducible and automatic pre-processing of ATAC-seq (Buenrostro *et al.*, 2013; Buenrostro *et al.*, 2015) samples.

## **2 I-ATAC**

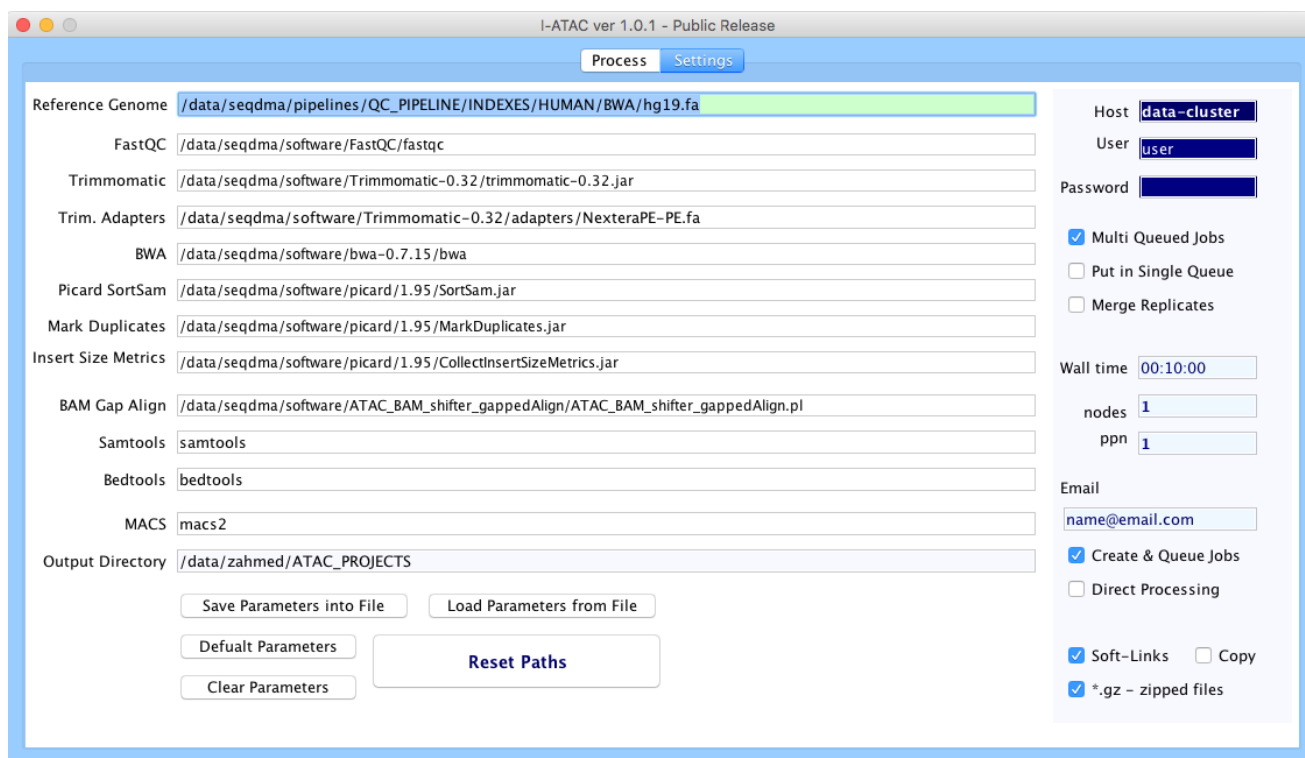
The targeted end users of I-ATAC are mainly the biologists, who are familiar and comfortable with interactive operating systems (e.g. Windows, Mac-OS-X) and applications (e.g. web based browsers or client based viewers), yet have limited experience with programming, shell scripting, and with the Unix environment. Moreover, I-ATAC could be a helpful tool for bioinformaticians, who are new to the field of epigenomic data analysis and are not familiar with ATAC-seq data processing steps.

The GUI of the I-ATAC (S-Fig. 2 A and B) is designed for simplicity and ease by following human computer interaction (HCI) guidelines (Ahmed *et al.*, 2014). The concept behind designing I-ATAC GUI was to implement "One Click Operations" concept, similar to a Google search that requires users to enter one natural language based query and click a search button. Similarly, along with the default or customized settings (S-Fig. 2 B), I-ATAC requires only path to the sample data files (zipped or unzipped "FASTQ" files), project name and pressing button "Run ATAC-Seq" (S-Fig. 2 A) to perform following tasks:

- Get user login credentials
- Connect to the data cluster or local computer
- Create output directory structure
- Locate input data
- Copy & paste or create soft links of data to process
- Load modules, compilers & interpreters
- Write command line instructions to integrate applications
- Compose shell scripts (pipeline)
- Create & queue jobs (Unix based Secure Shell Scripts) at cluster or execute instructions on local computer
- Place output files in created directory structure
- Start data processing
- Disconnect to the connected data cluster



S-Fig. 2 (A): Graphical User Interface of I-ATAC: Create and run data processing jobs.



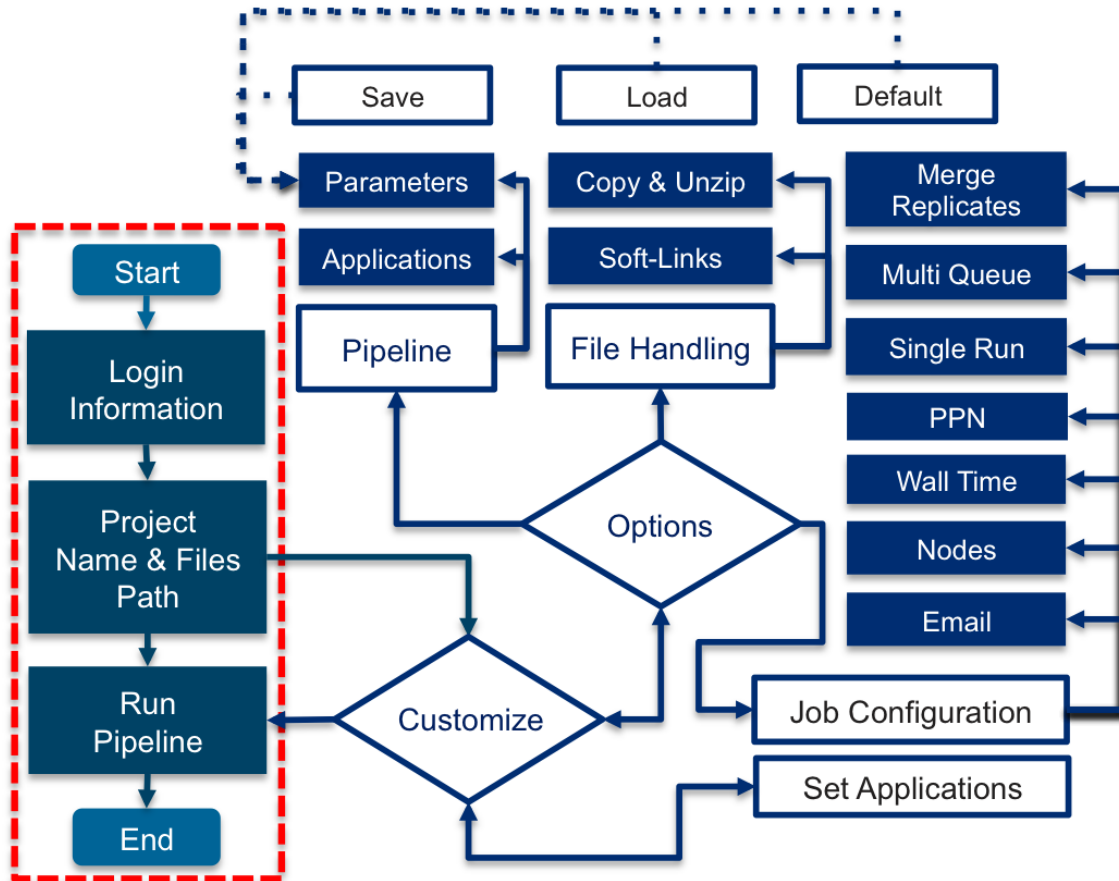
S-Fig. 2 (B): Graphical User Interface of I-ATAC: Set parameters and user credentials.

### 3 Design Description

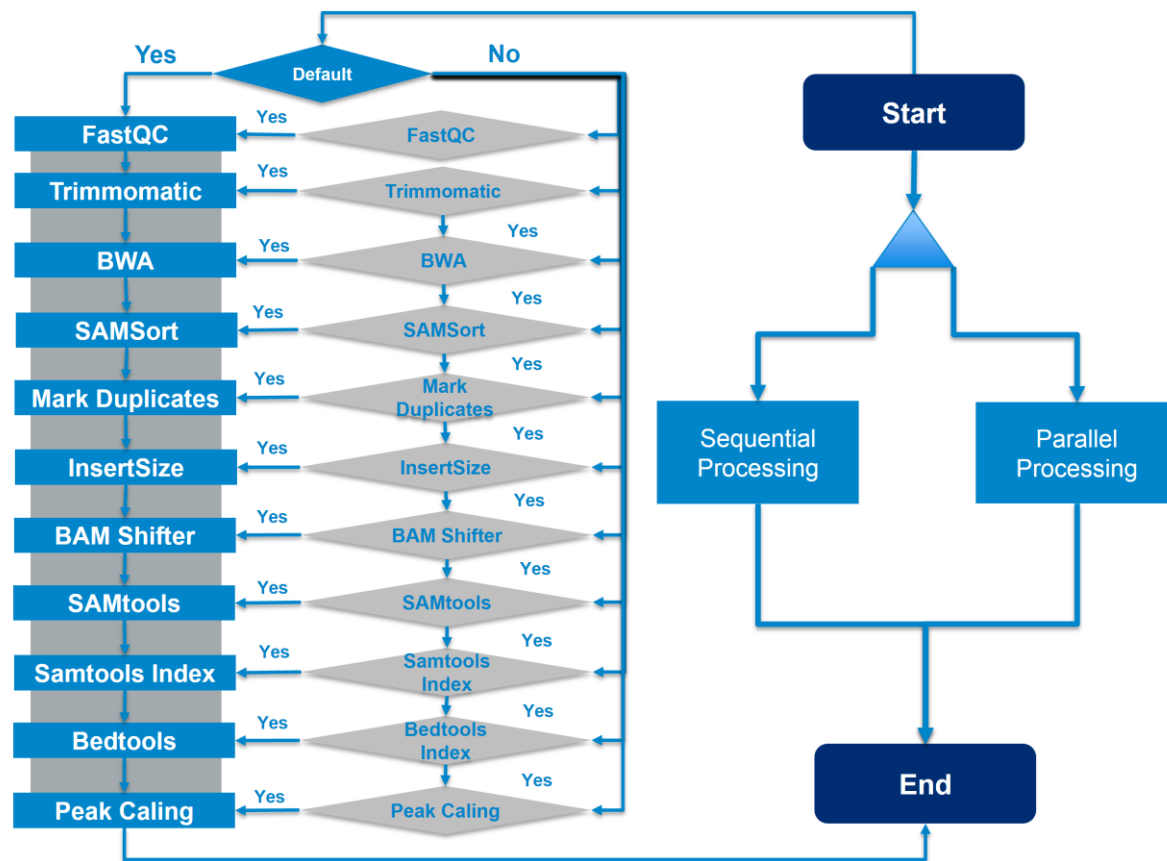
I-ATAC is a platform designed by following software engineering principles for the sustainable bioinformatics software implementation (Ahmed *et al.*, 2014). Here, we present its operational workflow, data structure and components' orientation.

#### 3.1 Operational Workflow of I-ATAC

Following default workflow (S-Fig. 3), user can process ATAC-seq samples with the application of complete pipeline, which involves the execution of all integrated applications (*FASTQC*, *Trimmomatic*, *BWA*, *Picard*, *ATAC\_BAM\_shiftt\_gappedAlign.pl*, *bedtools* and *macs2*) but user is not limited in the use of I-ATAC (S-Fig. 3). User can chose to run a single application as well as customize applications' workflow, following pre and post-requisites e.g. in case user is only interested in having FASTQC reports or trimming of low quality reads and adapters or user has already trimmed filtered FASTQ files but would like to map to reference genome only or may be only interested in generating BED files from BAM and peak calling etc. I-ATAC supported such customization and it can be very helpful, especially in trouble shooting situations, where due to any reason either pipeline could not fully execute or if there is already data exists in a form which does not require all steps of ATAC-seq pipeline. This customization can save time and computational resources.



S-Fig. 3: I-ATAC: Operational workflow of I-ATAC



147

148 S-Fig. 4: I-ATAC: Customization of ATAC-seq data pre-processing pipeline with sequential (multiple jobs in one script)  
 149 and parallel (multiple jobs in multiple scripts, one of each) processing.

150 User can remotely handle sample data files for processing by either keeping them in the same parent  
 151 directory and putting only pre-processed results in the main project and sub-project directories or by  
 152 first copying compressed files into the project directory, unzips them and then process them. User  
 153 can configure job (UNIX based Secure Shell Scripts) settings by processing one or multiple samples  
 154 at a time as one job or multiple jobs (one for each sample).

155 I-ATAC also enables users to customize parameters used for data pre-processing steps by letting the  
 156 user to choose between applications as well as by setting different parameters (S-Fig. 4), which  
 157 enables customizing this pipeline for the analyses of other data types, such as ChIP-seq data. As the  
 158 output, I-ATAC produces data quality reports that can be visualized within the platform. It also  
 159 outputs ATAC-seq reads that are filtered, trimmed and aligned as well as peak calls from these reads.

### 160 3.2 Applications integration, data processing pipeline and project's directory structure

161 ATAC-seq data processing pipeline starts with the quality check, then paired end reads are trimmed,  
 162 aligned, filtered, and sorted in a "sam" file. The "sam" file is compressed and indexed to a bam file,  
 163 which is then used as input for peak calling. To manage pre-processed data, proposed directory  
 164 structure is followed and automatically created in data cluster before data processing (S-Fig. 5).

