

Manuscript Number:	GIGA-D-18-00155	
Full Title:	Reproducible genomics analysis pipelines with GNU Guix	
Article Type:	Technical Note	
Funding Information:	German Federal Ministry of Education and Research (BMBF) (031 A538C RBC)	Dr Bora Uyar
	Berlin Institute for Health	Dr Katarzyna Wreczycka
Abstract:	<p>In bioinformatics, as well as other computationally-intensive research fields, there is a need for workflows that can reliably produce consistent output, independent of the software environment or configuration settings of the machine on which they are executed. Indeed, this is essential for controlled comparison between different observations or for the wider dissemination of workflows. Providing this type of reproducibility, however, is often complicated by the need to accommodate the myriad dependencies included in a larger body of software, each of which generally come in various versions. Moreover, in many fields (bioinformatics being a prime example), these versions are subject to continual change due to rapidly evolving technologies, further complicating problems related to reproducibility. Here, we propose a principled approach for building analysis pipelines and managing their dependencies. As a case study to demonstrate the utility of our approach, we present a set of highly reproducible pipelines for the analysis of RNA-seq, ChIP-seq, Bisulfite-seq, and single-cell RNA-seq. All pipelines process raw experimental data, and generate reports containing publication-ready plots and figures, with interactive report elements and standard observables. Users may install these highly reproducible packages and apply them to their own datasets without any special computational expertise beyond the use of the command line. We hope such a toolkit will provide immediate benefit to laboratory workers wishing to process their own data sets or bioinformaticians seeking to automate all, or parts of, their analyses. In the long term, we hope our approach to reproducibility will serve as a blueprint for reproducible workflows in other areas. Our pipelines, along with their corresponding documentation and sample reports, are available at http://bioinformatics.mdc-berlin.de/pigx</p>	
Corresponding Author:	Altuna Akalin	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:		
Corresponding Author's Secondary Institution:		
First Author:	Ricardo Wurmus	
First Author Secondary Information:		
Order of Authors:	Ricardo Wurmus	
	Bora Uyar	
	Brendan Osberg	
	Vedran Franke	
	Alexander Godschan	
	Katarzyna Wreczycka	
	Jonathan Ronen	
	Altuna Akalin	

Order of Authors Secondary Information:	
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
<p>Experimental design and statistics</p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	Yes
<p>Resources</p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	Yes
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist?</p>	Yes



[Click here to view linked References](#)

Reproducible genomics analysis pipelines with GNU Guix

Ricardo Wurmus^{1*}, Bora Uyar^{1*}, Brendan Osberg^{1*}, Vedran Franke^{1*}, Alexander Godtschan^{1*},
Katarzyna Wreczycka¹, Jonathan Ronen¹, Altuna Akalin^{1#}

¹The Bioinformatics Platform, The Berlin Institute for Medical Systems Biology, Max-Delbrück Center for
Molecular Medicine, Robert-Rössle-Strasse 10, 13125 Berlin, Germany

* Equal contributions

Corresponding author (e-mail: altuna.akalin@mdc-berlin.de)

Keywords: Pipelines in genomics, reproducible software, functional package management,
RNA-seq, single cell RNA-seq, ChIP-seq, Bisulfite-seq, differential expression, differential
binding, differential methylation.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Abstract

In bioinformatics, as well as other computationally-intensive research fields, there is a need for workflows that can reliably produce consistent output, independent of the software environment or configuration settings of the machine on which they are executed. Indeed, this is essential for controlled comparison between different observations or for the wider dissemination of workflows. Providing this type of reproducibility, however, is often complicated by the need to accommodate the myriad dependencies included in a larger body of software, each of which generally come in various versions. Moreover, in many fields (bioinformatics being a prime example), these versions are subject to continual change due to rapidly evolving technologies, further complicating problems related to reproducibility. Here, we propose a principled approach for building analysis pipelines and managing their dependencies. As a case study to demonstrate the utility of our approach, we present a set of highly reproducible pipelines for the analysis of RNA-seq, ChIP-seq, Bisulfite-seq, and single-cell RNA-seq. All pipelines process raw experimental data, and generate reports containing publication-ready plots and figures, with interactive report elements and standard observables. Users may install these highly reproducible packages and apply them to their own datasets without any special computational expertise beyond the use of the command line. We hope such a toolkit will provide immediate benefit to laboratory workers wishing to process their own data sets or bioinformaticians seeking to automate all, or parts of, their analyses. In the long term, we hope our approach to reproducibility will serve as a blueprint for reproducible workflows in other areas. Our pipelines, along with their corresponding documentation and sample reports, are available at <http://bioinformatics.mdc-berlin.de/pigx>

Introduction

Reproducibility of scientific workflows is a ubiquitous problem in science, and is particularly problematic in areas that depend heavily on computation and data analysis (see (Peng 2011)). For such work it is essential that installed software is identical to versions used in publication, in order to facilitate the reproduction of published data and the controlled manipulation or augmentation of these software systems. Unfortunately, this goal is often unattainable for a variety of related reasons: Research-oriented software may be hard to build and install due to unsatisfiable dependency constraints; non-trivial software may yield different results when built or used with different versions or variants of declared dependencies; on workstations and shared High Performance Computing (HPC) systems alike, it may be undesirable or even impossible to comply with version and variant requirements due to software deployment limitations. Moreover, It is unrealistic to expect users to manually recreate environments that match the system and binary substrate on which the software was developed. In the field of bioinformatics the above problem is exacerbated by the fact that data production technology moves extremely fast; existing software and data analysis workflows require frequent updates. Thus, it is paramount that multiple versions and variants of the same software can be automatically built, in order to ensure reproducibility of projects that are either in-progress, or are already published.

An important related issue is the reproducibility of workflows and pipelines across different machines. In addition to bioinformatics, many scientific fields require the researcher to prototype their code on local workstations with a custom software stack, and then later run it on shared HPC clusters for large data sets. The researcher must then be able to recreate their local environment on the cluster to ensure identical behavior. All of these concerns add to the burden on scientists, and valuable time that could be spent on research is wasted accommodating the limitations of system administration practices to ensure reproducibility. Even worse, reproducibility failures can be overlooked amid this complication, and publications could be accompanied with irreproducible analysis workflows or software. For these reasons, the scientific community in general -and fast evolving fields like bioinformatics in particular- need reliable and reproducible software package management systems.

In recent years, several tools have gained popularity among software developers and system administrators for wrapping Linux kernel features to accomplish process isolation, bind mounts, and user namespaces, or to deploy services in isolated environments (also called “containers”). Examples of such tools include: Docker, Singularity, and Ixc. These tools are sometimes also proposed as solutions to the reproducibility problem (Peng 2011; Boettiger 2015), because they provide a way to ship an application alongside all of its runtime dependencies. This approach necessitates the use of file system images that are modified using imperative statements, e.g. to run a package manager inside a namespace, with the goal of embedding all dependencies in an opaque binary image. Containers and binary disk images alone do not make traditional tooling

1
2
3
4 any more suitable for the purposes of reproducible science. Software deployment inside of the
5 container is still subject to the well-known limitations of traditional package managers, such as
6 intractable stateful behavior, time-dependent installation results, the inability to install and
7 control more than a handful application or library variants of packages on the same system, to
8 name a few. Container systems like Docker only shift the problem of reproducibility from the
9 package level to the level of binary disk images, which is a much less useful level of abstraction.
10 As such, they bring little more to the table than traditional virtual machine images, albeit with
11 different trade-offs. We claim that reproducibility takes a more rigorous, declarative approach to
12 software environment management and packaging itself. Other package and environment
13 managers (such as Conda, EasyBuild, Spack) fail to take the complete dependency graph into
14 account; instead, they make tacit assumptions about the deployment environment. As a result,
15 it is much harder to understand and exactly reproduce an environment as neither the full
16 complexity of the graph of transitive dependencies nor the configuration space is captured.
17
18
19
20
21
22

23 For all the above reasons, we propose functional package management -as implemented in
24 GNU Guix- as a way to mitigate or obviate these problems by allowing us to *declare* the
25 complete dependency graph of software packages (and all of their dependencies recursively).
26 One important feature of this approach is that it allows for bit-by-bit reproducibility. To illustrate
27 this, we created a set of analysis tools (or 'pipelines') for common genomics analysis data sets:
28 RNA-seq, CHIP-seq, BS-seq and scRNA-seq (for the sequencing of RNA, Chromatin
29 Immunoprecipitation, Bisulfite-treated DNA, and single-cell resolution RNA, respectively). Each
30 pipeline has a complex and large graph of dependencies, and each graph is comprehensively
31 declared as a GNU Guix package definition; the graph is then built reproducibly by relying on
32 Guix package manager features. Note that these pipelines also represent production-level
33 pipeline tools, rather than simply model examples -they come with a full set of features including
34 alignment, quality check, quantification, assay specific analysis and HTML reports.
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Results

Pipeline design and implementation philosophy

The pipelines provided here were designed with special focus on several key features: namely, that they be 1) easy to use, 2) easy to install, 3) easy to distribute, and, most importantly, 4) reproducible; all of which are inter-related constraints. Care was taken to ensure that all of the pipelines have a similar interface, so that familiarity with one pipeline would make for a gentler learning curve in learning to use the others. For the end-user, each pipeline has the same input types: a sample sheet and a settings file. The sample sheet contains information about samples such as names, labels, covariates etc. The settings file contains extra arguments related to the execution of the pipelines. The users can generally run pipelines as follows:

```
$ pigx [pipeline_name] [sample_sheet] -s [settings_file]
```

where [pipeline_name] can refer to any of the four pipelines: “rnaseq”, “chipseq”, “bsseq”, or “scrnaseq”. The resulting output provided to the users includes high quality reports and figures containing a standard set of results from basic analyses and data quality checks. Where appropriate, reports also contain certain interactive elements.

In implementing this toolset, one of our first design choices was to use a conventional build system, the GNU Autotools collection, to configure and build the pipelines as if they were first-class software packages in their own right rather than a mere collection of tools and “glue code”. Instead of assuming that a user will provide a suitable environment *at runtime*, the use of a build system allows us to capture the software environment *at configuration time*. This is achieved by explicitly checking for the presence of required tools in the build environment and recording their exact location in the pipeline's configuration file. At runtime, the pipeline refers only to tools through the configuration file and does not assume the availability of dependent software in the global environment. Moreover, using a well-established build system makes it easy to package the pipelines for any package manager. We chose GNU Autotools over other build systems for two reasons: it does not require users to have a copy of the build system software as it compiles to shell code (which is highly portable), and it has been established long enough to implement a conventional and flexible build interface with well-known behavior even in somewhat unusual circumstances, such as the installation of files into unique prefixes as is done when building with GNU Guix.