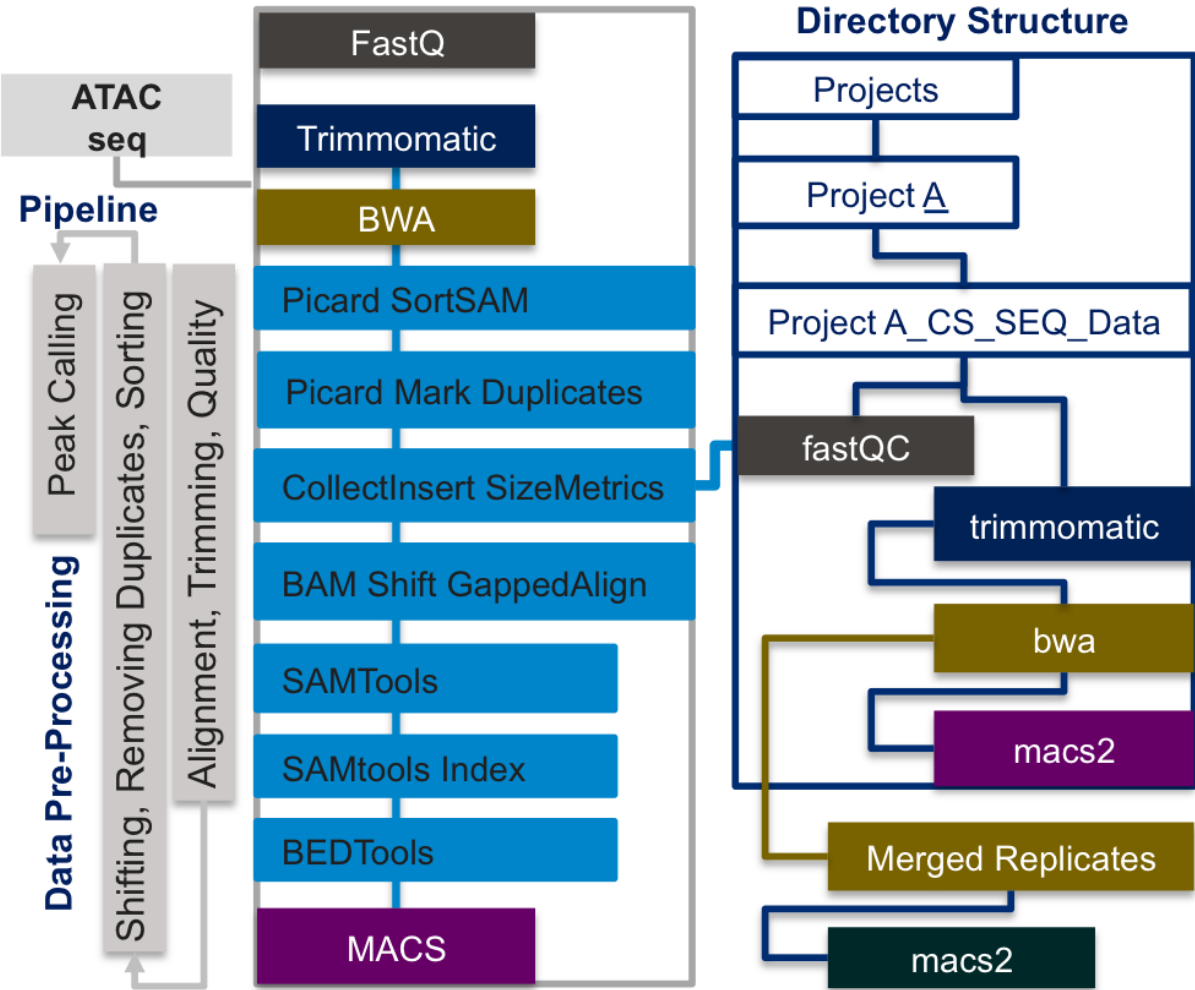
165

*A step-by-step guide to I-ATAC, validating pipeline with two case studies*

166  
167  
168  
169  
170  
171  
172  
173  
174

ATACseq\_Projects  
→ Project\_Name  
→ → Sample\_A\_R1\_Sample\_A\_R2  
→ → → fastQC  
→ → → trimmomatic  
→ → → → bwa  
→ → → → → macs2  
→ mergedreplicated  
→ → macs2

175 All the quality reports (“zip” and “html” files) are placed in “fastQC” sub-directory. Compressed  
176 files contain different output files including text (“txt”) and web page (“html”). Text file contains  
177 information about basic statistics, file name, file type, encoding, total sequences, sequence flagged  
178 quality, sequence length, base number, mean, median, lower, quartile, upper, quartile, 10th  
179 percentile, 90th percentile, quality, Count, per base sequence content, per sequence GC content, per  
180 base N content, sequence length distribution, sequence duplication levels, overrepresented sequences,  
181 adapter content and Kmer content. Whereas html file visualize quantitative results.

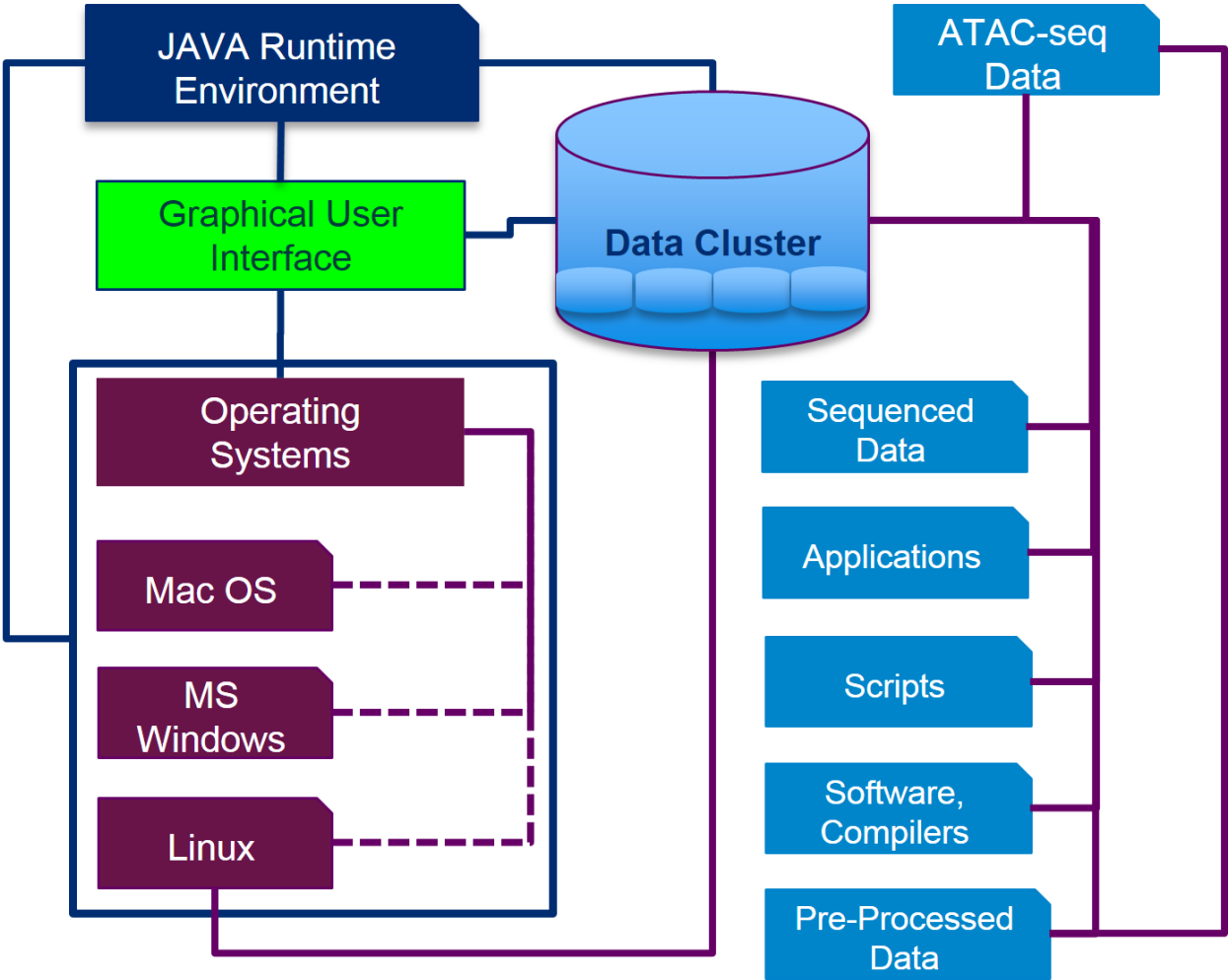


182

183 S-Fig. 5: I-ATAC: Applications and project directory structure.



All trimmed and filtered “FASTQ” files are placed in “trimmomatic” sub-directory, all the sorted, shifted “sam”, indexed “bam” and “bed” files are placed in “bwa” sub-directory. All the observed peak files are placed in the “macs2” sub-directory. The nested directory structure provides an organized and modular storage for multi-level ATAC-seq data analysis pipeline. Produced results in the form of sorted “sam” and “bam” files, as well as peaks can be visualized using available genome data browsers (e.g. UCSC, Chipster etc.) and viewers (e.g. IGV etc.).



S-Fig. 6: I-ATAC: Components workflow, operating systems and physical data storage in data cluster.

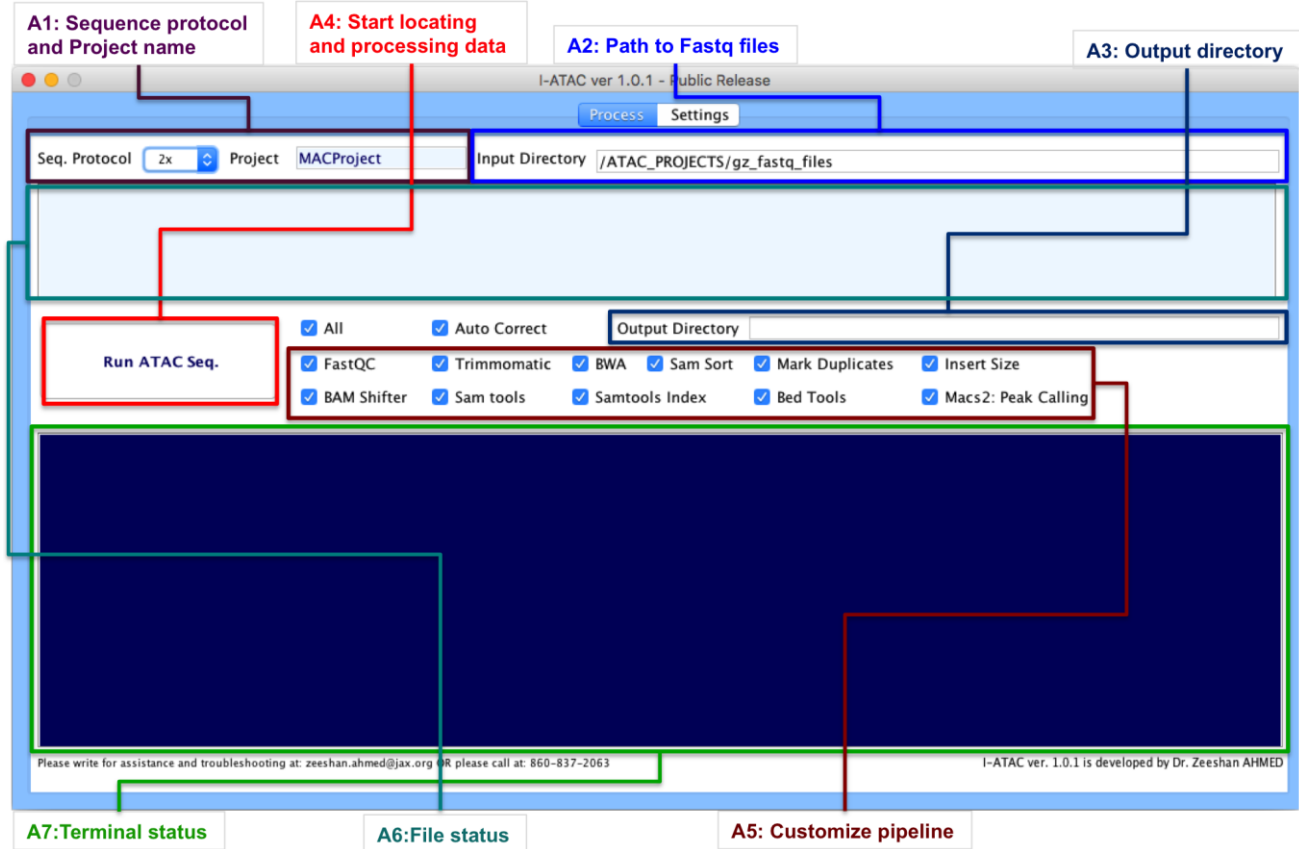
### 3.3 Comments workflow, operating systems and physical data storage in data cluster

The components workflow (S-Fig. 6) of I-ATAC depends on the Java Run Time Environment (<http://www.oracle.com/technetwork/java/javase/downloads/jre8-downloads-2133155.html>), to be installed at in-use operation system, which can be Mac-OS-X, Microsoft Windows and Linux etc. The sample, sequenced data files, applications (S-Table. 3), compilers and interpreters (S-Table. 4), pre-processed data and scripts are need to be placed in data cluster.

## 4 GUI Description

As shown in (S-Fig. 2 A and B), the overall GUI of the I-ATAC is divided in to two modules: *Process* and *Settings*.

The *Process* module is to generate and run pipeline. *Process* provides six major features: A1: *Sequence protocol and Project Name*, A2: *Output directory*, A3: *Path to FASTQ file*, A4: *Start locating and processing data*, A5: *Customize pipeline*, A6: *File Status*, and A7: *Terminal status* (S-Fig. 7 and S-Table 1).



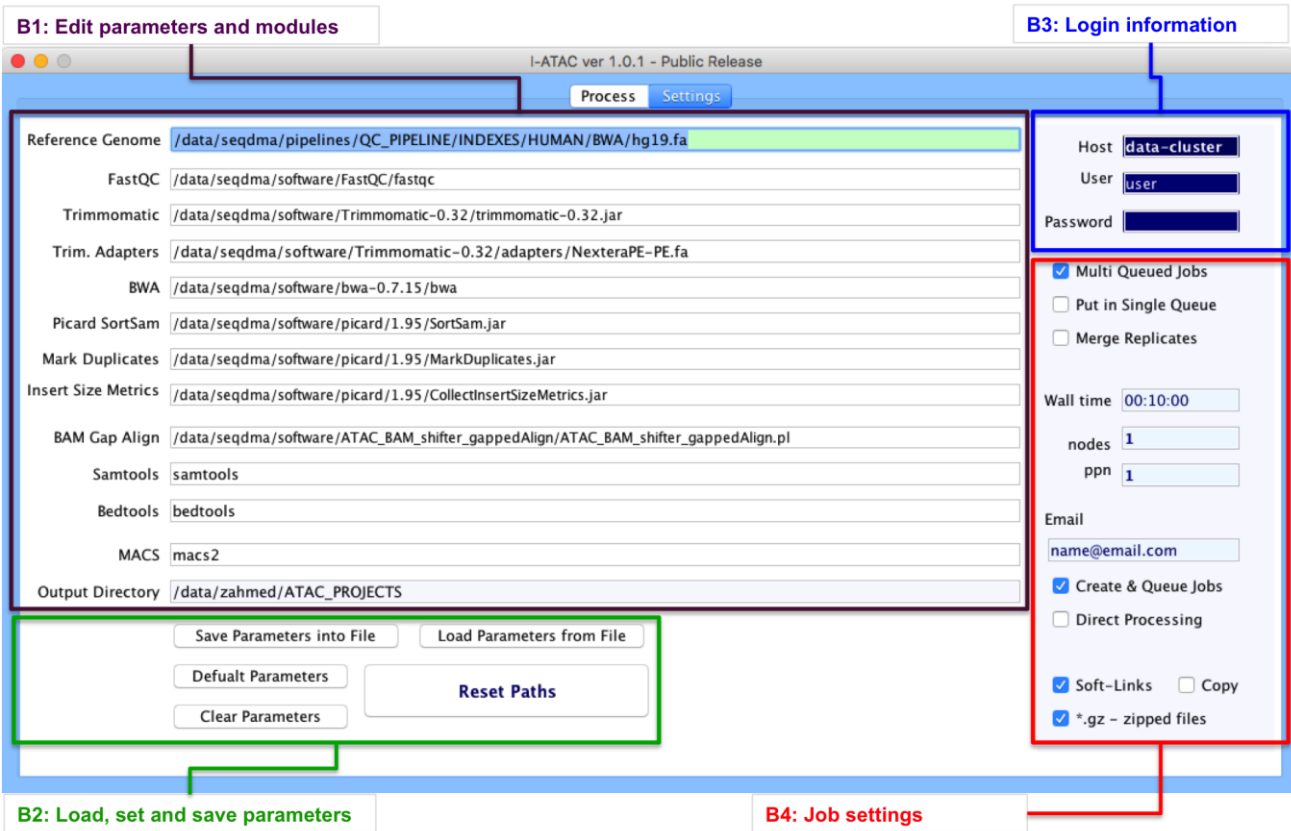
S-Fig. 7: GUI I-ATAC: Run ATAC-seq Pipeline (module: A).

Number	Feature	Description
A1	Sequence protocol and Project name	<ul style="list-style-type: none"> <li>Requires user to select either “1x” or “2x” for single and paired end data.</li> <li>Requires user to enter a Project Name, without any spaces and special characters. Reason to avoid spaces and special characters is, that, I-ATAC will automatically create a new directory, where data will be referenced (soft links) or copied for processing and results will be placed.</li> </ul>
A2	Path to the FASTQ files	<ul style="list-style-type: none"> <li>Requires path to the directory, where data (FASTQ files) are placed.</li> </ul>
A3	Output directory	<ul style="list-style-type: none"> <li>Requires path to the directory, where processed data (outcome/output files) will be placed.</li> </ul>
A4	Start locating and processing data	<ul style="list-style-type: none"> <li>Starts locating sample data files, copy from source to the main destination (project) director, unzip compressed sample data files,</li> </ul>

		automatically generate script, queue script(s) as job(s), perform data processing and place output in proposed directory structure.
A5	Customize Pipeline	<ul style="list-style-type: none"> <li>Default parameters include I/O redirected, sequential combination of integrated applications and parameters. Twelve different options are integrated: <ol style="list-style-type: none"> <li>1. <b>FASTQC</b></li> <li>2. <b>Trimmomatic</b></li> <li>3. <b>BWA</b></li> <li>4. <b>Sam Sort</b></li> <li>5. <b>Mark Duplicates</b></li> <li>6. <b>Insert Size</b></li> <li>7. <b>BAM Shifter</b></li> <li>8. <b>SAM tools</b></li> <li>9. <b>SAM tools index</b></li> <li>10. <b>BED tools</b></li> <li>11. <b>MACS2</b></li> <li>12. <b>All</b></li> <li>13. <b>Auto Correct</b></li> </ol> </li> <li>Option 12 is to select all options and perform data processing with default settings and option 13 is to correct the sequence I/O.</li> </ul>
A6	File Status	<ul style="list-style-type: none"> <li>Provides information about located data samples in the data cluster, using provided input path.</li> </ul>
A7	Terminal Status	<ul style="list-style-type: none"> <li>Provides information about execution of job in data cluster.</li> </ul>

210 S-Table. 1: Features description of GUI-A: Run ATAC-seq Pipeline

211 The GUI-B module is mainly used to set the parameters of the applications and directory paths. As  
212 shown in the figure (S-Fig. 8), it provides only four features: *Applications Parameters*, *Directory*  
213 *Paths*, *Save and load Parameters*, and *Reset Paths* (S-Fig. 8 and S-Table 2).



214 S-Fig. 8: GUI I-ATAC: Set Script Parameters (module: B).  
215

*A step-by-step guide to I-ATAC, validating pipeline with two case studies*

Number	Feature	Description
B1	Applications Parameters	<ul style="list-style-type: none"> <li>Paths and calling instructions to the following integrated applications and referenced genome and adapters: <ol style="list-style-type: none"> <li><b>Reference genome</b></li> <li><b>FASTQC</b></li> <li><b>Trimmomatic</b></li> <li><b>Adapters</b></li> <li><b>BWA</b></li> <li><b>SamSort</b></li> <li><b>Mark Duplicates</b></li> <li><b>Insert Size Metrics</b></li> <li><b>BAM Gap Align</b></li> <li><b>SAM tools</b></li> <li><b>BED tools</b></li> <li><b>Samtools Index</b></li> <li><b>MACS2</b></li> <li><b>Output Directory</b></li> </ol> </li> </ul>
B2	Load, set and save parameters	<p>Five options (buttons) are provided:</p> <ol style="list-style-type: none"> <li><b>Save Parameters:</b> To save the set parameters in the form of a text file, to reuse and share settings.</li> <li><b>Load Parameters:</b> To load saved settings in the form of a text file.</li> <li><b>Clear Parameters:</b> To clear all parameter fields.</li> <li><b>Default Parameters:</b> To load default parameters.</li> <li><b>Reset Paths:</b> To apply modifications.</li> </ol>
B3	User Login	<ul style="list-style-type: none"> <li>User requires entering name of the host (attached data cluster or name of the personal computer), user login name and password to let the I-ATAC successfully login into to host and access sample data files ("FASTQ") and applications to perform data processing.</li> </ul>
B4	Job Settings	<p>Default job (set of instructions, written in the form of a script and executed like a program (executable software) to perform certain set of operations) related parameters. Furthermore I-ATAC provides eight different options to customize script generation and job submission:</p> <ol style="list-style-type: none"> <li><b>Multi Queued Job:</b> Processes multiple samples at a time by generating and submitting parallel-multiple data processing jobs (one for each).</li> <li><b>Put in Single Queue:</b> Processes one or multiple samples at a time by generating and submitting one data processing job (one for all).</li> <li><b>Merge Replicates:</b> Applicable only in case of processing multiple samples at a time by submitting one data processing job for all. It enables selection of all generated "bam" files from all the pre-processed samples directories (bwa) and performs peak calling.</li> <li><b>Wall Time:</b> Sets time to be allocated for the processing of the queued job. In case of multiple-parallel jobs, it will set provided time for all jobs.</li> <li><b>Nodes:</b> Sets the number of nodes (connection points) requested for job. Default set node is 1.</li> <li><b>Processor per node (ppn):</b> Sets the number of cores (virtual processors) per node per. Default set ppn is 1.</li> <li><b>Email:</b> Sets to get notification (cancelled, completed) about the status of submitted job.</li> <li><b>Create &amp; Queue Jobs:</b> In case host is data cluster, then I-ATAC will prepare and submit jobs.</li> <li><b>Direct Processing:</b> In case host is personal computer, then I-ATAC will prepare and submit instructions.</li> <li><b>Creates soft links:</b> Having checked this option, I-ATAC will create soft links of FASTQ files in to output directory.</li> <li><b>Copy:</b> Having checked this option, I-ATAC will create copy FASTQ files in to output directory.</li> <li><b>*.gz zipped files:</b> Having checked this option, I-ATAC will expect input FASTQ files are zipped otherwise not.</li> </ol>

217 The sole objective of developing I-ATAC is to help with the provision of interactive ATAC-seq data  
218 processing pipeline that is why; we have not developed features for file handling between data cluster  
219 and operating systems. There are already some interactive tools available for such purposes e.g. File  
220 Zilla (<https://filezilla-project.org>), WinSCP (<http://winscp.net/eng/download.php>), Cyberduck  
221 (<https://cyberduck.io/?l=en>) etc.

## 222 **5 Integrated Applications Details**

223 ATAC-seq data processing pipeline consists of different third party applications (S-Table. 4); I/O  
224 (input/output) redirected (one's output is treated as another's input, in terms of both data analysis and  
225 processing) and integrated method (S-Fig. 6). Additionally, it requires all needed compilers and  
226 interpreters to be downloaded and installed as well (S-Table. 4).

### 227 **5.1 FASTQC:**

228 It is a command line based, non-interactive tool for the high throughput sequence data. It is  
229 programmed in Java and requires [Java Runtime Environment](#) and [Picard](#) BAM/SAM libraries to be  
230 installed in the data cluster. Its output is based on Basic Statistics, Per base sequence quality, Per tile  
231 sequence quality, Per sequence quality scores, Per base sequence content, Per sequence GC content,  
232 Per base N content, Sequence Length Distribution, Sequence Duplication Levels, Overrepresented  
233 sequences, Adapter Content and Kmer Content. FASTQC used version details, including input,  
234 output and download details are given in S-Table. 3, Row No.: 1.

### 235 **5.2 Trimmomatic**

236 It is a command line based, non-interactive tool for the trimming of reads (Bolger *et al.*, 2014) using  
237 paired-end and single ended data produced by the Illumina next generation sequencing technology  
238 (<http://www.illumina.com/>). It takes compressed or uncompressed FASTQ (phred-33 and phred-64  
239 quality scores) file as input and mainly performs adapter filtering, sliding window trimming, base  
240 cutting (start and end of reads, as well, at specific number) and removes below quality reads.  
241 Trimmomatic's used version details, including input, output and download details are given in S-  
242 Table. 3, Row No.: 2.

### 243 **5.3 BWA**

244 Burrows-Wheeler Alignment tool (BWA) is a software application for aligning short nucleotide  
245 sequences to a reference genome (Li and Durbin, 2009). It implements BWA-backtrack for reading  
246 sequence up to 100bp, and BWA-SW and BWA-MEM algorithms are for reading longer sequences  
247 between 70bp to 1Mbp. BWA's used version details, including input, output and download details  
248 are given in S-Table. 3, Row No.: 3.

### 249 **5.4 SAMtools**

250 Sequence Alignment/Map (SAM) tools is a software package with various utilities, mainly used for  
251 sequence data formatting (Li, 2011; Li, *et al.*, 2009). It helps in performing complex operations at  
252 sequence data files, including variant calling, alignment, sorting, indexing, viewing, data extraction  
253 and format conversion. SAMtools applied package's version details, including input, output and  
254 download details are given in S-Table. 3, Row No.: 4.

255

## 5.5 Picard

It is Java based non-interactive tool, which requires [Java Runtime Environment](#) to execute. It is mainly used for the sequence data manipulation in sam and bam files. Both sam and bam files contain same data structure and format, sam is human readable, whereas, bam is machine-readable format (binary). It's used version's details, including input, output and download details are given in S-Table. 3. It performs sorting in order and can read information about library, platform, sample, sequence, predicted insert size etc. Picard's used version details, including input, output and download details are given in S-Table. 3, Row No.: 5.

## 5.6 BEDtools

Browser Extensible Data (BED) tools (Quinlan and Hall, 2010) is a software application for converting "bam" to "bed" files and compare large sets of genomic features. Moreover, it can be used for converting BEDPE intervals to BAM and BAM to FASTQ, finding closest and potentially non-overlapping interval, creating HTML pages to link UCSC locations, finding pairs that overlap other pairs and intervals in various ways, randomly redistributed and adjust size of intervals and tag bam alignment etc. BEDtools used version details, including input, output and download details are given in S-Table. 3, Row No.: 6.

## 5.7 ATAC\_BAM\_shifter\_gappedAlign.pl

ATAC\_BAM\_shifter\_gappedAlign.pl is an open source Perl script, which can be used to perform read shifting based on the read quality. It takes aligned "bam" file as an input and offsets by 4bp for the positive strand (sequence containing instructions for building a protein) and -5bp for the negative strand (merely contains the complementary sequence and according to the base-pairing rules it is not normally transcribed into RNA nor translated into protein). Users can use any other tools for shifting the reads. ATAC\_BAM\_shifter\_gappedAlign version details, including input, output and download details are given in S-Table. 3, Row No.: 7.

## 5.8 MACS2

Model-based Analysis of ChIP-Seq (MACS) (Zhang, *et al.*, 2008) is a tool for analyzing short reads for the spatial resolution of the predicted sites, capturing local biases in the genome and generation of peaks with detailed information about length, genome coordinates, summit, p-value, q-values, false-discovery rate (FDR) and fold enrichment. MACS2's used version details, including input, output and download details are given in S-Table. 3, Row No.: 8.