Capturing the build-time environment alone is not enough to ensure reproducibility, nor is the use of a build system sufficient to make installation easy. Thus, our second design choice was to package the pipelines for the GNU Guix package manager. Like other package managers, GNU Guix allows users to install, upgrade and remove software without having to know the details of dependencies or the build procedure. Unlike traditional package managers, however, GNU Guix takes a rigorous, declarative approach to software environment management and

1
2
3
4 packaging called functional package management. This approach takes into account the
5 complete graph of dependencies and build-time configurations, and maximizes build
6 reproducibility by building binaries in fully declared isolated environments. Packages are
7 installed into paths with unique prefixes that are computed from the complete dependency
8 graph, allowing for the simultaneous installation of different versions or variants of applications
9 and libraries. With functional package management, a given software build will generally yield
10 bit-identical files when the build is performed on different machines or on the same machine at
11 different points in time, independent of the current state of the system (caveats to this
12 generalization are discussed below).
13
14
15
16

17 We consider software reproducibility an important asset in controlled experimentation.
18 Reproducing a software environment bit for bit is not a goal in itself, but it provides us with a
19 foundation upon which we can perform precise changes to the environment and assess the
20 impact of these changes. Without bit-for-bit reproducibility we cannot be certain of the nature
21 and impact of differences in the software environment. While virtual machines or binary
22 application bundles such as Docker images would be sufficient to freeze the state of our
23 software environment, relying on these tools would forgo the ability to recreate that same
24 environment from scratch; nor would it be possible to reason about the environment at the level
25 of software packages. The approach of functional package management as implemented in
26 GNU Guix preserves the relationships between software packages and ensures that differences
27 to the environment can be accounted for.
28
29
30
31
32

33 A further design choice remained regarding the workflow management system, which would
34 execute a series of tasks mostly in the form of scripts from different programming languages.
35 For this purpose, we used SnakeMake (Köster and Rahmann 2012), which provides
36 target-driven execution infrastructure similar to GNU Make but with Python syntax, along with
37 useful features such as parallel execution on HPC scheduling systems. However, we would like
38 to emphasize that the choice of workflow management system is not the most critical step for
39 reproducibility, but rather the management of dependencies. The different pipeline stages are
40 implemented with a workflow management system stitching together various bioinformatics
41 tools; they are made configurable with the GNU Autotools and packaged with GNU Guix. This
42 means they will be build-reproducible and can be installed via the one-liner:
43
44
45
46

```
47 guix package --install pigx.  
48  
49  
50  
51
```

52 **RNA-seq pipeline**

53 **General Description of PiGx-RNA-seq Pipeline**

54
55
56
57 PiGx RNA-seq provides an end-to-end preprocessing and analysis pipeline for RNA-seq
58 experiments. The pipeline takes a set of raw fastq read files and the experimental design as
59 described by the user, and produces differential expression reports with figures and tables of
60
61
62
63
64
65

1
2
3
4 differentially expressed genes, as well as GO term analysis thereof. Furthermore, it provides
5 quality control reports about the experiment. To use the pipeline, the user must provide two
6 files: the sample sheet describing the samples and corresponding fastq files, and a settings file
7 with configuration parameters related to the pipeline's execution. The settings file lists, among
8 other things, the location of a reference genome for alignment, a GTF file with genome
9 annotations, and a transcriptome reference, as well as a list of desired differential expression
10 analyses to be performed, specifying which samples to use as cases and controls --see
11 package documentation here http://bioinformatics.mdc-berlin.de/pigx_docs/pigx-rna-seq.html for
12 more details.
13
14
15

16
17 The pipeline can then be run with the command

18
19 `$ pigx rnaseq [sample_sheet] -s [settings_file]`, to generate the output --
20 which comes in several sequential steps (see [figure 1](#)).
21

22
23 PiGx RNA-seq uses the reference genome and transcriptome provided by the user to produce
24 indices using *STAR* (Dobin et al. 2013) and *Salmon* (Patro et al. 2017) respectively. It then uses
25 *Trim Galore!* (Babraham 2018b) to trim low quality reads and remove adapter sequences before
26 aligning the reads to the reference using *STAR*. At this point, PiGx RNA-seq uses *fastqc*
27 (Babraham 2018a) and *MultiQC* (Ewels et al. 2016) to generate comprehensive quality control
28 reports of the sequencing, trimming, and alignment steps. PiGx RNA-seq also uses *BEDTools*
29 (Quinlan and Hall 2010) to compute the depth of coverage in the experiment and outputs
30 convenient bedgraph files. Gene level expression quantification is obtained from *STAR*, and
31 transcript level quantification using *Salmon*. The gene expression count matrix is then used to
32 run differential expression analyses as specified by the user, using *DESeq2* (Love, Huber, and
33 Anders 2014) for statistical analysis and *g:ProfileR* (Reimand et al. 2007) for GO-term analysis.
34 Each differential expression analysis produces a self-containing HTML report.
35
36
37
38
39

40 The differential expression reports produced are comprehensive, including sortable tables for
41 differentially expressed genes for a detailed view, principal component analysis plots for a
42 birds-eye view of the experiment, as well as MA and volcano plots. In addition, the reports
43 include a section with GO term enrichment analysis.
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

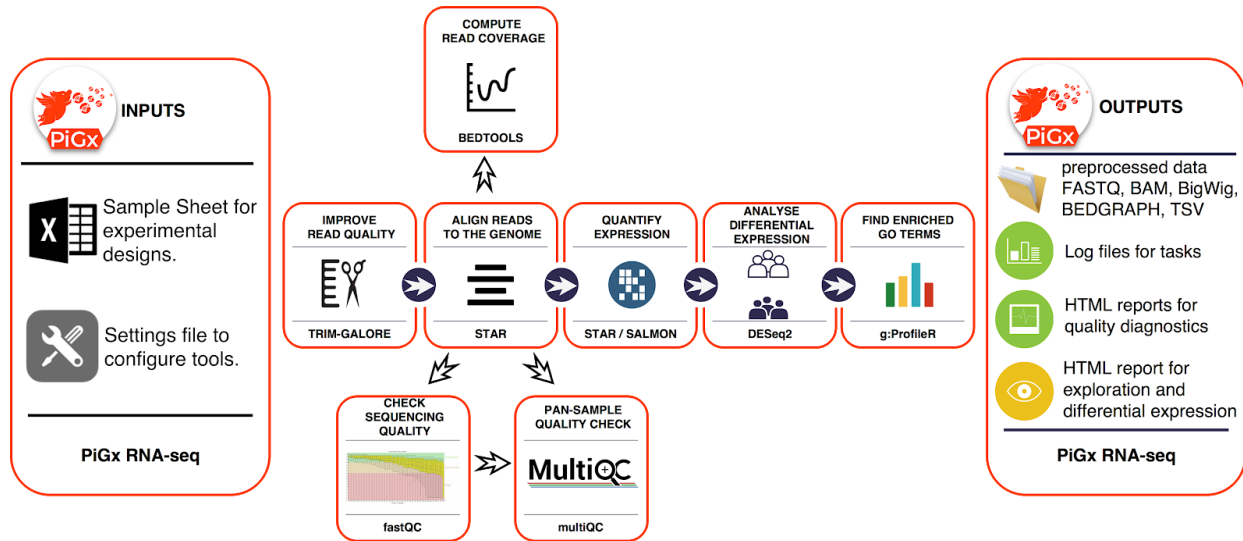


Figure 1
Workflow diagram of the PiGx-RNA-seq pipeline.

RNA-seq Use Case

The study by Hon *et al.* (2014) is motivated by several observations: DNA methyl-transferases (DNMTs) are the major mediators of cytosine methylation (producing 5-methyl-cytosine). 5hmC (5-hydroxy-methyl-cytosine) is a product of oxydation of 5mC's, and the TET family of proteins mediate 5mC oxydation. It has been established that DNA demethylation consists of the sequence of chemical reactions that convert 5mC into 5hmC, which is subsequently converted into 5fC (5-formyl-cytosine) and 5caC (5-carboxyl-cytosine). Active enhancers are depleted for 5mC but are enriched for 5hmC marks (Rampal *et al.* 2014), suggesting that an interplay between DNMTs and TET proteins could determine the activity level of enhancers. Mutating DNMTs or TET proteins in mouse embryonic stem cells (mESCs) perturbs global DNA methylation status, however cells do not lose the ability to regenerate. Moreover, mutating TET proteins and perturbing the oxydation levels have previously been shown to skew the differentiation of mESCs. Based on these facts, the authors address the following question: Can the skewed differentiation in mESCs be explained by deregulated balance of 5mC / 5hmC levels at active enhancers following the loss of activity of TET proteins?

The authors of the above study use TAB-Seq, Bisulfite-Seq, ChIP-seq and RNA-seq methods to profile genome-wide methylation, demethylation, histone modifications and gene expression levels to address these questions. They find that *Tet2* has the biggest role in enhancer demethylation in mESCs. Deletion of *Tet2* leads to enhancer hypermethylation, which in turn reduces enhancer activity. The reduced enhancer activity leads to a disruption in the activation of more than 300 genes in the early stages of differentiation, however the activity levels of these

1
2
3
4 genes are restored to wild-type levels at the later stages of differentiation. Reduced enhancer
5 activity followed by delayed gene activation explains the skew observed in mESC differentiation.
6
7

8 The authors of the above study profile the transcriptomes of mESCs as they differentiate into
9 neural progenitor cells (NPCs) within a six day period. They quantify gene expression levels of
10 wild-type, *Tet1* *-/-* and *Tet2* *-/-* cells on day zero, day three, and day six and sequenced two
11 biological replicates per sample. Thus, they obtained 18 samples in total (3 genotypes x 2
12 replicates x 3 days). In figure 5 of the original manuscript, the authors summarise the results of
13 the RNA-seq analysis. Here, we use the PiGx-RNA-seq pipeline to pre-process the raw fastq
14 files downloaded from the GEO archive (GEO accession: GSE48519), map the reads to the
15 *Mus musculus* genome (GRCM38 (mm10) build), and finally quantify the expression levels of
16 genes using both Salmon (Patro et al. 2017) and STAR (Dobin et al. 2013). We then use
17 DESeq2 (Love, Huber, and Anders 2014) to perform multiple differential expression analyses as
18 described in the original publication. Based on the processed and normalized count tables and
19 differential expression analysis results produced by the PiGx pipeline, we have written a small
20 custom script to reproduce the panels in figure 5 of Hon *et al.* In order to reproduce this figure,
21 we needed to perform seven differential expression analyses as described in Table 1. HTML
22 reports for each differential expression analysis (based on read counts computing using STAR)
23 can be found here: <http://bioinformatics.mdc-berlin.de/pigx/supplementary-materials.html>.
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41

Analysis	Case Sample	Control Sample	Description
tet2_diff_day3	day3_tet2_KO	day0_tet2_KO	<i>Tet2</i> <i>-/-</i> cells on day 3 are compared to <i>Tet2</i> <i>-/-</i> cells on day 0.
tet2_diff_day6	day6_tet2_KO	day0_tet2_KO	<i>Tet2</i> <i>-/-</i> cells on day 6 are compared to <i>Tet2</i> <i>-/-</i> cells on day 0.
WT_diff_day3	day3_WT	day0_WT	Wild-type cells on day 3 are compared to wild-type cells on day 0.
WT_diff_day6	day6_WT	day0_WT	Wild-type cells on day 6 are compared to wild-type cells on day 0.