No.	Applications	Versions	Download Web links	Input File Formats	Outputs File Formats
1	FASTQC	0.11.2	<a href="http://www.bioinformatics.braham.ac.uk/projects/fastqc/">http://www.bioinformatics.braham.ac.uk/projects/fastqc/</a>	FASTQ	html, zip, txt.
2	Trimmomatic (Bolger <i>et al.</i> , 2014)	0.32	<a href="http://www.usadellab.org/cms/?page=trimmomatic">http://www.usadellab.org/cms/?page=trimmomatic</a>	FASTQ	FASTQ, trimU.fastq,
3	Burrows-Wheeler Alignment tool (BWA) (Li and Durbin, 2009)	0.7.10	<a href="http://bio-bwa.sourceforge.net">http://bio-bwa.sourceforge.net</a>	fastq_filtered	sam
4	Sequence Alignment/Map (SAM) tools (Li, 2011; Li, <i>et al.</i> , 2009)	0.1.19	<a href="http://samtools.sourceforge.net">http://samtools.sourceforge.net</a>	sam	Sam. txt, pdf
5	Picard	1.95	<a href="http://broadinstitute.github.io">http://broadinstitute.github.io</a>	sam, bam	sam, bam

			o/picard/		
6	Browser Extensible Data (BED) tools (Quinlan and Hall, 2010)	2.22.0	http://bedtools.readthedocs.org/en/latest/content/overview.html	bam	bed
7	ATAC_BAM_shifter_gappedAlign	1	https://github.com/acdaugherly/scripts/blob/master/MostUsed/ATAC_BAM_add1.pl	bam	bam
8	Model-based Analysis of ChIP-Seq (MACS) (Zhang, <i>et al.</i> , 2008)	2.1.0.20151222	http://liulab.dfci.harvard.edu/MACS/	bed	bed, bdg, broadPeak, gappedPeak

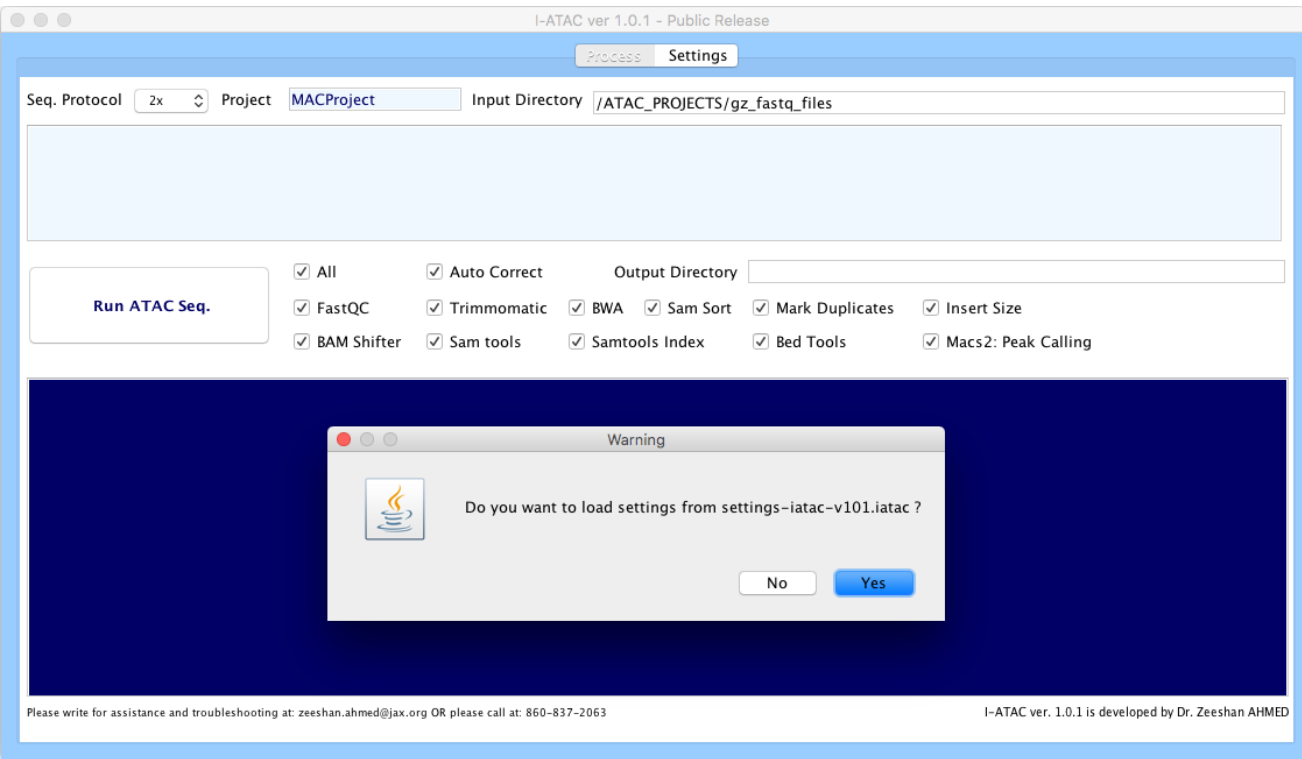
S-Table. 3: Integrated applications in I-ATAC data processing pipeline

No.	Compiler / Interpreter	Versions	Download Web links
1	JAVA	8	https://www.java.com/en/download/
2	R	3.2.3	https://www.r-project.org
3	Perl	5.10.1	https://www.perl.org
4	Python	2.7.3	https://www.python.org

S-Table. 4: Needed compilers and interpreters

## 6 Installation and Configuration

The software executable (JAR file) is open source and freely available and to execute I-ATAC, major requirement is the installation of Java Runtime Environment (<http://www.oracle.com/technetwork/java/javase/downloads/jre8-downloads-2133155.html>).



S-Fig. 9A: I-ATAC, parameters and protocol settings.

## A step-by-step guide to I-ATAC, validating pipeline with two case studies

I-ATAC ver 1.0.1 - Public Release

Process Settings

Reference Genome:

FastQC:

Trimmomatic:

Trim. Adapters:

BWA:

Picard SortSam:

Mark Duplicates:

Insert Size Metrics:

BAM Gap Align:

Samtools:

Bedtools:

MACS:

Output Directory:

Buttons: Save Parameters into File, Load Parameters from File, Default Parameters, Reset Paths, Clear Parameters

Host:

User:

Password:

☒ Multi Queued Jobs

☐ Put in Single Queue

☐ Merge Replicates

Wall time:

nodes:

ppn:

Email:

☒ Create & Queue Jobs

☐ Direct Processing

☒ Soft-Links ☐ Copy

☒ \*.gz - zipped files

S-Fig. 9B: I-ATAC, default setting.

I-ATAC ver 1.0.1 - Public Release

Process Settings

Reference Genome:

FastQC:

Trimmomatic:

Trim. Adapters:

BWA:

Picard SortSam:

Mark Duplicates:

Insert Size Metrics:

BAM Gap Align:

Samtools:

Bedtools:

MACS:

Output Directory:

Buttons: Save Parameters into File, Load Parameters from File, Default Parameters, Reset Paths, Clear Parameters

Host:

User:

Password:

☐ Multi Queued Jobs

☒ Put in Single Queue

☒ Merge Replicates

Wall time:

nodes:

ppn:

Email:

☒ Create & Queue Jobs

☐ Direct Processing

☒ Soft-Links ☐ Copy

☒ \*.gz - zipped files

S-Fig. 9C: I-ATAC, user setting, loaded from "settings-iatac-v101.iatac".



299 After executing I-ATAC and before starting data processing, it is important to set valid applications  
300 paths and calling protocols (section: Graphical User Interface of I-ATAC). Our default parameters  
301 (S-Fig. 9A, 9B and 9C) are set according to our data cluster and installed versions of application (S-  
302 Table. 3), and Compilers/Interpreters (S-Table. 4).

303 Using default configuration settings; I-ATAC will consider logged-in user with a default directory of  
304 same name as of user in the data cluster (e.g. Zeeshan → “d:/data/Zeeshan/ATAC\_PROJECTS/”).  
305 However, user can alter, reset and save default project directory settings.

## 306 **7 Case Studies**

307 In order to validate the performance of I-ATAC and to guide the users, we present two case studies.  
308 First involves using the example data; where we have created small size example dataset (provided in  
309 supplementary material and can be downloaded from the following web link:  
310 <https://zenodo.org/record/46079#.VsJMg7S5LHM>) with artificial names (to explain the process,  
311 execution steps in simpler way.). The reason for giving example study is to let the user, use the  
312 application and observe results in possible shortest time. Moreover, it will also help in figuring out  
313 and resolving trouble shooting conditions (e.g. could be due to inappropriate installation of  
314 downloaded application and compilers/interpreter or any other exceptional reason etc.). Second study  
315 is using publically available data (GM12878, CD4); where we have processed publically available  
316 data, which a trained user can download and process using I-ATAC. In both case studies, I-ATAC is  
317 run at the Mac-OS-X-Yosemite 10.10.5 platform.

### 318 **7.1 Example Dataset**

#### 319 **7.1.1 Dataset Details**

320 Raw dataset and produced results mentioned in this example case study, which can be downloaded  
321 from the provided project web link. Sequenced, paired sample data (“FASTQ” or “FASTQ.gz”) files  
322 are need to be collected and placed in the attached data cluster.

#### 323 **7.1.2 Input**

324 The input to I-ATAC is the path to ATAC-seq sample data, which in our case is:

325 *“/data/zahmed/ATAC\_PROJECTS/gz\_fastq\_files”*

326 As shown in S-Fig. 10, there are two samples available (paired data, four “FASTQ” zipped files) in  
327 the above-mentioned directory i.e. “gz\_fastq\_files”, which are:

328 *Firt\_SampleData\_R1.fastq.gz*  
329 *Firt\_SampleData\_R2.fastq.gz*  
330 *Second\_SampleData\_R1.fastq.gz*  
331 *Second\_SampleData\_R2.fastq.gz*

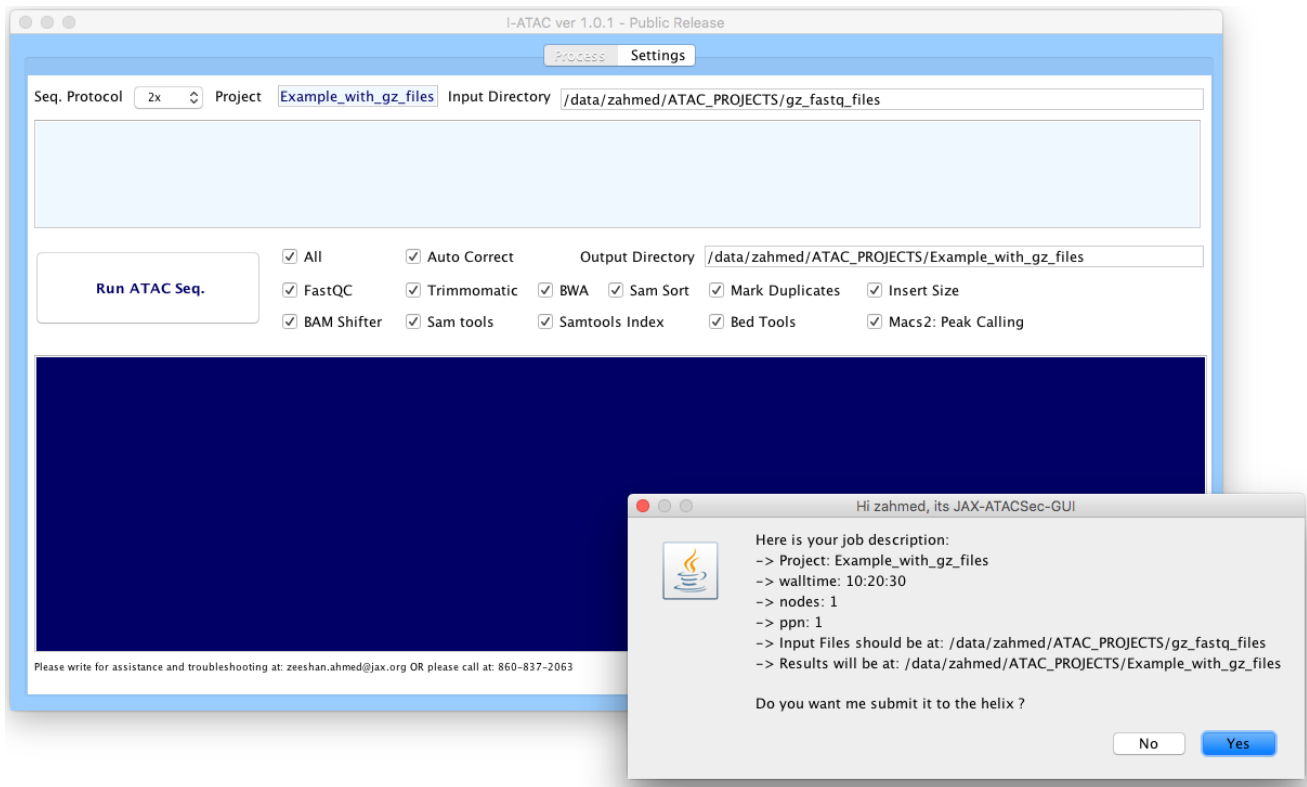
## A step-by-step guide to I-ATAC, validating pipeline with two case studies

```
zahmed — zahmed@helix:/data/zahmed/ATAC_PROJECTS/gz_fastq_files — ss...
[zahmed@helix gz_fastq_files]$ pwd
/data/zahmed/ATAC_PROJECTS/gz_fastq_files
[zahmed@helix gz_fastq_files]$ ls -l
total 1558
-rw-r--r-- 1 zahmed jaxuser 126452 Jan 13 14:17 Firt_SampleData_R1.fastq.gz
-rw-r--r-- 1 zahmed jaxuser 131611 Jan 13 14:17 Firt_SampleData_R2.fastq.gz
-rw-r--r-- 1 zahmed jaxuser 126454 Jan 13 14:17 Second_SampleData_R1.fastq.gz
-rw-r--r-- 1 zahmed jaxuser 131613 Jan 13 14:17 Second_SampleData_R2.fastq.gz
[zahmed@helix gz_fastq_files]$
```

332

333 S-Fig. 10: Screen shot (Linux Terminal, using Mac-OS-X) of compressed sample data files

334 After setting parameters and input path to the I-ATAC-seq, pressed button “Run ATAC-seq”, an  
335 information message will appear (S-Fig. 11) to verify the input sample data source location, output  
336 directory location and set job parameters.

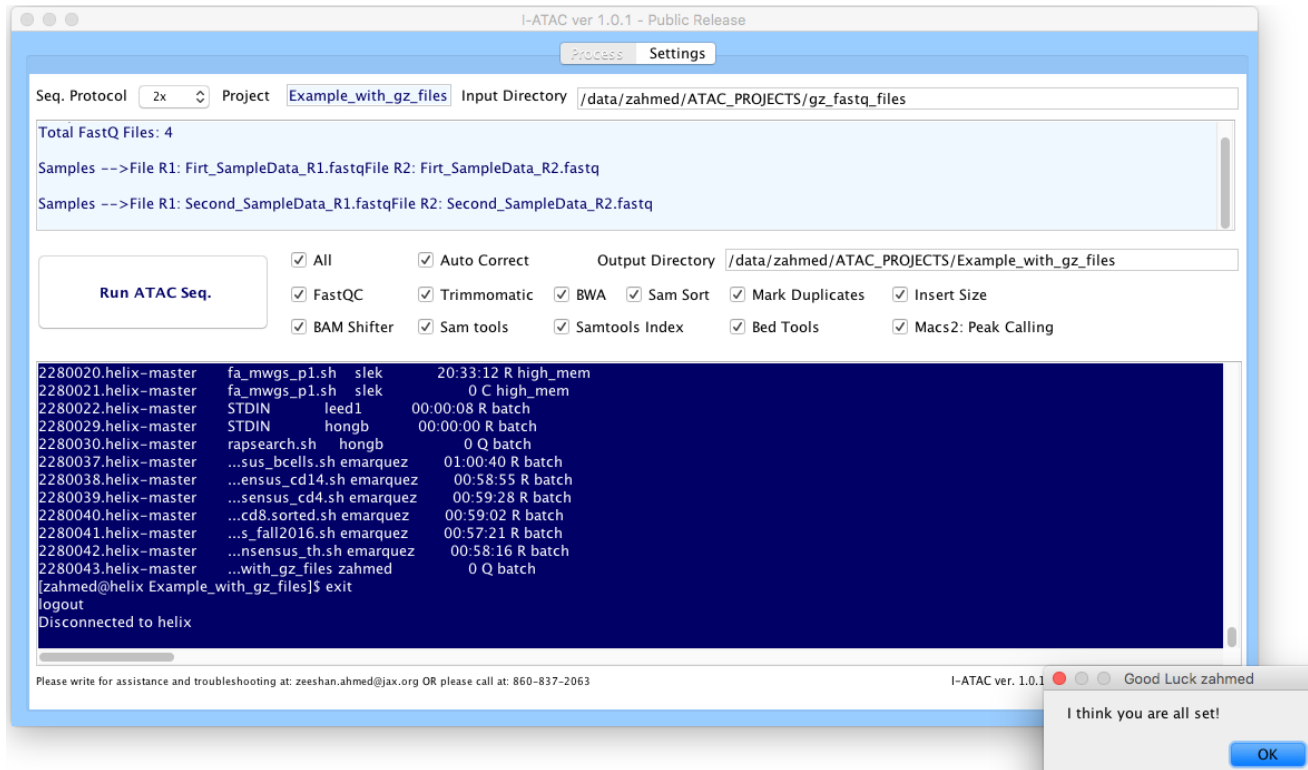


337

338 S-Fig. 11: I-ATAC, input sample data and set parameters' verification



## A step-by-step guide to I-ATAC, validating pipeline with two case studies



S-Fig. 14: I-ATAC job queued.