42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

tet2_vs_WT_day0	day0_tet2_KO	day0_WT	<i>Tet2</i> <i>-/-</i> cells on day 0 are compared to wild-type cells on day 0.
tet2_vs_WT_day3	day3_tet2_KO	day3_WT	<i>Tet2</i> <i>-/-</i> cells on day 3 are compared to wild-type cells on day 3.
tet2_vs_WT_day6	day6_tet2_KO	day6_WT	<i>Tet2</i> <i>-/-</i> cells on day 6 are compared to wild-type cells on day 6.

Table 1

Differential expression analyses performed by PiGx-RNA-seq.

Having performed the above analysis, we first took a global look at how all sequenced samples cluster. Using a table of TPM (transcripts per million reads) counts generated by Salmon at the gene level, we selected the top 100 most variable genes and plotted a heatmap of all the samples using pheatmap package (Kolde 2018). We observed that the samples mainly cluster by the differentiation stage rather than genotype, which confirms the authors' findings (figure 2A). Next, again using the same TPM counts table, we plotted the expression levels of a select list of genes (*Nes6*, *Pax6*, *Sox1*, *Tet1*, *Tet2*, *Tet3*, *Slit3*, *Lmo4*, *Irx3*) on day 0, day 3, and day 6 (figure 2B). The changes in the expression levels of these genes perfectly match the patterns as described by Hon et al. At this point the authors recognise that some neural marker genes such as *slit3* and *lmo4* show discordant expression patterns between WT and *Tet2* *-/-* samples particularly on day 3, which are restored back to WT levels on day 6. The authors then investigated whether such a delayed induction mechanism can be observed globally. It was shown that the percentage of genes that are differentially expressed in both *Tet2* *-/-* and WT cells (compared to the undifferentiated samples of the corresponding genotypes on day 0), is significantly higher on day 6 than on day 3. We also observe a similar pattern, however the difference we observe is somewhat reduced. Our findings are reproduced based on gene counts quantified by both STAR and Salmon (figure 2C).

In figure 5F of the original publication, the authors take a closer look into the list of discordantly induced genes on day 3 in *Tet2* *-/-* samples. There it is shown that the majority of the genes that get induced in WT samples by day 3, don't get induced in the *Tet2* *-/-* samples as highly as they do in the WT samples. On the other hand, these numbers are comparable on day 6. We also observe the same difference and reproduce the findings using both Salmon and STAR-based gene counts (figure 2D). This suggests that there must be a list of genes that get activated in WT, but lag behind in *Tet2* *-/-* samples at the early stage of differentiation, however they catch up later with the WT levels. The authors call these genes '*delayed induction genes*' and find 333 genes that fit such a description. In figure 5G, the authors show the relative expression of these genes in *Tet2* *-/-* samples compared to WT samples throughout differentiation and compare it to the remaining list of genes in the genome. We have successfully reproduced the same patterns based on 357 delayed induction genes detected by Salmon-based gene counts (282 genes

detected by STAR-based gene counts) (Figure 2E). In figure 5H, the authors show the most significant GO terms enriched for the delayed induction genes. Although we don't observe the same set of terms as reported by the authors, we found seven development-related GO terms including 'tissue development' and 'nervous system development' as enriched terms (figure 2F).

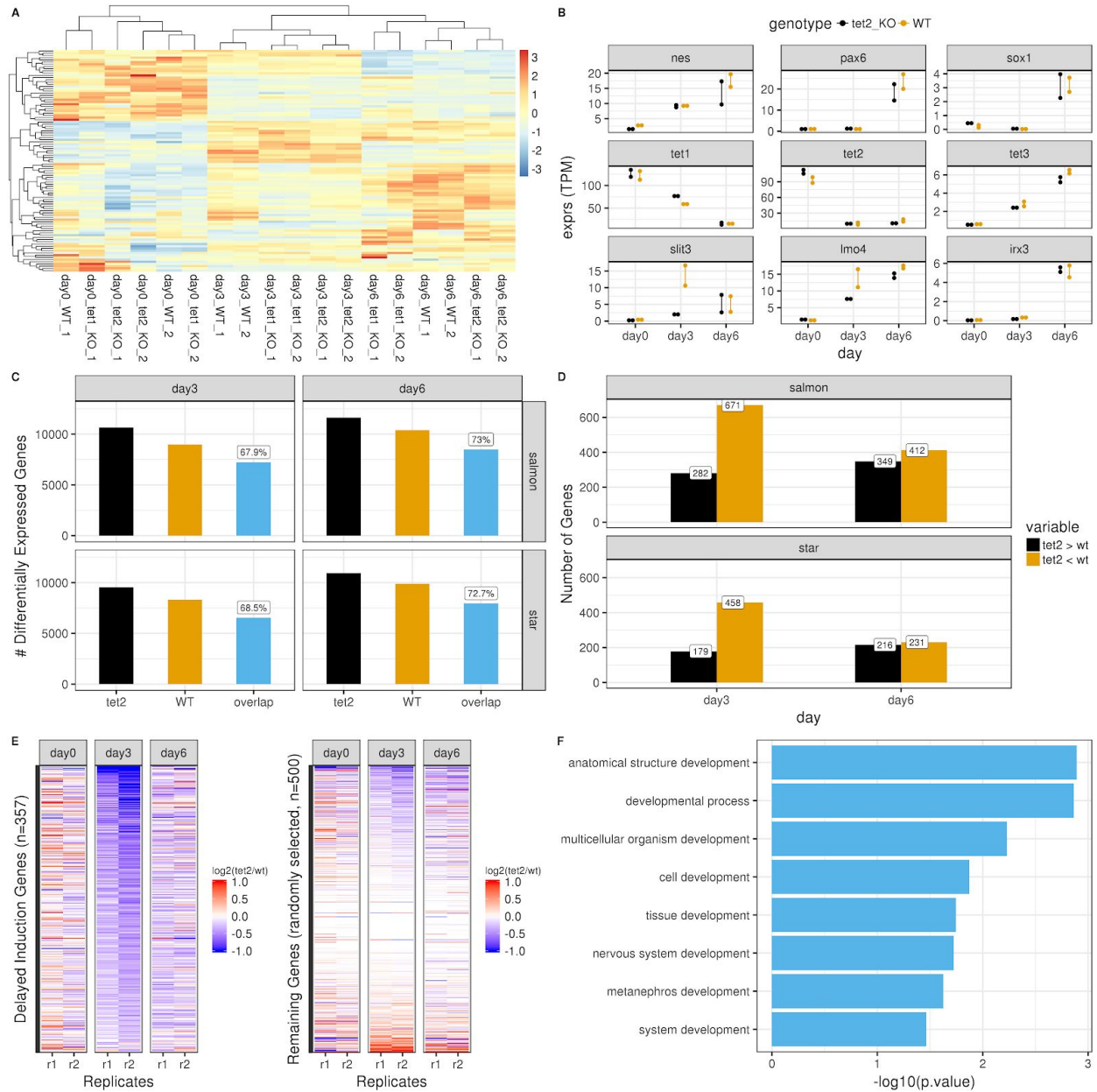


Figure 2

Reproduction of figure 5 from (Hon et al. 2014) using datasets processed by PiGx-RNA-seq pipeline. **A**) Hierarchically clustered heatmap of the top 100 most variable genes across all samples (transcripts per million (TPM) aggregated on the gene level, produced with Salmon).

Each row represents a gene and each column represents a sequenced sample (See Table 1 for

1
2
3
4 descriptions of the samples). The expression values are scaled by 'row'. **B)** Changes in the
5 expression levels of a selected list of genes throughout differentiation period on day 0, day 3,
6 and day 6. The y-axis shows the normalised expression levels (TPM at gene-level). The
7 expression patterns of samples with *Tet2* *-/-* background are depicted in black and wild type
8 background in orange. **C)** Abundance of differentially expressed genes (adjusted p-value < 0.1)
9 (on y-axis) when comparing samples on day 3 or day 6 with the samples on day 0 with
10 corresponding genotypes (*Tet2* *-/-* or wildtype). The bar labeled 'overlap' represents the number
11 of differentially expressed genes in both genotypes. The percentage is calculated by dividing the
12 value of 'overlap' with the value of *Tet2* *-/-*. The results are reproduced by both Salmon-based
13 gene-level read counts (top row) and STAR-based gene-level read counts (bottom row). **D)**
14 Genes that are up-regulated (induced) in wild-type samples on day 3 (or day 6) compared to
15 wild-type samples on day 0, are intersected with genes that are differentially expressed between
16 wild-type samples and *Tet2* *-/-* samples at the same stage of differentiation, and classified as
17 '*Tet2* > wt' (the gene is up-regulated in the *Tet2* *-/-* sample moreso than in the wild-type sample)
18 or '*Tet2* < wt' (the gene is upregulated in *Tet2* *-/-* sample less than in the wild-type sample). The
19 plot is reproduced using both Salmon-based gene counts and STAR-based gene counts. **E)**
20 Heatmaps for delayed induction genes (on the left) and 500 genes randomly selected from the
21 remainder (on the right). The colors of the heatmap represent the log₂ scale ratio of normalised
22 expression value (gene-level TPM counts obtained using Salmon) of each delayed induction
23 gene between *Tet2* *-/-* sample and the wild-type sample of the corresponding replicates (r1:
24 replicate-1, r2: replicate-2) on the corresponding stages of differentiation (day 0, day 3, and day
25 6). The rows of the heatmap are ordered in increasing order based on the average values of the
26 two replicates on day 3. The color scales range between -1 and 1 before reaching saturation. **F)**
27 Top GO terms for biological processes (on the y-axis) enriched among the delayed induction
28 genes. The GO terms are detected using g:ProfileR tool (Reimand 2016). The resulting terms
29 are filtered for p-value<0.05 and further filtered for the keyword 'development'. On the x-axis, the
30 p-values are depicted at log₁₀ scale.
31
32
33
34
35
36
37
38
39
40
41
42
43

44 ChIP-seq pipeline

47 General Description of PiGx-ChIP-seq Pipeline

48 PiGx ChIP-seq is an end-to-end processing and analysis pipeline for ChIP-seq experiments.
49 From the input fastq files, the pipeline produces sequencing quality control, ChIP quality control,
50 peak calling, and IDR (Q. Li *et al.* 2011) estimation. PiGx ChIP-seq also prepares the data for
51 visualization in a genome browser. The pipeline execution is highly customizable - the user can
52 specify which parts of the pipeline to execute, and which parameter settings to use. As in the
53 other pipelines, to use PiGx ChIP-seq, the user must provide two files: a sample sheet
54 containing the names of the fastq files with a descriptive label, and a settings file. The settings
55 file contains the locations of the reference genome, and the GTF file with genome annotations,
56
57
58
59
60
61
62
63
64
65

as well as a list of configurations for each executable step. Upon completion, the user is provided with quality reports, and all of the pre-processed data, which substantially facilitates downstream analysis and visualization.

PiGx ChIP-seq pipeline aligns the reads to the genome using Bowtie2 (Langmead and Salzberg 2012), does peak calling using MACS2 (Zhang et al. 2008), calculates the irreproducibility rate and outputs a series of quality statistics, such as: GC content, strand cross correlation, distribution of reads and peaks over annotated genomic features, and clustering of samples based on their similarity (Landt et al. 2012). The pipeline also produces UCSC Track hub for exploration of the dataset. The purpose of the pipeline is to improve the routine processing steps for ChIP-seq experiments and enable the user to focus on data quality control and biologically relevant data exploration. The pipeline heavily depends on Bioconductor (Huber et al. 2015) packages such as GenomicRanges (Lawrence et al. 2013) and Genomation (Akalin et al. 2015) for annotating peaks and summarizing ChIP-seq scores over regions of interest.

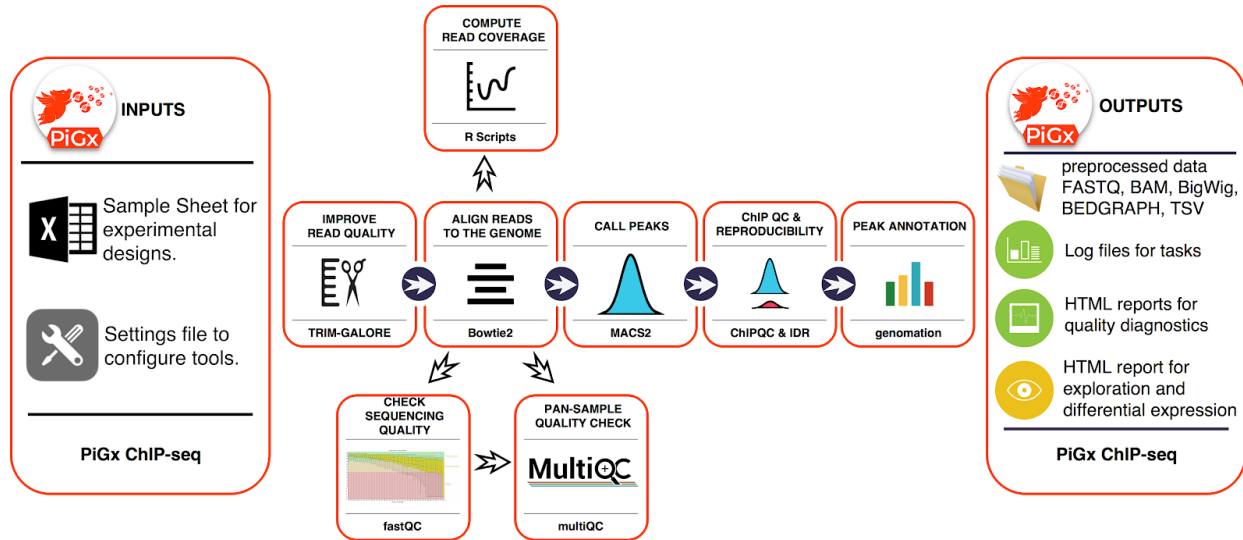


Figure 3
Workflow diagram for ChIP-seq pipeline

ChIP-seq Use Case

For consistency, we applied the ChIP-seq pipeline to data from the same study as in the section "RNA-seq Use Case" above (Hon et al. 2014); for the biological underpinnings of this experiment, please see the description provided there. Figure 4 shows part of the ChIP-seq quality control output performed on untreated, wild type ChIP samples, of various activating and repressing histone marks, and the corresponding input samples. One standard procedure is to validate the consistency of results with known biological priors, in order to quickly find samples

1
2
3
4 with outlying properties, and to discover batch effects. For example, figure 4A shows the
5 expected clustering of repressive (H3k27me3, H3k9me3) and activating (H3k4me3, H3k4me1,
6 H3k27ac, and H4k36ac) histone marks. Upon closer inspection, however, it becomes clear that
7 the activating histone marks cluster by their corresponding *batches*, and not by their biological
8 functionality. Figure 4B shows the cross-correlation between the signal on the plus and minus
9 genomic strands, shifted within a defined range (usually 1 - 400 nucleotides). The maximum
10 intensity in each row indicates the average DNA fragment size in each corresponding ChIP
11 experiment. Large discrepancies in the cross correlation profile, between experiments, can
12 indicate problems with fragmentation, fixation, or chromatin immunoprecipitation. The figure
13 shows that most of the samples have an average fragment size between 100 - 150 bp. One of
14 the H3k27me3 replicates, however, shows aberrant fragment size profile (second sample in the
15 plot). Upon visual inspection, the sample had extremely low signal to noise ratio and the peak
16 calling resulted in zero enriched regions. Such samples should either be repeated or omitted
17 from the downstream analysis. Figure 4C represents the relationship between the GC content of
18 one kilobase genomic bins and the ChIP signal; this plot is used as a diagnostics tool for
19 enrichment of fragments with extreme nucleotide content (enrichment of fragments with GC
20 content strongly deviating from the genomic mean), which can indicate problems with
21 PCR-based fragment amplification, and chromatin immunoprecipitation. Figure 4D represents
22 the distribution of reads over functional genomic features. It is used to observe whether the
23 experimental results conform to known expectations, based on previous experiments - i.e.
24 H3k4me3 should show strong enrichment over transcription start sites, while the H3k36me3
25 should show an enrichment over exonic and intronic regions. Non-conforming experiments can
26 indicate a weak ChIP, or antibody cross reactivity with unexpected epitopes. Figure 4
27 represents just a subset of quality control metrics implemented as a standard output from the
28 PiGx- ChIP-seq pipeline. The full set can be found here:

29
30
31
32
33
34
35
36
37 <http://bioinformatics.mdc-berlin.de/pigx/supplementary-materials.html>
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

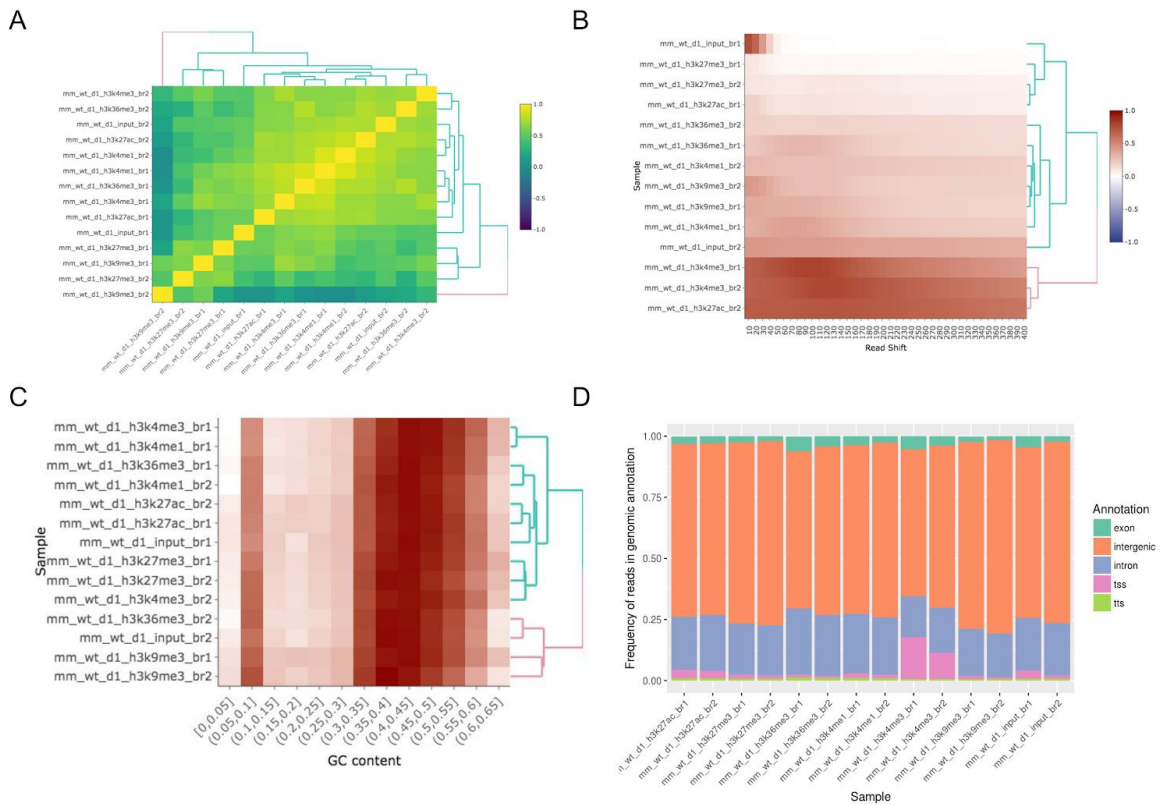


Figure 4

Example ChIP-seq quality control output. **A)** Clustering of samples based on correlation of normalized ChIP reads in one kilobase bins. **B)** Cross correlation between coverage profiles on Watson and Crick strands, shifted by the amount specified on the x axis. **C)** Relationship between read count and GC content in 1 kb bins. **D)** Distribution of reads in functional genomic features.