### 7.1.3 Output

After the successful execution of the ATAC-seq data processing pipeline, the system's generated output can be located in the mentioned output directory (S-Fig. 15). The project directory contains automatically generated and run script:

*Example\_with\_gz\_files\_Example\_with\_gz\_files.sh*

copied, pasted and unzipped "FASTQ" files:

*Firt\_SampleData\_R1.fastq*  
*Firt\_SampleData\_R2.fastq*  
*Second\_SampleData\_R1.fastq*  
*Second\_SampleData\_R2.fastq*

and sub-directories:

*Example\_with\_gz\_files\_Firt\_SampleData\_R1\_Firt\_SampleData\_R2,*  
*Example\_with\_gz\_files\_Second\_SampleData\_R1\_Second\_SampleData\_R2,*  
*MergedSamples*

The output files were placed in proposed system's automatically created sub-directory structure (Section: Applications integration, data processing pipeline and project's directory structure), as shown in S-Fig. 14. We also input two samples and asked system to produce merged replicates as well. So, we observed results for both samples as well as merged replicates.

```

[ zahmed@helix Example_with_gz_files ]$ ls -l
total 4153
-rw-r--r-- 1 zahmed jakuser 29832 Jan 13 17:20 Example_with_gz_files.e981108
-rw-r--r-- 1 zahmed jakuser 8998 Jan 13 17:19 Example_with_gz_files.Example_with_gz_files.sh
drwxr-sr-x 4 zahmed jakuser 53 Jan 13 17:19 Example_with_gz_files.First_SampleData_R1_First_SampleData_R2
-rw-r--r-- 1 zahmed jakuser 192 Jan 13 17:19 Example_with_gz_files.o981108
drwxr-sr-x 4 zahmed jakuser 53 Jan 13 17:19 Example_with_gz_files.Second_SampleData_R1_Second_SampleData_R2
-rw-r--r-- 1 zahmed jakuser 714798 Jan 13 17:18 First_SampleData_R1.fastq
-rw-r--r-- 1 zahmed jakuser 714798 Jan 13 17:18 First_SampleData_R2.fastq
drwxr-sr-x 3 zahmed jakuser 328 Jan 13 17:20 MergedSamples
-rw-r--r-- 1 zahmed jakuser 714798 Jan 13 17:18 Second_SampleData_R1.fastq
-rw-r--r-- 1 zahmed jakuser 714798 Jan 13 17:18 Second_SampleData_R2.fastq
[ zahmed@helix Example_with_gz_files ]$

```

S-Fig. 15: Screen shot (Linux Terminal, using Mac-OS-X) of produced I-ATAC output project directory and files

The produced results from First\_SampleData are shown in S-Fig. 16, including quality reports in FASTQC directory, which are:

```

First_SampleData_R1_fastqc.html
First_SampleData_R1_fastqc.zip
First_SampleData_R2_fastqc.html
First_SampleData_R2_fastqc.zip

```

trimmed “FASTQ” files in trimmomatic directory, which are:

```

First_SampleData_R1.fastq_filtered
First_SampleData_R1.trimU.fastq
First_SampleData_R2.fastq_filtered
First_SampleData_R2.trimU.fastq

```

all sorted, shifted and indexed “sam” and “bam”, “bed” and related files are placed in “bwa” directory, which are:

```

Example_with_gz_files_First_SampleData_R1_First_SampleData_R2.sam
Example_with_gz_files_First_SampleData_R1_First_SampleData_R2_sorted.sam
Example_with_gz_files_First_SampleData_R1_First_SampleData_R2_rmdup.sam
Example_with_gz_files_First_SampleData_R1_First_SampleData_R2_rmdup_shifted.bam
Example_with_gz_files_First_SampleData_R1_First_SampleData_R2_rmdup_shifted_sorted.bam
Example_with_gz_files_First_SampleData_R1_First_SampleData_R2_rmdup_shifted_sorted.bam.bai
Example_with_gz_files_First_SampleData_R1_First_SampleData_R2_rmdup_shifted_sorted.bam.sorted.bed
Example_with_gz_files_First_SampleData_R1_First_SampleData_R2_rmdup_metrics.txt
Example_with_gz_files_First_SampleData_R1_First_SampleData_R2_rmdup_metrics.pdf
Example_with_gz_files_First_SampleData_R1_First_SampleData_R2_rmdup_insertSize.txt

```

and all produced results at peak calling were placed in “macs2” directory, which are:

```

Example_with_gz_files_First_SampleData_R1_First_SampleData_R2_rmdup_shifted_sorted.bam.sorted_control_lambda
.bdg
Example_with_gz_files_First_SampleData_R1_First_SampleData_R2_rmdup_shifted_sorted.bam.sorted_peaks.broadPeak
Example_with_gz_files_First_SampleData_R1_First_SampleData_R2_rmdup_shifted_sorted.bam.sorted_peaks.gappedPeak
Example_with_gz_files_First_SampleData_R1_First_SampleData_R2_rmdup_shifted_sorted.bam.sorted_peaks.xls
Example_with_gz_files_First_SampleData_R1_First_SampleData_R2_rmdup_shifted_sorted.bam.sorted_treat_pileup.bdg
g

```



## A step-by-step guide to I-ATAC, validating pipeline with two case studies

```
zahmed@helix:~/data/zahmed/ATAC_PROJECTS/Example_with_gz_files/Example_with_gz_files_Firt_SampleData_R1_Firt_SampleData_R2/trimmomatic/bwa/mac2 -- ssh -- 183x45
[zahmed@helix Example_with_gz_files]$ cd Example_with_gz_files_Firt_SampleData_R1_Firt_SampleData_R2/
[zahmed@helix Example_with_gz_files_Firt_SampleData_R1_Firt_SampleData_R2]$ ls -l
total 4
drwxr-sr-x 2 zahmed jaxuser 198 Jan 13 17:19 fastQC
drwxr-sr-x 3 zahmed jaxuser 219 Jan 13 17:19 trimmomatic
[zahmed@helix Example_with_gz_files_Firt_SampleData_R1_Firt_SampleData_R2]$ cd fastQC/
[zahmed@helix fastQC]$ ls -l
total 2174
-rw-r--r-- 1 zahmed jaxuser 259607 Jan 13 17:19 Firt_SampleData_R1_fastqc.html
-rw-r--r-- 1 zahmed jaxuser 273565 Jan 13 17:19 Firt_SampleData_R1_fastqc.zip
-rw-r--r-- 1 zahmed jaxuser 256000 Jan 13 17:19 Firt_SampleData_R2_fastqc.html
-rw-r--r-- 1 zahmed jaxuser 267338 Jan 13 17:19 Firt_SampleData_R2_fastqc.zip
[zahmed@helix fastQC]$ cd ..
[zahmed@helix Example_with_gz_files_Firt_SampleData_R1_Firt_SampleData_R2]$ cd trimmomatic/
[zahmed@helix trimmomatic]$ ls -l
total 2032
drwxr-sr-x 3 zahmed jaxuser 983 Jan 13 17:19 bwa
-rw-r--r-- 1 zahmed jaxuser 694814 Jan 13 17:19 Firt_SampleData_R1_fastq_filtered
-rw-r--r-- 1 zahmed jaxuser 15594 Jan 13 17:19 Firt_SampleData_R1_trimU_fastq
-rw-r--r-- 1 zahmed jaxuser 694662 Jan 13 17:19 Firt_SampleData_R2_fastq_filtered
-rw-r--r-- 1 zahmed jaxuser 3374 Jan 13 17:19 Firt_SampleData_R2_trimU_fastq
[zahmed@helix trimmomatic]$ cd bwa/
[zahmed@helix bwa]$ ls -l
total 7497
-rw-r--r-- 1 zahmed jaxuser 4365 Jan 13 17:19 Example_with_gz_files_Firt_SampleData_R1_Firt_SampleData_R2_rmdup_insertSize.txt
-rw-r--r-- 1 zahmed jaxuser 7993 Jan 13 17:19 Example_with_gz_files_Firt_SampleData_R1_Firt_SampleData_R2_rmdup_metrics.pdf
-rw-r--r-- 1 zahmed jaxuser 2672 Jan 13 17:19 Example_with_gz_files_Firt_SampleData_R1_Firt_SampleData_R2_rmdup_metrics.txt
-rw-r--r-- 1 zahmed jaxuser 1165195 Jan 13 17:19 Example_with_gz_files_Firt_SampleData_R1_Firt_SampleData_R2_rmdup.sam
-rw-r--r-- 1 zahmed jaxuser 229248 Jan 13 17:19 Example_with_gz_files_Firt_SampleData_R1_Firt_SampleData_R2_rmdup_shifted.bam
-rw-r--r-- 1 zahmed jaxuser 229028 Jan 13 17:19 Example_with_gz_files_Firt_SampleData_R1_Firt_SampleData_R2_rmdup_shifted_sorted.bam
-rw-r--r-- 1 zahmed jaxuser 1481296 Jan 13 17:19 Example_with_gz_files_Firt_SampleData_R1_Firt_SampleData_R2_rmdup_shifted_sorted.bam.bai
-rw-r--r-- 1 zahmed jaxuser 281340 Jan 13 17:19 Example_with_gz_files_Firt_SampleData_R1_Firt_SampleData_R2_rmdup_shifted_sorted.bam.sorted.bed
-rw-r--r-- 1 zahmed jaxuser 1070086 Jan 13 17:19 Example_with_gz_files_Firt_SampleData_R1_Firt_SampleData_R2_rmdup_shifted_sorted.bam.sorted.peaks.gappedPeak
-rw-r--r-- 1 zahmed jaxuser 1070031 Jan 13 17:19 Example_with_gz_files_Firt_SampleData_R1_Firt_SampleData_R2_rmdup_shifted_sorted.bam.sorted.peaks.xls
drwxr-sr-x 2 zahmed jaxuser 624 Jan 13 17:19 macs2
[zahmed@helix bwa]$ cd macs2/
[zahmed@helix macs2]$ ls -l
total 944
-rw-r--r-- 1 zahmed jaxuser 176901 Jan 13 17:19 Example_with_gz_files_Firt_SampleData_R1_Firt_SampleData_R2_rmdup_shifted_sorted.bam.sorted_control_lambda.bdg
-rw-r--r-- 1 zahmed jaxuser 0 Jan 13 17:19 Example_with_gz_files_Firt_SampleData_R1_Firt_SampleData_R2_rmdup_shifted_sorted.bam.sorted_peaks.broadPeak
-rw-r--r-- 1 zahmed jaxuser 0 Jan 13 17:19 Example_with_gz_files_Firt_SampleData_R1_Firt_SampleData_R2_rmdup_shifted_sorted.bam.sorted_peaks.gappedPeak
-rw-r--r-- 1 zahmed jaxuser 1427 Jan 13 17:19 Example_with_gz_files_Firt_SampleData_R1_Firt_SampleData_R2_rmdup_shifted_sorted.bam.sorted_peaks.xls
-rw-r--r-- 1 zahmed jaxuser 200932 Jan 13 17:19 Example_with_gz_files_Firt_SampleData_R1_Firt_SampleData_R2_rmdup_shifted_sorted.bam.sorted_treat_pileup.bdg
[zahmed@helix macs2]$
```

S-Fig. 16: Screen shot (Linux Terminal, using Mac-OS-X) of produced I-ATAC output files for Firt\_SampleData

Likewise, First\_SampleData, the produced results from Second\_SampleData are shown in S-Fig. 17.

```
zahmed@helix:~/data/zahmed/ATAC_PROJECTS/Example_with_gz_files/Example_with_gz_files_Second_SampleData_R1_Second_SampleData_R2/trimmomatic/bwa/mac2 -- ssh -- 183x45
[zahmed@helix Example_with_gz_files_Second_SampleData_R1_Second_SampleData_R2]$ cd Example_with_gz_files_Second_SampleData_R1_Second_SampleData_R2/
[zahmed@helix Example_with_gz_files_Second_SampleData_R1_Second_SampleData_R2]$ ls -l
total 4
drwxr-sr-x 2 zahmed jaxuser 198 Jan 13 17:19 fastQC
drwxr-sr-x 3 zahmed jaxuser 227 Jan 13 17:19 trimmomatic
[zahmed@helix Example_with_gz_files_Second_SampleData_R1_Second_SampleData_R2]$ cd fastQC/
[zahmed@helix fastQC]$ ls -l
total 2174
-rw-r--r-- 1 zahmed jaxuser 259613 Jan 13 17:19 Second_SampleData_R1_fastqc.html
-rw-r--r-- 1 zahmed jaxuser 273652 Jan 13 17:19 Second_SampleData_R1_fastqc.zip
-rw-r--r-- 1 zahmed jaxuser 256006 Jan 13 17:19 Second_SampleData_R2_fastqc.html
-rw-r--r-- 1 zahmed jaxuser 267428 Jan 13 17:19 Second_SampleData_R2_fastqc.zip
[zahmed@helix fastQC]$ cd ..
[zahmed@helix Example_with_gz_files_Second_SampleData_R1_Second_SampleData_R2]$ cd trimmomatic/
[zahmed@helix trimmomatic]$ ls -l
total 2032
drwxr-sr-x 3 zahmed jaxuser 1023 Jan 13 17:20 bwa
-rw-r--r-- 1 zahmed jaxuser 694814 Jan 13 17:19 Second_SampleData_R1_fastq_filtered
-rw-r--r-- 1 zahmed jaxuser 15594 Jan 13 17:19 Second_SampleData_R1_trimU_fastq
-rw-r--r-- 1 zahmed jaxuser 694662 Jan 13 17:19 Second_SampleData_R2_fastq_filtered
-rw-r--r-- 1 zahmed jaxuser 3374 Jan 13 17:19 Second_SampleData_R2_trimU_fastq
[zahmed@helix trimmomatic]$ cd bwa/
[zahmed@helix bwa]$ ls -l
total 7497
-rw-r--r-- 1 zahmed jaxuser 4377 Jan 13 17:20 Example_with_gz_files_Second_SampleData_R1_Second_SampleData_R2_rmdup_insertSize.txt
-rw-r--r-- 1 zahmed jaxuser 7993 Jan 13 17:20 Example_with_gz_files_Second_SampleData_R1_Second_SampleData_R2_rmdup_metrics.pdf
-rw-r--r-- 1 zahmed jaxuser 2684 Jan 13 17:20 Example_with_gz_files_Second_SampleData_R1_Second_SampleData_R2_rmdup_metrics.txt
-rw-r--r-- 1 zahmed jaxuser 1165121 Jan 13 17:20 Example_with_gz_files_Second_SampleData_R1_Second_SampleData_R2_rmdup.sam
-rw-r--r-- 1 zahmed jaxuser 229251 Jan 13 17:20 Example_with_gz_files_Second_SampleData_R1_Second_SampleData_R2_rmdup_shifted.bam
-rw-r--r-- 1 zahmed jaxuser 229031 Jan 13 17:20 Example_with_gz_files_Second_SampleData_R1_Second_SampleData_R2_rmdup_shifted_sorted.bam
-rw-r--r-- 1 zahmed jaxuser 1481296 Jan 13 17:20 Example_with_gz_files_Second_SampleData_R1_Second_SampleData_R2_rmdup_shifted_sorted.bam.bai
-rw-r--r-- 1 zahmed jaxuser 281340 Jan 13 17:20 Example_with_gz_files_Second_SampleData_R1_Second_SampleData_R2_rmdup_shifted_sorted.bam.sorted.bed
-rw-r--r-- 1 zahmed jaxuser 1070018 Jan 13 17:20 Example_with_gz_files_Second_SampleData_R1_Second_SampleData_R2_rmdup_shifted_sorted.bam.sorted.peaks.gappedPeak
-rw-r--r-- 1 zahmed jaxuser 1070035 Jan 13 17:20 Example_with_gz_files_Second_SampleData_R1_Second_SampleData_R2_rmdup_shifted_sorted.bam.sorted.peaks.xls
drwxr-sr-x 2 zahmed jaxuser 644 Jan 13 17:20 macs2
[zahmed@helix bwa]$ cd macs2/
[zahmed@helix macs2]$ ls -l
total 944
-rw-r--r-- 1 zahmed jaxuser 176901 Jan 13 17:20 Example_with_gz_files_Second_SampleData_R1_Second_SampleData_R2_rmdup_shifted_sorted.bam.sorted_control_lambda.bdg
-rw-r--r-- 1 zahmed jaxuser 0 Jan 13 17:20 Example_with_gz_files_Second_SampleData_R1_Second_SampleData_R2_rmdup_shifted_sorted.bam.sorted_peaks.broadPeak
-rw-r--r-- 1 zahmed jaxuser 0 Jan 13 17:20 Example_with_gz_files_Second_SampleData_R1_Second_SampleData_R2_rmdup_shifted_sorted.bam.sorted_peaks.gappedPeak
-rw-r--r-- 1 zahmed jaxuser 1447 Jan 13 17:20 Example_with_gz_files_Second_SampleData_R1_Second_SampleData_R2_rmdup_shifted_sorted.bam.sorted_peaks.xls
-rw-r--r-- 1 zahmed jaxuser 200932 Jan 13 17:20 Example_with_gz_files_Second_SampleData_R1_Second_SampleData_R2_rmdup_shifted_sorted.bam.sorted_treat_pileup.bdg
[zahmed@helix macs2]$
```

S-Fig. 17: Screen shot (Linux Terminal, using Mac-OS-X) of produced I-ATAC output files for Second\_SampleData

The produced results from merged samples are shown in S-Fig. 18.

```

zahmed — zahmed@helix:/data/zahmed/ATAC_PROJECTS/Example_with_gz_files/MergedSamples/mac2 — ssh — 183x19
[zahmed@helix Example_with_gz_files]$ cd MergedSamples/
[zahmed@helix MergedSamples]$ ls -l
total 5866
-rw-r--r-- 1 zahmed jaxuser 252113 Jan 13 17:20 Example_with_gz_files_mergedSample.bam
-rw-r--r-- 1 zahmed jaxuser 1973149 Jan 13 17:20 Example_with_gz_files_mergedSample.sam
-rw-r--r-- 1 zahmed jaxuser 252113 Jan 13 17:20 Example_with_gz_files_mergedSample_sorted.bam
-rw-r--r-- 1 zahmed jaxuser 1481296 Jan 13 17:20 Example_with_gz_files_mergedSample_sorted.bam.bai
-rw-r--r-- 1 zahmed jaxuser 562680 Jan 13 17:20 Example_with_gz_files_mergedSample_sorted.bed
drwxr-xr-x 2 zahmed jaxuser 374 Jan 13 17:20 macs2
[zahmed@helix MergedSamples]$ cd macs2/
[zahmed@helix macs2]$ ls -l
total 968
-rw-r--r-- 1 zahmed jaxuser 176901 Jan 13 17:20 Example_with_gz_files_mergedSample_sorted_control_lambda.bdg
-rw-r--r-- 1 zahmed jaxuser 0 Jan 13 17:20 Example_with_gz_files_mergedSample_sorted_peaks.broadPeak
-rw-r--r-- 1 zahmed jaxuser 0 Jan 13 17:20 Example_with_gz_files_mergedSample_sorted_peaks.gappedPeak
-rw-r--r-- 1 zahmed jaxuser 1165 Jan 13 17:20 Example_with_gz_files_mergedSample_sorted_peaks.xls
-rw-r--r-- 1 zahmed jaxuser 200932 Jan 13 17:20 Example_with_gz_files_mergedSample_sorted_treat_pileup.bdg
[zahmed@helix macs2]$

```

S-Fig. 18: Screen shot (Linux Terminal, using Mac-OS-X) of produced I-ATAC output files for Merged Samples

## 7.2 Case Study 2: Using GM12878 – CD4 T- Cells

### 7.2.1 Dataset Details

Information about Raw dataset (GM12878 – CD4 T- Cell, Day 1, Rep1 SRR891275 and Rep2 SRR891276) is available at web link [https://catalog.coriell.org/0/Sections/Search/Sample\\_Detail.aspx?Ref=GM12878&product=CC](https://catalog.coriell.org/0/Sections/Search/Sample_Detail.aspx?Ref=GM12878&product=CC) and produced results, which are mentioned in this case study can be downloaded from the following project web link.

### 7.2.2 Input

The input to I-ATAC is the path to ATAC-seq sample data (S-Fig. 19), which in our case is:

`"/data/zahmed/ATAC_seq_data/CD4"`

```

zahmed — zahmed@helix:/data/zahmed/ATAC_seq_data/CD4 — ssh — 111x10
[zahmed@helix CD4]$ pwd
/data/zahmed/ATAC_seq_data/CD4
[zahmed@helix CD4]$ ls -l
total 2448782
-rwxr-xr-x 1 zahmed jaxuser 469964798 Jan 15 14:49 CD4+_ATACseq_Day1_Rep1_SRR891275_R1_ALL.fastq.gz
-rwxr-xr-x 1 zahmed jaxuser 458991263 Jan 15 14:49 CD4+_ATACseq_Day1_Rep1_SRR891275_R2_ALL.fastq.gz
-rwxr-xr-x 1 zahmed jaxuser 653655528 Jan 15 14:50 CD4+_ATACseq_Day1_Rep2_SRR891276_R1_ALL.fastq.gz
-rwxr-xr-x 1 zahmed jaxuser 640307650 Jan 15 14:50 CD4+_ATACseq_Day1_Rep2_SRR891276_R2_ALL.fastq.gz
[zahmed@helix CD4]$

```

S-Fig. 19: Screen shot (Linux Terminal, using Mac-OS-X) of compressed sample data files of GM12878 – CD4 T- Cell, Day 1, Rep1 SRR891275 and Rep2 SRR891276.

Likewise, earlier discussed case study, at successful identification and verification of sample data files (S-Fig. 20), data processing job was created (S-Fig. 21) and successfully queued (S-Fig. 22).



24



## 439 7.2.3 Output

440 As in the earlier discussed case study, all the produced results were placed in the proposed and auto  
441 generated directory structure (S-Fig. 22, 23and 24).

```

zahmed@helix GM12878_CD4_Day1$ ls -l
total 12489801
-rw-r--r-- 1 zahmed jakuser 2351167270 Jan 15 14:52 CD4+_ATACseq_Day1_Rep1_SRR891275_R1_ALL.fastq
-rw-r--r-- 1 zahmed jakuser 2351167270 Jan 15 14:52 CD4+_ATACseq_Day1_Rep1_SRR891275_R2_ALL.fastq
-rw-r--r-- 1 zahmed jakuser 3279168866 Jan 15 14:52 CD4+_ATACseq_Day1_Rep2_SRR891276_R1_ALL.fastq
-rw-r--r-- 1 zahmed jakuser 3279168866 Jan 15 14:52 CD4+_ATACseq_Day1_Rep2_SRR891276_R2_ALL.fastq
drwxr-sr-x 4 zahmed jakuser 53 Jan 15 14:55 GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep1_SRR891275_R1_ALL_CD4+_ATACseq_Day1_Rep1_SRR891275_R2_ALL
drwxr-sr-x 4 zahmed jakuser 53 Jan 15 15:49 GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep2_SRR891276_R1_ALL_CD4+_ATACseq_Day1_Rep2_SRR891276_R2_ALL
-rw-r--r-- 1 zahmed jakuser 4966515 Jan 15 17:03 GM12878_CD4_Day1.o985775
-rw-r--r-- 1 zahmed jakuser 18964 Jan 15 14:55 GM12878_CD4_Day1_GM12878_CD4_Day1.sh
-rw-r--r-- 1 zahmed jakuser 272 Jan 15 15:50 GM12878_CD4_Day1.o985775
drwxr-sr-x 3 zahmed jakuser 383 Jan 15 17:00 MergedSamples
[zahmed@helix GM12878_CD4_Day1]$ cd GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep1_SRR891275_R1_ALL_CD4+_ATACseq_Day1_Rep1_SRR891275_R2_ALL
[zahmed@helix GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep1_SRR891275_R1_ALL_CD4+_ATACseq_Day1_Rep1_SRR891275_R2_ALL]$ cd fastQC/
[zahmed@helix fastQC]$ ls -l
total 2878
-rw-r--r-- 1 zahmed jakuser 383986 Jan 15 14:56 CD4+_ATACseq_Day1_Rep1_SRR891275_R1_ALL_fastqc.html
-rw-r--r-- 1 zahmed jakuser 518046 Jan 15 14:56 CD4+_ATACseq_Day1_Rep1_SRR891275_R1_ALL_fastqc.zip
-rw-r--r-- 1 zahmed jakuser 376161 Jan 15 14:56 CD4+_ATACseq_Day1_Rep1_SRR891275_R2_ALL_fastqc.html
-rw-r--r-- 1 zahmed jakuser 584211 Jan 15 14:56 CD4+_ATACseq_Day1_Rep1_SRR891275_R2_ALL_fastqc.zip
[zahmed@helix fastQC]$ cd ..
[zahmed@helix GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep1_SRR891275_R1_ALL_CD4+_ATACseq_Day1_Rep1_SRR891275_R2_ALL]$ cd trimmomatic/
[zahmed@helix trimmomatic]$ ls -l
total 5019864
drwxr-sr-x 3 zahmed jakuser 1353 Jan 15 15:47 bwa
-rw-r--r-- 1 zahmed jakuser 2222603420 Jan 15 14:58 CD4+_ATACseq_Day1_Rep1_SRR891275_R1_ALL_fastq_filtered
-rw-r--r-- 1 zahmed jakuser 61380938 Jan 15 14:58 CD4+_ATACseq_Day1_Rep1_SRR891275_R1_ALL_trim0.fastq
-rw-r--r-- 1 zahmed jakuser 2225050298 Jan 15 14:58 CD4+_ATACseq_Day1_Rep1_SRR891275_R2_ALL_fastq_filtered
-rw-r--r-- 1 zahmed jakuser 49900792 Jan 15 14:58 CD4+_ATACseq_Day1_Rep1_SRR891275_R2_ALL_trim0.fastq
[zahmed@helix trimmomatic]$ cd bwa/
[zahmed@helix bwa]$ ls -l
total 10945305
-rw-r--r-- 1 zahmed jakuser 8621 Jan 15 15:42 GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep1_SRR891275_R1_ALL_CD4+_ATACseq_Day1_Rep1_SRR891275_R2_ALL_rmdup_insertSize.txt
-rw-r--r-- 1 zahmed jakuser 11473 Jan 15 15:42 GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep1_SRR891275_R1_ALL_CD4+_ATACseq_Day1_Rep1_SRR891275_R2_ALL_rmdup_metrics.pdf
-rw-r--r-- 1 zahmed jakuser 2538 Jan 15 15:41 GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep1_SRR891275_R1_ALL_CD4+_ATACseq_Day1_Rep1_SRR891275_R2_ALL_rmdup_metrics.txt
-rw-r--r-- 1 zahmed jakuser 1900204325 Jan 15 15:41 GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep1_SRR891275_R1_ALL_CD4+_ATACseq_Day1_Rep1_SRR891275_R2_ALL_rmdup.sam
-rw-r--r-- 1 zahmed jakuser 378142444 Jan 15 15:45 GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep1_SRR891275_R1_ALL_CD4+_ATACseq_Day1_Rep1_SRR891275_R2_ALL_rmdup_shifted.bam
-rw-r--r-- 1 zahmed jakuser 378832925 Jan 15 15:47 GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep1_SRR891275_R1_ALL_CD4+_ATACseq_Day1_Rep1_SRR891275_R2_ALL_rmdup_shifted_sorted.bam
-rw-r--r-- 1 zahmed jakuser 5847092 Jan 15 15:47 GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep1_SRR891275_R1_ALL_CD4+_ATACseq_Day1_Rep1_SRR891275_R2_ALL_rmdup_shifted_sorted.bam.bai
-rw-r--r-- 1 zahmed jakuser 483741438 Jan 15 15:47 GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep1_SRR891275_R1_ALL_CD4+_ATACseq_Day1_Rep1_SRR891275_R2_ALL_rmdup_shifted_sorted.bam.sorted.bed
-rw-r--r-- 1 zahmed jakuser 3396628083 Jan 15 15:48 GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep1_SRR891275_R1_ALL_CD4+_ATACseq_Day1_Rep1_SRR891275_R2_ALL_rmdup_shifted_sorted.bam.sorted.bed
-rw-r--r-- 1 zahmed jakuser 3396628028 Jan 15 15:39 GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep1_SRR891275_R1_ALL_CD4+_ATACseq_Day1_Rep1_SRR891275_R2_ALL_rmdup_shifted_sorted.bam.sorted.bed
drwxr-sr-x 2 zahmed jakuser 809 Jan 15 15:49 macs2
[zahmed@helix bwa]$ cd macs2/
[zahmed@helix macs2]$ ls -l
total 555680
-rw-r--r-- 1 zahmed jakuser 168991285 Jan 15 15:48 GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep1_SRR891275_R1_ALL_CD4+_ATACseq_Day1_Rep1_SRR891275_R2_ALL_rmdup_shifted_sorted.bam.sorted_control_lambda.bdg
-rw-r--r-- 1 zahmed jakuser 11495026 Jan 15 15:49 GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep1_SRR891275_R1_ALL_CD4+_ATACseq_Day1_Rep1_SRR891275_R2_ALL_rmdup_shifted_sorted.bam.sorted_peaks.broadPeak
-rw-r--r-- 1 zahmed jakuser 13531859 Jan 15 15:49 GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep1_SRR891275_R1_ALL_CD4+_ATACseq_Day1_Rep1_SRR891275_R2_ALL_rmdup_shifted_sorted.bam.sorted_peaks.gappedPeak
-rw-r--r-- 1 zahmed jakuser 11743053 Jan 15 15:48 GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep1_SRR891275_R1_ALL_CD4+_ATACseq_Day1_Rep1_SRR891275_R2_ALL_rmdup_shifted_sorted.bam.sorted_peaks.xls
-rw-r--r-- 1 zahmed jakuser 298027077 Jan 15 15:48 GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep1_SRR891275_R1_ALL_CD4+_ATACseq_Day1_Rep1_SRR891275_R2_ALL_rmdup_shifted_sorted.bam.sorted_treat_pileup.bdg
[zahmed@helix macs2]$