BS-seq pipeline

General description of the PiGx BS-seq pipeline

PiGx BS-seq is a bisulfite sequencing processing pipeline used to detect genome-wide methylation patterns and to perform differential methylation calling for case-control settings. It produces individual reports for each sample provided by the user, in addition to differential-methylation reports for arbitrarily many pairs of treatment conditions provided by the user. PiGx BS-seq uses *Trim Galore!* (Babraham 2018b) to trim reads for adapter sequences and quality, and *fastqc* (Babraham 2018a) for quality control (both before and after trimming). PiGx BS-seq produces GA- and CT- converted versions of the reference genome, if necessary,

using bismark_genome_preparation (Krueger and Andrews 2011). Reads are then mapped to the reference using Bowtie2 (Langmead and Salzberg 2012), before being sorted by location in the genome and filtered for uniqueness using samtools (Krueger and Andrews 2011; H. Li *et al.* 2009). The corresponding reports and .bam files for each of these steps are saved to their respective directories.

As in the other pipelines, to use PiGx BS-seq, the user must provide two input files: a sample sheet containing the paths to the fastq files with a descriptive label, and a settings file. The pipeline is robust to paired-end or single-end input data, and processing of each case is initiated automatically, based on whether the user supplies only a single input file, or a pair of files, for each sample. The settings file contains the locations of the reference genome, among other directories, as well as a list of configuration steps for each executable step. The pipeline can then be run with the command:

```
$ pigx bsseq [sample_sheet] -s [settings_file],
```

Post-mapping analysis steps performed automatically by PiGx BS-seq include tabulation of the fractional methylation of CpG sites, the segmentation of genomic methylation patterns across the genome, and the selection of differentially methylated sites between pairs of treatments provided in the settings file above. Furthermore, the final reports include genomic annotation of differentially methylated regions and methylome segments. A single execution of the pipeline can perform differential methylation analysis between a sample and arbitrarily many references; each comparison will have its own dedicated report, in addition to the final report for the sample itself. For traceability, direct links to input files, and various execution tools are saved directly within the output folder. Finally, a copy of the full methylome for each sample is also saved in BigWig (.bw) format, compatible with visualization in an online genome browser.

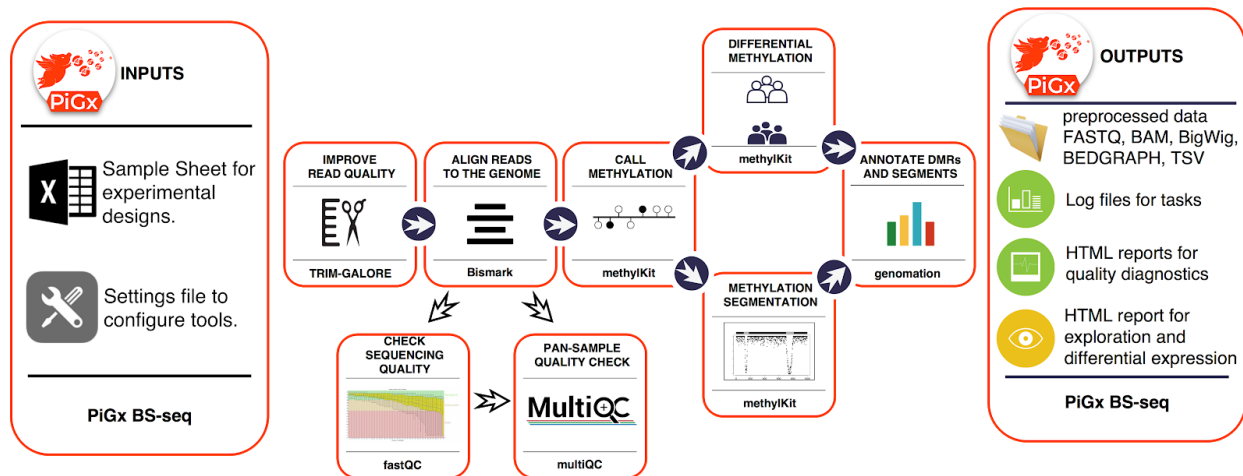


Figure 5
Workflow diagram for PiGX BS-seq pipeline

BS-seq Use Case

We applied the BS-seq pipeline to data from embryonic stem cells in mice, comparing wild type and *Tet2* deletion experiments (accessions SRX317877, and SRX317883 respectively). These data sets derive from the same study as was used for controlled comparison in the section “RNA-seq Use Case” above (Hon *et al.* 2014); for a biological description of this experiment, please refer to that section. HTML reports for each of the performed analyses can be found here: <http://bioinformatics.mdc-berlin.de/pigx/supplementary-materials.html>

Figure 6 shows a standard set of data analysis metrics generated automatically by the pipeline. For example, methylation levels near the promoter region of a list of annotated genes for each sample are shown in figures (A) and (B). For generality, figure 6 averages over all known genes; the user may freely probe for more specific results by supplying any arbitrary set of genes under investigation (in the absence of such an annotation file, this figure is simply omitted from the final report). A coarse map of the genome is provided in (C), which, for some datasets, may serve to highlight differential methylation localized to particular regions or chromosomes. In this particular use-case it is more useful as a null control showing that these regions are uniformly distributed throughout the genome. In addition, a histogram for differential methylation status of CpGs throughout the genome is provided in (D) using the same colour-code as in (C). The methylation differences of hyper-methylated, hypo-methylated and non-differentially methylated CpGs are shown as histogram with the color-code as in Figure 6C. The latter is shown as a distribution of methylation differences deemed to be not statistically significant (in black); since these are generally far more numerous than the former, the two curves are normalized independently. Note also that since these curves represent *relative* distributions, the vertical axis is of arbitrary units and tick marks are omitted. Finally, a screenshot of data-visualization from the genome browser (Robinson *et al.* 2011; Thorvaldsdóttir, Robinson, and Mesirov 2013) is provided in (E), here, regions of interest can be inspected manually at arbitrary precision.

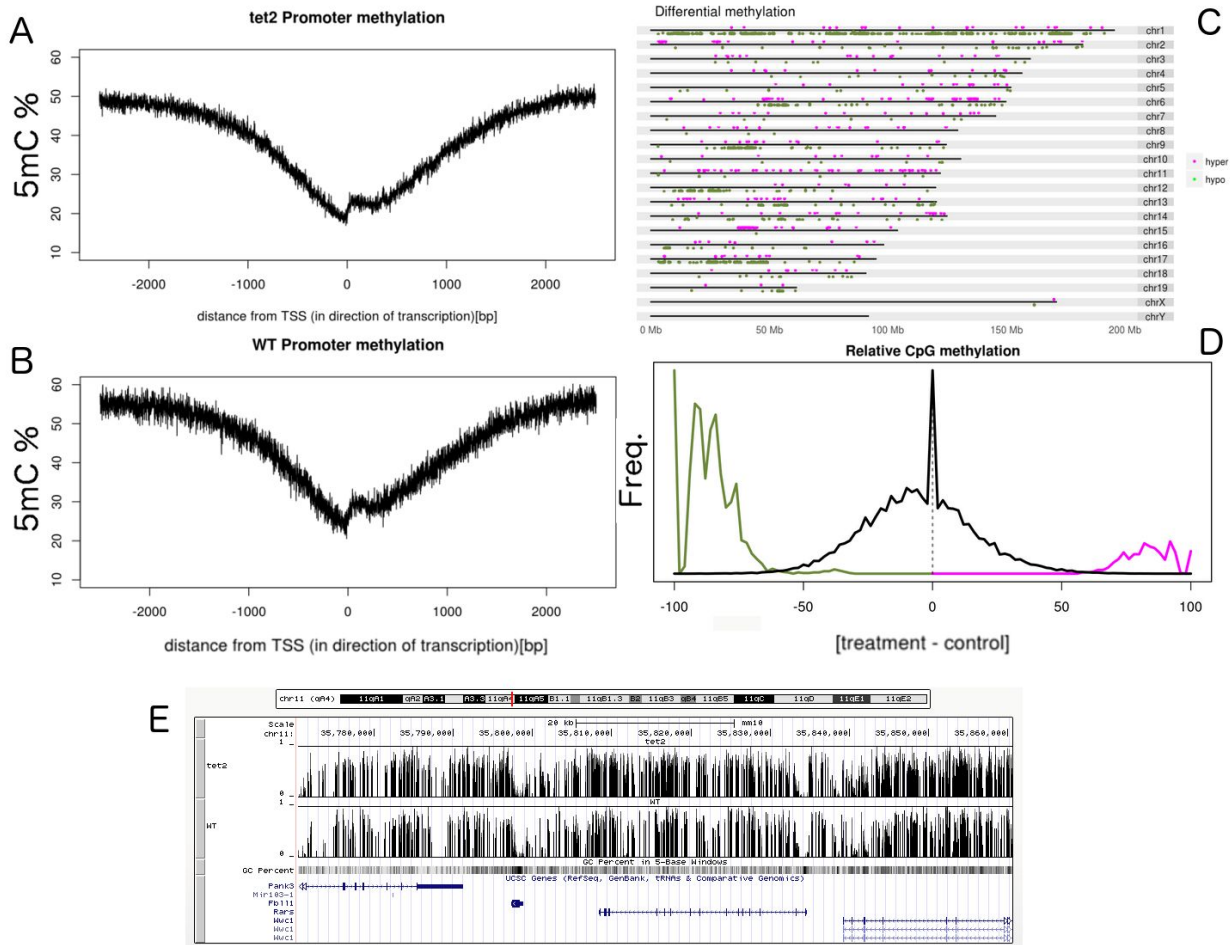


Figure 6

Output from the PiGx BS-seq pipeline. **(A,B)**: average CpG methylation throughout the promoter regions of the mm10 genome for *Tet2* $-/-$ and WT respectively. **(C)** Whole-genome map of differentially methylated CpGs, with colour-code to indicate hyper- and hypo- methylation of the treatment (*Tet2* $-/-$) relative to the control (Wild-type). **(D)** Histogram of the difference in average CpG methylation between *Tet2* $-/-$ and wild-type. For differentially-methylated cytosines, colors are consistent with (C), while CpGs with statistically insignificant difference in methylation are provided in black. Normalization of these two curves is performed independently (since the latter are generally far more numerous than the former), and the graph conveys only relative proportions (thus, as the absolute y-axis is of arbitrary scale, units and tick marks are omitted). **(E)** Screenshot of the genome browser using bigwig data from PiGx; here the data can be examined in much finer detail than in C).

scRNA-seq pipeline

General description of the PiGx scRNA-seq pipeline

Single cell RNA-seq is an extremely powerful technology, that is becoming increasingly prevalent in biological studies. The rapid development of UMI based methods, along with droplet based cell separation (Macosko et al. 2015; Klein et al. 2015), has enabled even simple experiments to quantify expression in several tens of thousand of cells. **PiGx scRNA-seq** is a pipeline for pre-processing of UMI based single-cell experiments. The purpose of the pipeline is to enable seamless integration and quality control of multiple single cell data sets. The pipeline works with minimal user input. As in the other pipelines, the user must provide a sample sheet with a basic experimental description, and a settings file which defines, among other parameters, the location of the input data and reference sequence and annotation. The pipeline can then be run with the command:

```
$ pigx scrnaseq [sample_sheet] -s [settings_file]
```

The pipeline does preliminary read processing, maps the reads with the STAR (Dobin et al. 2013) aligner, and assigns reads to gene models. It also separates cells from background barcodes (Alles et al. 2017), and constructs digital expression matrices for each sample (each saved in loom format); loom files from all samples are then merged into one large loom file using the loompy package (Linnarsson 2018). The expression data are subsequently processed into a SingleCellExperiment (Aaron Lun and Risso 2018) object. SingleCellExperiment is a Bioconductor class for storing expression values, along with the cell, and gene data, and experimental meta data in a single container. It is constructed on top of hdf5 file based arrays (Pagès 2018), which enables exploration even on systems with limited RAM (random access memory).