```

442 S-Fig. 22: Screen shot (Linux Terminal, using Mac-OS-X) of produced I-ATAC output main project directory and files,  
443 and for sample CD4+\_ATACseq\_Day1\_Rep1\_SRR891275  
444

```

zahmed@helix GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep2_SRR891276_R1_ALL_CD4+_ATACseq_Day1_Rep2_SRR891276_R2_ALL$ cd fastQC/
[zahmed@helix fastQC]$ ls -l
total 2934
-rw-r--r-- 1 zahmed jakuser 389214 Jan 15 15:50 CD4+_ATACseq_Day1_Rep2_SRR891276_R1_ALL_fastqc.html
-rw-r--r-- 1 zahmed jakuser 524579 Jan 15 15:50 CD4+_ATACseq_Day1_Rep2_SRR891276_R1_ALL_fastqc.zip
-rw-r--r-- 1 zahmed jakuser 390892 Jan 15 15:50 CD4+_ATACseq_Day1_Rep2_SRR891276_R2_ALL_fastqc.html
-rw-r--r-- 1 zahmed jakuser 527829 Jan 15 15:50 CD4+_ATACseq_Day1_Rep2_SRR891276_R2_ALL_fastqc.zip
[zahmed@helix fastQC]$ cd ..
[zahmed@helix GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep2_SRR891276_R1_ALL_CD4+_ATACseq_Day1_Rep2_SRR891276_R2_ALL]$ cd trimmomatic/
[zahmed@helix trimmomatic]$ ls -l
total 6993240
drwxr-sr-x 3 zahmed jakuser 1353 Jan 15 16:51 bwa
-rw-r--r-- 1 zahmed jakuser 3892680856 Jan 15 15:51 CD4+_ATACseq_Day1_Rep2_SRR891276_R1_ALL_fastq_filtered
-rw-r--r-- 1 zahmed jakuser 92162516 Jan 15 15:51 CD4+_ATACseq_Day1_Rep2_SRR891276_R1_ALL_trim0.fastq
-rw-r--r-- 1 zahmed jakuser 3895898034 Jan 15 15:51 CD4+_ATACseq_Day1_Rep2_SRR891276_R2_ALL_fastq_filtered
-rw-r--r-- 1 zahmed jakuser 69102608 Jan 15 15:51 CD4+_ATACseq_Day1_Rep2_SRR891276_R2_ALL_trim0.fastq
[zahmed@helix trimmomatic]$ cd bwa/
[zahmed@helix bwa]$ ls -l
total 14884585
-rw-r--r-- 1 zahmed jakuser 9560 Jan 15 16:44 GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep2_SRR891276_R1_ALL_CD4+_ATACseq_Day1_Rep2_SRR891276_R2_ALL_rmdup_insertSize.txt
-rw-r--r-- 1 zahmed jakuser 12203 Jan 15 16:44 GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep2_SRR891276_R1_ALL_CD4+_ATACseq_Day1_Rep2_SRR891276_R2_ALL_rmdup_metrics.pdf
-rw-r--r-- 1 zahmed jakuser 2723 Jan 15 16:43 GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep2_SRR891276_R1_ALL_CD4+_ATACseq_Day1_Rep2_SRR891276_R2_ALL_rmdup_metrics.txt
-rw-r--r-- 1 zahmed jakuser 2455630664 Jan 15 16:43 GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep2_SRR891276_R1_ALL_CD4+_ATACseq_Day1_Rep2_SRR891276_R2_ALL_rmdup.sam
-rw-r--r-- 1 zahmed jakuser 492114851 Jan 15 16:49 GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep2_SRR891276_R1_ALL_CD4+_ATACseq_Day1_Rep2_SRR891276_R2_ALL_rmdup_shifted.bam
-rw-r--r-- 1 zahmed jakuser 491905411 Jan 15 16:50 GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep2_SRR891276_R1_ALL_CD4+_ATACseq_Day1_Rep2_SRR891276_R2_ALL_rmdup_shifted_sorted.bam
-rw-r--r-- 1 zahmed jakuser 5908968 Jan 15 16:51 GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep2_SRR891276_R1_ALL_CD4+_ATACseq_Day1_Rep2_SRR891276_R2_ALL_rmdup_shifted_sorted.bam.bai
-rw-r--r-- 1 zahmed jakuser 625972032 Jan 15 16:51 GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep2_SRR891276_R1_ALL_CD4+_ATACseq_Day1_Rep2_SRR891276_R2_ALL_rmdup_shifted_sorted.bam.sorted.bed
-rw-r--r-- 1 zahmed jakuser 4721832832 Jan 15 16:36 GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep2_SRR891276_R1_ALL_CD4+_ATACseq_Day1_Rep2_SRR891276_R2_ALL_rmdup_shifted_sorted.bam.sorted.bed
-rw-r--r-- 1 zahmed jakuser 4721832857 Jan 15 16:48 GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep2_SRR891276_R1_ALL_CD4+_ATACseq_Day1_Rep2_SRR891276_R2_ALL_rmdup_shifted_sorted.bam.sorted.bed
drwxr-sr-x 2 zahmed jakuser 809 Jan 15 16:53 macs2
[zahmed@helix bwa]$ cd macs2/
[zahmed@helix macs2]$ ls -l
total 785152
-rw-r--r-- 1 zahmed jakuser 228334931 Jan 15 16:53 GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep2_SRR891276_R1_ALL_CD4+_ATACseq_Day1_Rep2_SRR891276_R2_ALL_rmdup_shifted_sorted.bam.sorted_control_lambda.bdg
-rw-r--r-- 1 zahmed jakuser 8091139 Jan 15 16:53 GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep2_SRR891276_R1_ALL_CD4+_ATACseq_Day1_Rep2_SRR891276_R2_ALL_rmdup_shifted_sorted.bam.sorted_peaks.broadPeak
-rw-r--r-- 1 zahmed jakuser 9480827 Jan 15 16:53 GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep2_SRR891276_R1_ALL_CD4+_ATACseq_Day1_Rep2_SRR891276_R2_ALL_rmdup_shifted_sorted.bam.sorted_peaks.gappedPeak
-rw-r--r-- 1 zahmed jakuser 8278153 Jan 15 16:53 GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep2_SRR891276_R1_ALL_CD4+_ATACseq_Day1_Rep2_SRR891276_R2_ALL_rmdup_shifted_sorted.bam.sorted_peaks.xls
-rw-r--r-- 1 zahmed jakuser 39353911 Jan 15 16:53 GM12878_CD4_Day1_CD4+_ATACseq_Day1_Rep2_SRR891276_R1_ALL_CD4+_ATACseq_Day1_Rep2_SRR891276_R2_ALL_rmdup_shifted_sorted.bam.sorted_treat_pileup.bdg
[zahmed@helix macs2]$

```

445 S-Fig. 23: Screen shot (Linux Terminal, using Mac-OS-X) of produced I-ATAC output files for sample  
446 CD4+\_ATACseq\_Day1\_Rep2\_SRR891276  
447

448

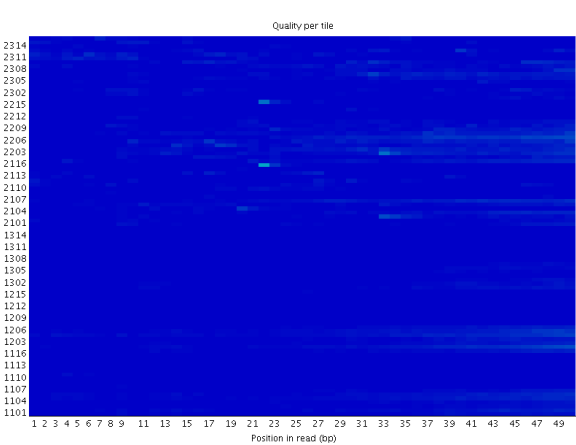
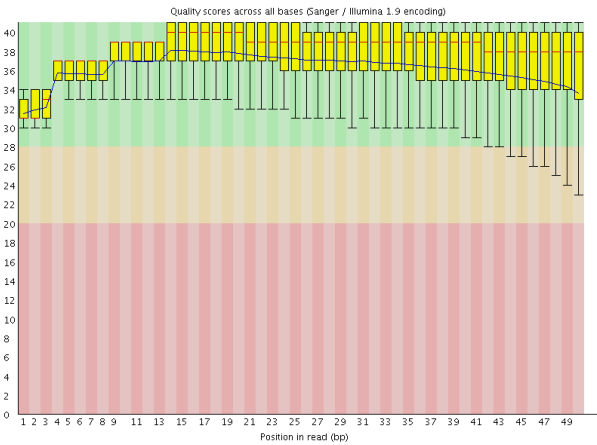
A step-by-step guide to I-ATAC, validating pipeline with two case studies

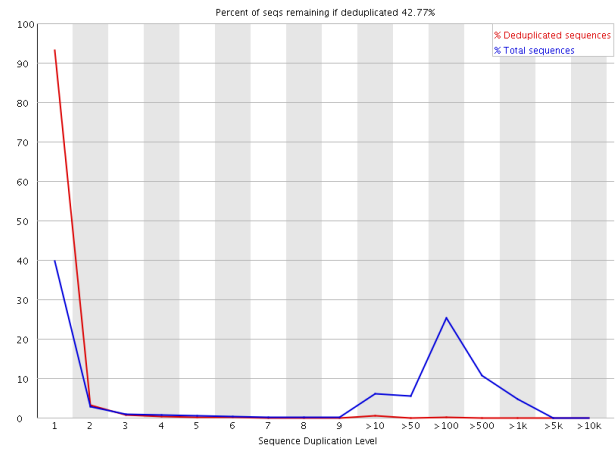
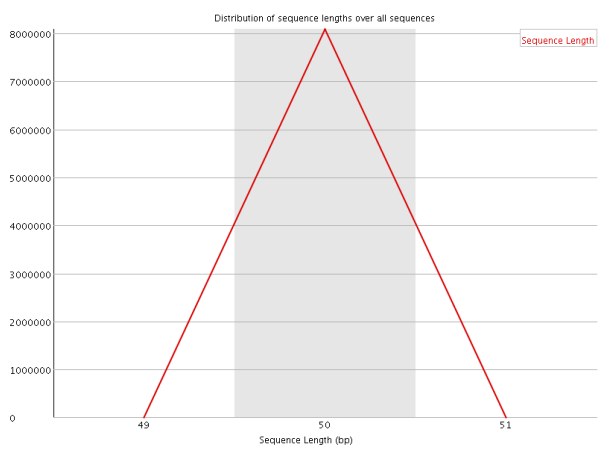
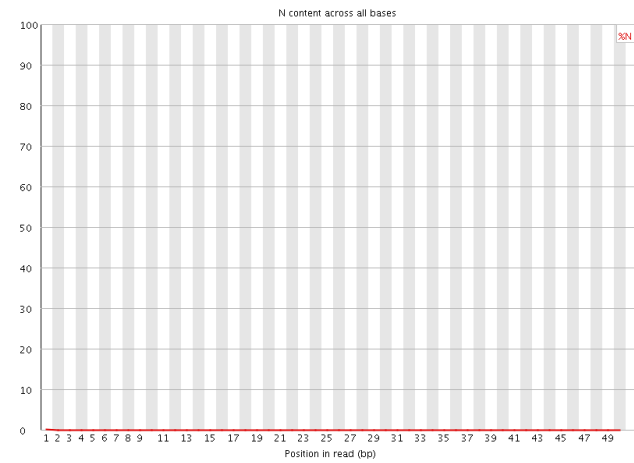
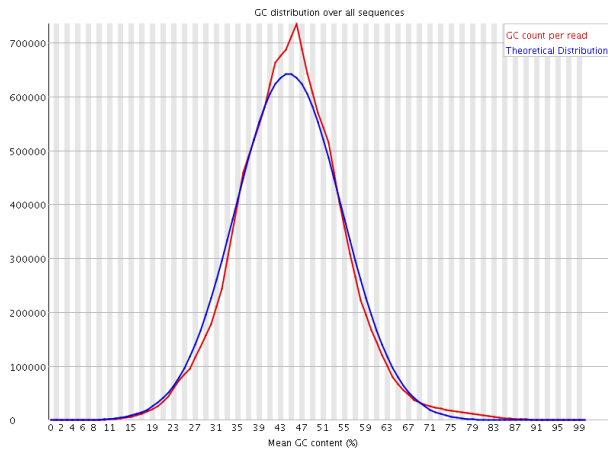
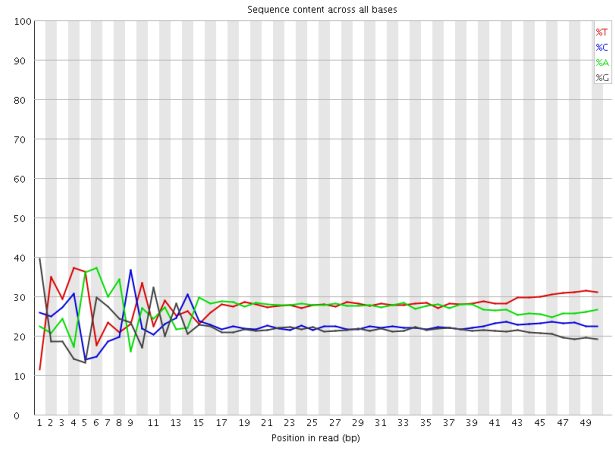
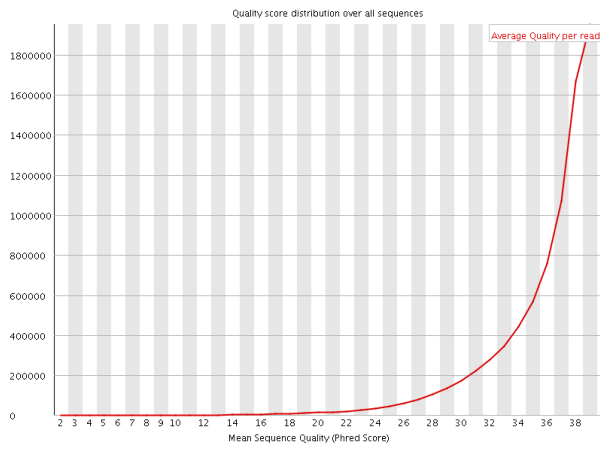
```
zahmed — zahmed@helix:/data/zahmed/ATAC_PROJECTS/GM12878_CD4_Day1/MergedSamples/mac2 — ssh — 142x20
[zahmed@helix GM12878_CD4_Day1]$ cd MergedSamples/
[zahmed@helix MergedSamples]$ ls -l
total 7175522
-rw-r--r-- 1 zahmed jaxuser 861237223 Jan 15 16:55 GM12878_CD4_Day1_mergedSample.bam
-rw-r--r-- 1 zahmed jaxuser 3676571000 Jan 15 17:00 GM12878_CD4_Day1_mergedSample.sam
-rw-r--r-- 1 zahmed jaxuser 861237223 Jan 15 16:59 GM12878_CD4_Day1_mergedSample_sorted.bam
-rw-r--r-- 1 zahmed jaxuser 6057328 Jan 15 16:59 GM12878_CD4_Day1_mergedSample_sorted.bam.bai
-rw-r--r-- 1 zahmed jaxuser 1109713470 Jan 15 17:00 GM12878_CD4_Day1_mergedSample_sorted.bed
drwxr-sr-x 2 zahmed jaxuser 349 Jan 15 17:03 macs2
[zahmed@helix MergedSamples]$ cd macs2/
[zahmed@helix macs2]$ ls -l
total 693456
-rw-r--r-- 1 zahmed jaxuser 206569472 Jan 15 17:02 GM12878_CD4_Day1_mergedSample_sorted_control_lambda.bdg
-rw-r--r-- 1 zahmed jaxuser 5015956 Jan 15 17:03 GM12878_CD4_Day1_mergedSample_sorted_peaks.broadPeak
-rw-r--r-- 1 zahmed jaxuser 6688771 Jan 15 17:03 GM12878_CD4_Day1_mergedSample_sorted_peaks.gappedPeak
-rw-r--r-- 1 zahmed jaxuser 5231644 Jan 15 17:03 GM12878_CD4_Day1_mergedSample_sorted_peaks.xls
-rw-r--r-- 1 zahmed jaxuser 405274624 Jan 15 17:02 GM12878_CD4_Day1_mergedSample_sorted_treat_pileup.bdg
[zahmed@helix macs2]$
```

S-Fig. 24: Screen shot (Linux Terminal, using Mac-OS-X) of produced I-ATAC output files for Merged Samples

Basic Statistics

Measure	Value
Filename	CD4+_ATACseq_Day1_Repl_SRR891275_R1_ALL.fastq
File type	Conventional base calls
Encoding	Sanger / Illumina 1.9
Total Sequences	8086940
Sequences flagged as poor quality	0
Sequence length	50
%GC	44





## 461



(J)

(K)

464

471

471

471

471

471

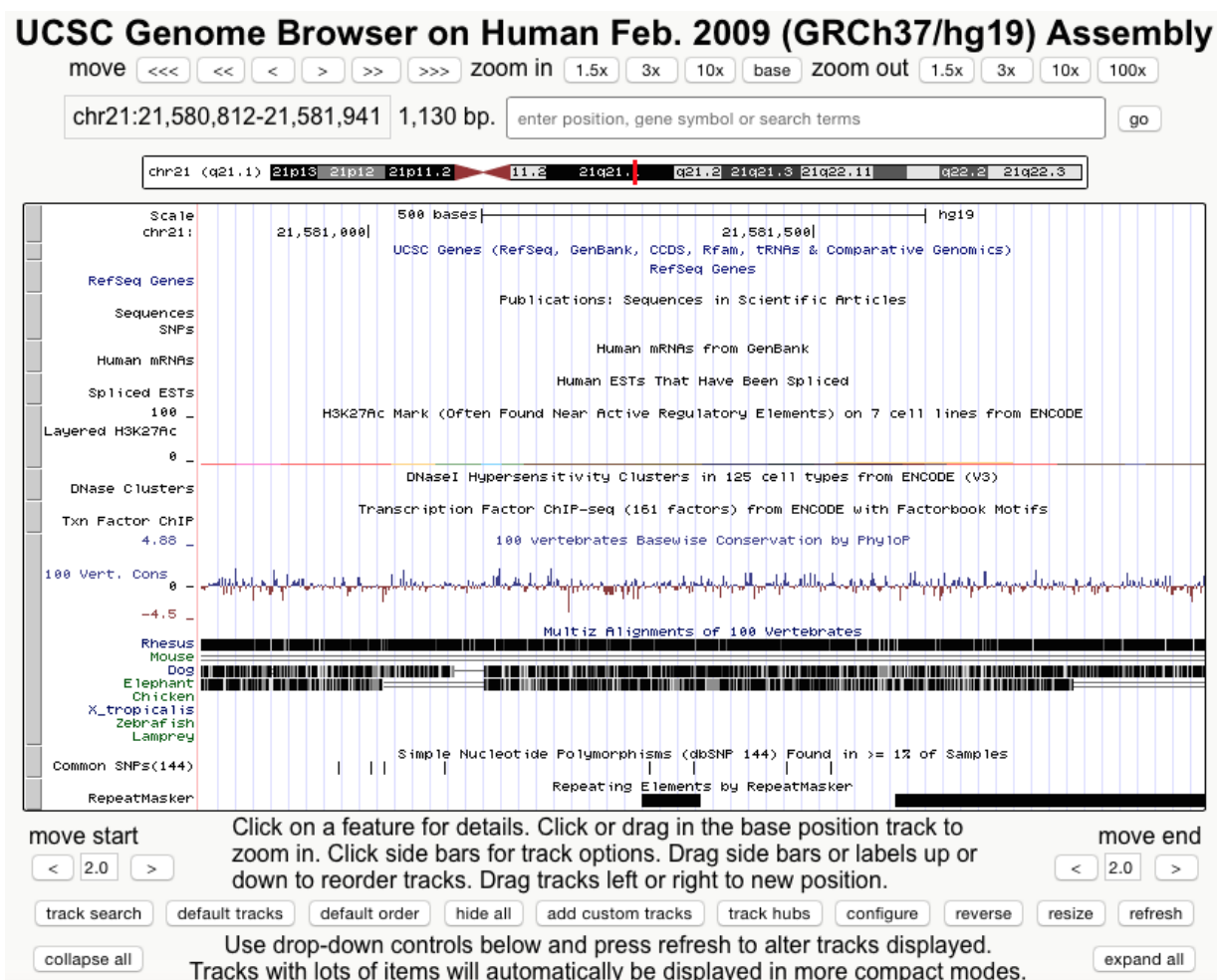
471

471



473

474



S-Fig. 27 Example Visualization using “bdg” file (GM12878\_CD4\_Day1\_mergedSample\_sorted\_control\_lambda.bdg) using USCS Genome Browser.

The detailed output of samples a CD4+\_ATACseq\_Day1\_Rep1\_SRR891275 and CD4+\_ATACseq\_Day1\_Rep2\_SRR891276 are attached in the supplementary material. Example visualization of produced results is created by visualizing sorted „bam“ files (CD4+\_ATACseq\_Day1\_Rep1\_SRR891275, CD4+\_ATACseq\_Day1\_Rep2\_SRR891276 and GM12878\_CD4\_Day1\_mergedSample) using IGV (S-Fig. 26) and peak file (CD4+\_ATACseq\_Day1\_Rep1\_SRR891275) using USCS Genome browser (S-Fig. 27).

## 8 Conclusions

To the best of our knowledge, I-ATAC platform is the first desktop tool that is specialized to processing and analysis of ATAC-seq data. I-ATAC provides a flexible algorithm and parameter setting GUI for non-computational scientists and a time-efficient parallel data analysis environment for computational scientists. Future work includes incorporating visualization and differential analysis modules in I-ATAC platform.

## 492    **9    Acknowledgments**

493    The Jackson Laboratory (JAX) supports and owns this project. Special thanks to Ucar and  
494    Banchereau labs as well as the Genome Technologies and Computational Sciences cores at JAX,  
495    who provided insight and expertise that greatly, assisted the research and development of this  
496    platform.

## 497    **10    Funding:**

498    This work was supported by The Jackson Laboratory, USA.

## 499    **11    Conflict of Interests:**

500    The authors declare that they have no competing interests.

## 501    **12    Additional Requirements**

502    For additional information, please refer to the project webpage: [https://www.jax.org/research-and-](https://www.jax.org/research-and-faculty/tools/i-atac)  
503    [faculty/tools/i-atac](https://www.jax.org/research-and-faculty/tools/i-atac)

504    Source code, JAR files for MAC OS X and Windows, and complete source code package for Eclipse  
505    IDE is available at <https://github.com/UcarLab/I-ATAC>

506    Example dataset is available at: <https://zenodo.org/record/46079#.WAe3l5MrK7Y>

507    Supporting software and dependencies are available at:  
508    <https://zenodo.org/record/162023#.WAe3dJMrK7Y>

## 509    **13    References**

510    Ahmed Z, Ucar D. (2017) A standalone software platform for the interactive management and pre-  
511    processing of ATAC-seq samples. *PeerJ Preprints.*, **5**:e2942v1  
512    <https://doi.org/10.7287/peerj.preprints.2942v1>

513    Ahmed Z, Zeeshan S, Dandekar T. (2014) Developing sustainable software solutions for  
514    bioinformatics by the “Butterfly” paradigm. *F1000Res.*, **3**, 71.

515    Bauch A. (2011) openBIS: a flexible framework for managing and analyzing complex data in  
516    biology research. *BMC Bioinf.*, **12**, 468.

517    Buenrostro JD, Giresi PG, Zaba LC, Chang HY, Greenleaf WJ. *et al.* (2013) Transposition of native  
518    chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-binding proteins and  
519    nucleosome position. *Nat Methods.*, **10**, 1213-1218.

520    Bolger AM, Lohse M, Usadel B. (2014) Trimmomatic: a flexible trimmer for Illumina sequence data.  
521    *Bioinformatics.*, **30**, 2114-20.

522    Dander A, Pabinger S, Sperk M, Fischer M, Stocker G, Trajanoski Z. (2014) SeqBench: integrated  
523    solution for the management and analysis of exome sequencing data. *BMC Res Notes.*, **7**, 43.

524    Giardine B. (2005) Galaxy: a platform for interactive large-scale genome analysis. *Genome Res.* **15**,  
525    1451-1455.

526 Horak CE, Snyder M. (2002) ChIP-chip: A genomic approach for identifying transcription factor  
 527 binding sites. *Methods Enzymol.*, **350**, 469–483  
 528 Li H, Durbin R. (2009) Fast and accurate short read alignment with Burrows-Wheeler Transform.  
 529 *Bioinformatics.*, **25**, 1754-60.  
 530 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R; 1000  
 531 Genome Project Data Processing Subgroup. (2009) The Sequence alignment/map (SAM) format and  
 532 SAMtools. *Bioinformatics.*, **25**, 2078-2079.  
 533 Mariette J, Escudié F, Allias N, Salin G, Noirot C, Thomas S, Klopp C. (2012) NG6: Integrated next  
 534 generation sequencing storage and processing environment. *BMC Genomics.* **13**, 462.  
 535 McLellan A S, Dubin RA, Jing Q, Broin PÓ, Moskowitz D, Suzuki M, Calder RB, Hargitai J,  
 536 Golden A, Greally JM. (2012) The Wasp System: an open source environment for managing and  
 537 analyzing genomic data. *Genomics.*, **100**, 345-51.  
 538 Mount DM. (2004). *Bioinformatics: Sequence and Genome Analysis* (2nd ed.). Cold Spring Harbor  
 539 Laboratory Press: Cold Spring Harbor, NY. ISBN 0-87969-608-7.  
 540 Orvis J, Crabtree J, Galens K, Gussman A, Inman JM, Lee E, Nampally S, Riley D, Sundaram JP,  
 541 Felix V, Whitty B, Mahurkar A, Wortman J, White O, Angiuoli SV. (2010) Ergatis: a web interface  
 542 and scalable software system for bioinformatics workflows. *Bioinformatics*, **26**, 1488-92.  
 543 Quinlan, A R, Hall IM. (2010) BEDTools: a flexible suite of utilities for comparing genomic features.  
 544 *Bioinformatics.*, **26**, 841-842.  
 545 Scholtalbers J. (2013) Galaxy LIMS for next-generation sequencing. *Bioinformatics.* **29**, 1233-1234.  
 546 Tsompana M, Buck MJ. (2014) Chromatin accessibility: a window into the genome. *Epigenetics &*  
 547 *Chromatin.*, **7**, 33.  
 548 Venco F, Vaskin Y, Ceol A, Muller H. (2014) SMITH: a LIMS for handling next-generation  
 549 sequencing workflows. *BMC Bioinf.*, **15**, (Suppl 14):S3.  
 550 Zhang Y, Liu T, Meyer CA, Eeckhoutte J, Johnson DS, Bernstein BE, Nusbaum C, Myers RM, Brown M, Li W, Liu XS.  
 551 (2008) Model-based analysis of ChIP-Seq (MACS). *Genome Biol.*, **9**, R137.  
 552