During the object construction, the pipeline performs expression normalization, dimensionality reduction, and identification of significantly variable genes. Then, it classifies cells by cell cycle phase and calculates the quality statistics. The SingleCellExperiment object contains all of the necessary data needed for further exploration. The object connects the PiGx pipeline with the Bioconductor single cell computing environment, and enables integration with state of the art statistical, and machine learning methods (scrn (A. T. L. Lun, McCarthy, and Marioni 2016), zinbwave (Risso et al. 2018), netSmooth (Ronen and Akalin 2018), iSEE (Aaran Lun et al. 2018), etc.).

The pipeline produces an HTML report containing quality controls, labeled by input covariates, which can be used for detecting batch effects.

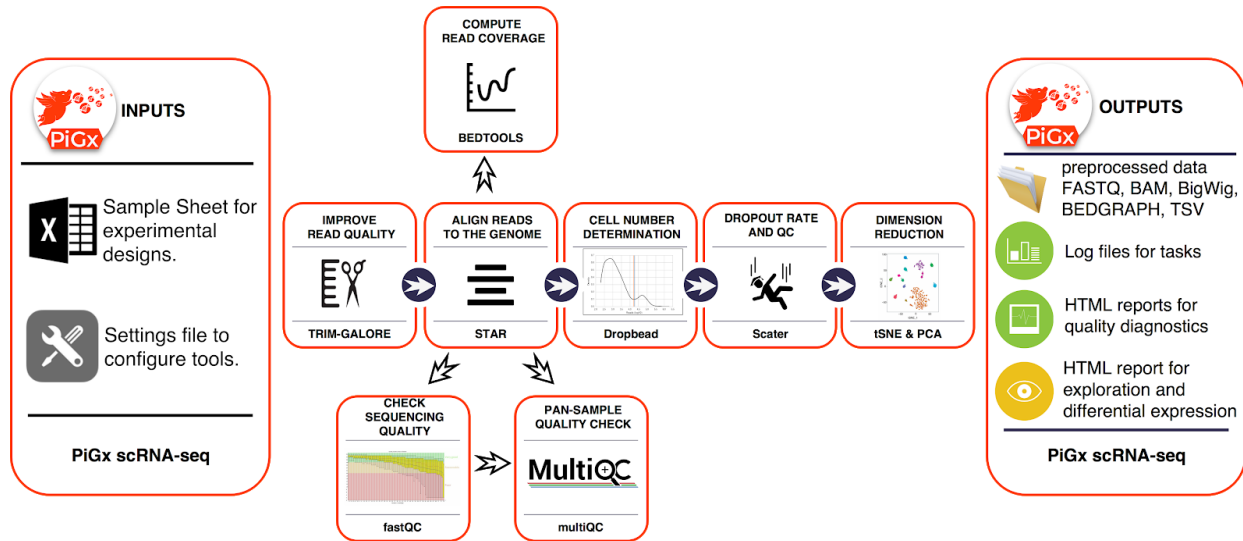


Figure 7
Workflow diagram for PiGx-scRNA-seq pipeline.

scRNA-seq Use Case

To showcase the capabilities of PiGx scRNA-seq, we ran the pipeline on isolated single nuclei from the mouse brain (Hu *et al.* 2017). In this study, the authors developed a gradient-based method for nucleus separation, and used it in combination with Drop-seq to profile the transcriptomes of more than 18,000 single nuclei. Figure 8 shows a part of the quality control output from the PiGx scRNA-seq pipeline. Figure 8A shows the per sample number of total and uniquely mapped reads. Figure 8B visualizes the cells on the first two principal components. The color gradient corresponds to the number of detected genes per cells. The figure shows that the total number of detected genes strongly correlates with the first two principal components. Figure 8C is analogous to figure 7B of the original publication, with the color scheme representing labeling each cell with its respective stage of the cell-cycle. Thus, figure 8C shows that the first two principal components correlate with the stage of the cell cycle. The heatmap in figure 8D shows scaled normalized expression values for genes that contribute the most to the first principle component. High read-count variability in a small number of genes drives the variation around the first principle component. The column-wise annotations show that the variation is driven mainly by cells in the G1 phase of the cell-cycle from the second biological replicate. The HTML report for this analysis can be accessed here:

<http://bioinformatics.mdc-berlin.de/pigx/supplementary-materials.html>

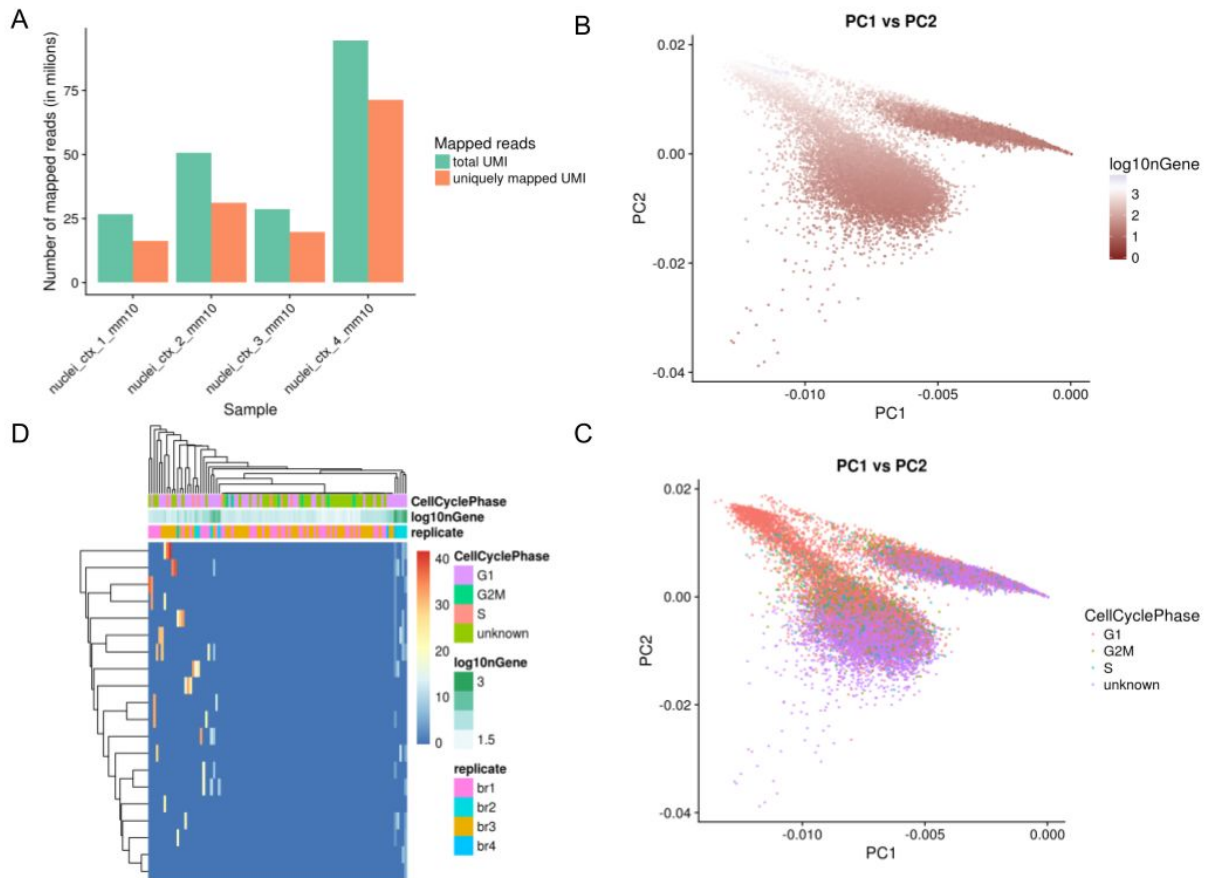


Figure 8

Sample output from the PiGx scRNA-seq pipeline. **A)** Abundance of total uniquely mapping UMIs per sample. **B)** Visualization of cells on the first and second principal component calculated from the normalized expression values. The gradient shows the total number of UMIs per cell. **C)** Same data representation as in B, but colored based on the cell cycle assignment. Cell cycle was assigned using the cyclone function from the scran Bioconductor package (A. T. Lun, McCarthy, and Marioni 2016). **D)** Expression heatmap of genes contributing most to the first principle component. Genes are ordered in rows, while cell are in columns. Color bars above the heatmap show relevant experimental variables.

Reproducibility metrics of the pipelines in different systems

We define the complete software environment needed for each of the pipelines using Guix package definitions. These package specifications not only outline the immediate dependencies of the pipelines, but extend to the full software stack recursively. The dependency graph is

1
2
3
4 rooted in a handful of bootstrap binaries. Apart from these binary roots, every application or
5 library in the graph is built from source. Guix ensures that packages are built in an isolated
6 environment in which nothing but the specified dependencies are available. This is a
7 precondition for bit-reproducible builds, i.e. repeatable package builds that yield the very same
8 binary output for the same set of inputs. Under ideal circumstances, a Guix specification for the
9 complete dependency graph and the set of all source code would be sufficient to exactly
10 reproduce the very same binaries of the pipelines presented in this paper.
11
12
13
14

15 Unfortunately, there are additional obstacles to bit-reproducibility that cannot be avoided purely
16 by the functional package management model. Examples for sources of irreproducibility in build
17 artefacts include embedded timestamps, non-deterministic sorting of strings, non-deterministic
18 compiler output, and the like. While some of these obstacles can be removed by deliberate
19 patching of compilers or applications, others are harder to diagnose and can thus lead to failure
20 to reproduce the same arrangement of bits in independent builds, be that on the same machine
21 at different points in time or on different systems.
22
23
24

25 To estimate the level of bit reproducibility in our pipelines, we checked out version
26 v0.14.0-3597-g17967d1 of GNU Guix, repeatedly built the pipeline packages and their direct
27 dependencies on three different systems (an office workstation, a virtual machine, and a build
28 farm consisting of 20 heterogeneous build nodes), and recorded the hashes of the produced
29 package trees. Whenever the hashes of any two builds differed we looked at the exact
30 differences with diffoscope (<https://diffoscope.org/>). Upon closer inspection we identified a
31 number of common issues in non-deterministic builds, such as timestamps embedded in
32 compiled binaries and text files, or randomized file names in files generated by test suites.
33
34
35
36

37 Python dependencies are of particular note here, because they are generally not reproducible
38 due to the fact that the byte compiler records the timestamp of the source file in the compiled
39 binary. This means that all compiled Python files will differ when they are compiled at different
40 points in time. (This problem will be addressed in the upcoming Python 3.7, which will
41 implement PEP 552 for deterministic compilation.) To avoid this problem and increase the
42 number of packages that could be made reproducible, we patched our variant of Python 3.6
43 such that it resets the embedded timestamp in compiled files to the Unix epoch. This allowed us
44 to greatly increase the number of fully bit-reproducible packages. As can be seen in Table 2,
45 only a total of 8 out of 355 packages (or only about 2.2%) were not bit-reproducible for as yet
46 unknown reasons.
47
48
49
50

51
52 Figure 9 visualizes the degree of bit-reproducibility for the direct dependencies of each of the
53 individual pipeline packages. Dependent packages whose files differed compared to builds on
54 other systems fell either in the category of “minor problems” or “not reproducible”, dependent on
55 the source and magnitude of non-determinism. The exact dependency counts for each category
56 and pipeline package are listed in Table 2. A comprehensive list of all dependent packages that
57 were categorized as having “minor problems” is contained in Table 3. This table shows that the
58
59
60
61
62
63
64
65

reproducibility problems of these packages are of negligible magnitude and could be corrected with minor patches to the package definitions in Guix.

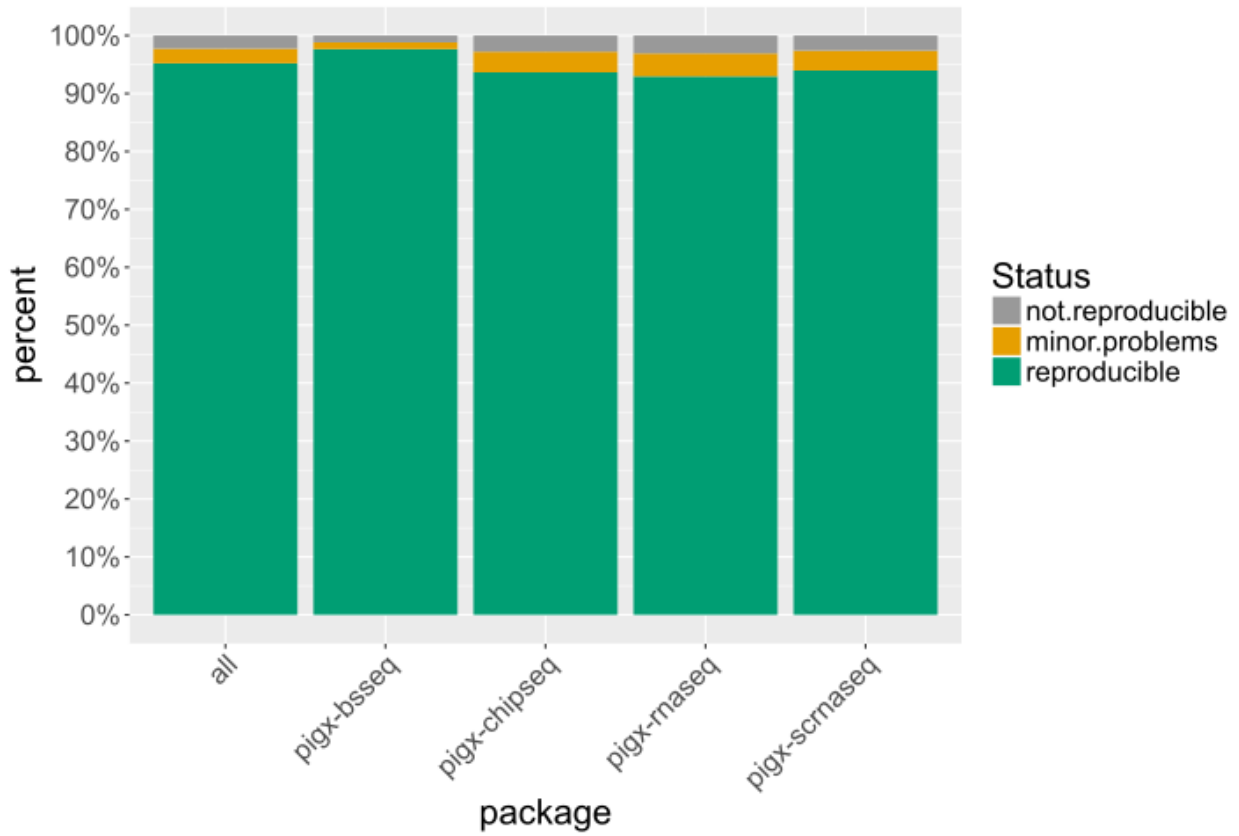


Figure 9
Percentage of directly-dependent packages building in a bit-reproducible fashion across different systems for each of the pipelines.

Package	Not reproducible	Minor problems	Reproducible
pigx-bsseq	2	2	167
pigx-chipseq	7	9	236
pigx-maseq	7	9	211
pigx-scrnaseq	6	8	218
All pipelines	8	9	338

Table 2

Number of dependent packages and their reproducibility status. See Table 3 for more details about packages with minor problems.

Package	Magnitude	Notes
r-minimal	2 bytes	non-deterministic line break
python	~ 6%	timestamp byte in header of bytecode files
python-matplotlib	~ 1.7%	single file difference
python-pycparser	~ 3%	single file with timestamp
python-cffi	~ 1.8%	recorded random test file names
python-numpy	< 0.5%	six bytecode files differ
python-simplejson	2 bytes	two files have single byte differences
gtk+	< 1%	single file (icon cache)
glib	< 0.1%	single file difference

Table 3

Table of packages with minor reproducibility problems and the magnitude of irreproducible files.

Alternative ways to install the pipelines:

We provide a generated application bundle containing all pipelines for use with Docker. The Docker image was generated by exporting the "closure" (i.e. the package and all packages it references, recursively) of the pigx package from the declarative Guix package definition instead of iteratively modifying a base image containing a GNU+Linux operating system in a series of imperative steps. The Docker image is merely a translation of a functional description of the desired environment; consequently, it is independent to global state, such as the contents of third-party package repositories or build time. The Docker image can be obtained at <https://hub.docker.com/r/bimsbbioinfo/pigx/>

Since the pipelines use the well-known GNU build system as implemented by the Autotools suite, the pipelines can be configured and built in any environment providing the required dependencies. The portable configure script detects and records references to needed software

1
2
3
4 in the environment and reuses them at runtime using their absolute file names. Any package
5 manager (such as Conda) can be used to fashion such a build-time environment. With regards
6 to reproducibility, however, we recommend that a package manager be used that can provide
7 separate, immutable, and uniquely prefixed environments to ensure that references to tools that
8 are recorded at configuration time are identical to the variants that are used at runtime.
9
10

11 Discussion

12
13
14
15 Computation is becoming an increasingly essential part of the biological sciences as the field
16 becomes more data intensive. The diversity and amount of data requires a multitude of tools
17 being used for analysis. Thus, the published software or workflows often come with complex
18 dependencies. Even if sensible guidelines (e.g. “Software with Impact” 2014), such as sharing
19 code online and providing documentation, are employed, sometimes it is impossible to recreate
20 the software used for analysis. Providing the code and documentation alone does not guarantee
21 reproducibility or usability, nor do Docker containers completely remedy this problem. We
22 propose GNU Guix and principled pipeline-as-software implementation as a solution to
23 reproducibility problems in complex bioinformatics workflows. Here, we demonstrated the utility
24 and the reproducibility of PiGx pipelines for genomics data analysis using GNU Guix.
25
26
27
28
29

30 Our decision to treat pipelines as first-class software packages and to adopt a conventional
31 build system with Autotools made it possible to reduce the installation of complex software
32 environments to a simple one-line command. By recording the exact locations of runtime
33 dependencies of the pipeline packages during the configuration stage, we were able to
34 eliminate ambiguity at runtime. When configuring the pipeline packages in an environment that
35 ensures that different versions or variants of applications and libraries are stored in unique
36 locations (such as an environment provided by GNU Guix), recording the exact location of
37 dependencies at *configuration time* allows us to reproduce the detected environment at *runtime*.
38
39
40

41 We have shown that with a recursive definition of software dependencies using the framework
42 provided by the functional package management paradigm as implemented in GNU Guix, it is
43 possible to fully and exhaustively describe complex real-world bioinformatics software
44 environments. The software environments were fully specified at the level of declarative,
45 stateless package abstractions instead of using an imperative, stateful approach. We have also
46 shown that the principled declarative approach to the management of software environments
47 lays a solid foundation for bit-reproducibility. The higher-level definitions of software
48 environments can be translated in an automated fashion to lower-level application bundles such
49 as Docker images. In contrast with container systems like Docker or Singularity, Guix encloses
50 the complete software environment and enables users to transparently rebuild it reproducibly
51 from source without having to trust a binary application bundle. Due to referential transparency,
52 binaries in Guix can only be the result of their corresponding sources.
53
54
55
56

57 Functional package management as implemented by GNU Guix significantly reduces the
58 complexity of, and lowers the barrier to, managing bit-reproducible software environments.
59 Users are freed from menial bookkeeping tasks such as keeping track of the origin of package
60
61
62
63
64
65

1
2
3
4 binaries, the time of installation, the order of installation instructions, the state of the operating
5 system at the time of installation, or any other runtime state. As far as users are concerned, it is
6 enough to know the names of the packages that should be installed (in our case, simply “pigx”)
7 and the current version of Guix; everything else such as source code provenance tracking,
8 dependency management, package configuration, and compilation in isolated environments is
9 handled by Guix. The guarantees provided by Guix enable users to contemplate obstacles to
10 experimental reproducibility beyond the software environment, such as sources of
11 non-determinism at *runtime*.
12
13

14
15 In our attempts to analyze the degree of repeatability of the HTML reports produced by PiGx,
16 we identified a number of such sources of non-determinism. The Salmon aligner, for example,
17 has a random component and does not provide a way for users to specify a seed for the
18 pseudo-random number generators. This makes it impossible to exactly repeat an analysis and
19 may require patching of the Salmon source code or virtualization of the random number
20 generator facilities of the host system. Other tools are sensitive to the user's locale settings and
21 may generate output in non-deterministic order. We were also surprised to find that an
22 increasingly large number of tools rely on a connection to the Internet, either directly or indirectly
23 through dependent packages. This can be a great source of non-determinism if the
24 experimental setup does not take the volatile nature of networked resources into account.
25 Another important obstacle to reproducibility is the large kernel binary at runtime. Although the
26 GNU C library provides a unified interface for all applications to use, the features that are
27 actually implemented by the kernel at runtime may differ vastly. For example, the variant of
28 Linux provided by Red Hat for their series 6 of operating systems reports its version as the
29 obsolete and unsupported 2.6.32, but it contains many backported features from much newer
30 kernel versions. Although this is usually not a problem, the kernel version and the implemented
31 features should be taken into account. Our use of version 2.26 of the GNU C library, for
32 example, necessitates either the use of Linux version 3.10 or higher, or a patched C library.
33
34
35
36
37
38

39 The use of a principled, declarative mechanism to managing software environments is a
40 fundamental component in a holistic approach to reproducibility at all levels: repeatable builds,
41 bit-reproducible binaries, software and data provenance, control over the configuration space,
42 and deterministic runtime behavior. We argue that this approach can serve as a template for
43 reproducible computational workflows.
44
45
46

47 Acknowledgements

48 We are grateful to the many volunteer contributors to GNU Guix who keep improving the
49 system.
50
51
52

53 Funding

54 B.U acknowledges funding by the German Federal Ministry of Education and Research (BMBF)
55 as part of the RNA Bioinformatics Center of the German Network for Bioinformatics
56
57
58
59
60
61
62
63
64
65

1
2
3
4 Infrastructure (de.NBI) [031 A538C RBC (de.NBI)]. We also acknowledge support for K.W from
5 Berlin Institute of Health (BIH).
6
7
8

9 10 **References**

- 11
12
13
14 Akalin, Altuna, Vedran Franke, Kristian Vlahoviček, Christopher E. Mason, and Dirk Schübeler.
15 2015. “Genomation: A Toolkit to Summarize, Annotate and Visualize Genomic Intervals.”
16 *Bioinformatics* 31 (7): 1127–29.
17
18 Alles, Jonathan, Nikos Karaïskos, Samantha D. Praktijnjo, Stefanie Grosswendt, Philipp Wahle,
19 Pierre-Louis Ruffault, Salah Ayoub, et al. 2017. “Cell Fixation and Preservation for
20 Droplet-Based Single-Cell Transcriptomics.” *BMC Biology* 15 (1): 44.
21
22 Babraham, Bioinformatics. 2018a. “fastQC.” 2018.
23 <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>.
24
25 ———. 2018b. “Trim Galore!” 2018.
26 https://www.bioinformatics.babraham.ac.uk/projects/trim_galore/.
27
28 Boettiger, Carl. 2015. “An Introduction to Docker for Reproducible Research.” *ACM SIGOPS*
29 *Operating Systems Review* 49 (1): 71–79.
30
31 Dobin, Alexander, Carrie A. Davis, Felix Schlesinger, Jorg Drenkow, Chris Zaleski, Sonali Jha,
32 Philippe Batut, Mark Chaisson, and Thomas R. Gingeras. 2013. “STAR: Ultrafast Universal
33 RNA-Seq Aligner.” *Bioinformatics* 29 (1): 15–21.
34
35 Ewels, Philip, Måns Magnusson, Sverker Lundin, and Max Käller. 2016. “MultiQC: Summarize
36 Analysis Results for Multiple Tools and Samples in a Single Report.” *Bioinformatics* 32
37 (19): 3047–48.
38
39 Hon, Gary C., Chun-Xiao Song, Tingting Du, Fulai Jin, Siddarth Selvaraj, Ah Young Lee,
40 Chia-An Yen, et al. 2014. “5mC Oxidation by Tet2 Modulates Enhancer Activity and Timing
41 of Transcriptome Reprogramming during Differentiation.” *Molecular Cell* 56 (2): 286–97.
42
43 Huber, Wolfgang, Vincent J. Carey, Robert Gentleman, Simon Anders, Marc Carlson, Benilton
44 S. Carvalho, Hector Corrada Bravo, et al. 2015. “Orchestrating High-Throughput Genomic
45 Analysis with Bioconductor.” *Nature Methods* 12 (2): 115–21.
46
47 Hu, Peng, Emily Fabyanic, Deborah Y. Kwon, Sheng Tang, Zhaolan Zhou, and Hao Wu. 2017.
48 “Dissecting Cell-Type Composition and Activity-Dependent Transcriptional State in
49 Mammalian Brains by Massively Parallel Single-Nucleus RNA-Seq.” *Molecular Cell* 68 (5):
50 1006–15.e7.
51
52 Klein, Allon M., Linas Mazutis, Ilke Akartuna, Naren Tallapragada, Adrian Veres, Victor Li,
53 Leonid Peshkin, David A. Weitz, and Marc W. Kirschner. 2015. “Droplet Barcoding for
54 Single-Cell Transcriptomics Applied to Embryonic Stem Cells.” *Cell* 161 (5): 1187–1201.
55
56 Kolde, Raivo. 2018. “Pheatmap: Pretty Heatmaps. R Package Version 1.0.8.” *CRAN*.
57 <https://CRAN.R-project.org/package=pheatmap>.
58
59 Köster, Johannes, and Sven Rahmann. 2012. “Snakemake—a Scalable Bioinformatics Workflow
60 Engine.” *Bioinformatics* 28 (19): 2520–22.
61
62 Krueger, Felix, and Simon R. Andrews. 2011. “Bismark: A Flexible Aligner and Methylation
63 Caller for Bisulfite-Seq Applications.” *Bioinformatics* 27 (11): 1571–72.
64
65 Landt, Stephen G., Georgi K. Marinov, Anshul Kundaje, Pouya Kheradpour, Florencia Pauli,
66 Serafim Batzoglou, Bradley E. Bernstein, et al. 2012. “ChIP-Seq Guidelines and Practices

- 1
2
3
4 of the ENCODE and modENCODE Consortia.” *Genome Research* 22 (9): 1813–31.
- 5 Langmead, Ben, and Steven L. Salzberg. 2012. “Fast Gapped-Read Alignment with Bowtie 2.”
6 *Nature Methods* 9 (4): 357–59.
- 7 Lawrence, Michael, Wolfgang Huber, Hervé Pagès, Patrick Aboyoun, Marc Carlson, Robert
8 Gentleman, Martin T. Morgan, and Vincent J. Carey. 2013. “Software for Computing and
9 Annotating Genomic Ranges.” *PLoS Computational Biology* 9 (8): e1003118.
- 10 Li, Heng, Bob Handsaker, Alec Wysoker, Tim Fennell, Jue Ruan, Nils Homer, Gabor Marth,
11 Goncalo Abecasis, Richard Durbin, and 1000 Genome Project Data Processing Subgroup.
12 2009. “The Sequence Alignment/Map Format and SAMtools.” *Bioinformatics* 25 (16):
13 2078–79.
- 14 Linnarsson. 2018. “Loompy: Python Implementation of the Loom File Format.” 2018.
15 <http://loompy.org>.
- 16 Li, Qunhua, James B. Brown, Haiyan Huang, and Peter J. Bickel. 2011. “Measuring
17 Reproducibility of High-Throughput Experiments.” *The Annals of Applied Statistics* 5 (3):
18 1752–79.
- 19 Love, Michael I., Wolfgang Huber, and Simon Anders. 2014. “Moderated Estimation of Fold
20 Change and Dispersion for RNA-Seq Data with DESeq2.” *Genome Biology* 15 (12): 550.
- 21 Lun, Aaran, Kevin Rue, Federico Marini, C. Sonesson, and Mark Robinson. 2018. “iSEE -
22 Interactive SummarizedExperiment/SingleCellExperiment Explorer.” 2018.
23 <https://github.com/csoneson/iSEE>.
- 24 Lun, Aaron, and Davide Risso. 2018. “Single Cell Experiment: S4 Classes for Single Cell Data.”
25 *Bioconductor*.
- 26 Lun, Aaron T. L., Davis J. McCarthy, and John C. Marioni. 2016. “A Step-by-Step Workflow for
27 Low-Level Analysis of Single-Cell RNA-Seq Data with Bioconductor.” *F1000Research* 5
28 (August): 2122.
- 29 Macosko, Evan Z., Anindita Basu, Rahul Satija, James Nemesh, Karthik Shekhar, Melissa
30 Goldman, Itay Tirosh, et al. 2015. “Highly Parallel Genome-Wide Expression Profiling of
31 Individual Cells Using Nanoliter Droplets.” *Cell* 161 (5): 1202–14.
- 32 Pagès, Hervé. 2018. “DelayedArray: Delayed Operations on Array-like Objects.” *Bioconductor*.
- 33 Patro, Rob, Geet Duggal, Michael I. Love, Rafael A. Irizarry, and Carl Kingsford. 2017. “Salmon
34 Provides Fast and Bias-Aware Quantification of Transcript Expression.” *Nature Methods* 14
35 (4): 417–19.
- 36 Peng, Roger D. 2011. “Reproducible Research in Computational Science.” *Science* 334 (6060):
37 1226–27.
- 38 Quinlan, Aaron R., and Ira M. Hall. 2010. “BEDTools: A Flexible Suite of Utilities for Comparing
39 Genomic Features.” *Bioinformatics* 26 (6): 841–42.
- 40 Rampal, Raajit, Altuna Alkalin, Jozef Madzo, Aparna Vasanthakumar, Elodie Pronier, Jay Patel,
41 Yushan Li, et al. 2014. “DNA Hydroxymethylation Profiling Reveals That WT1 Mutations
42 Result in Loss of TET2 Function in Acute Myeloid Leukemia.” *Cell Reports* 9 (5): 1841–55.
- 43 Reimand, Jüri, Meelis Kull, Hedi Peterson, Jaanus Hansen, and Jaak Vilo. 2007. “g:Profiler—a
44 Web-Based Toolset for Functional Profiling of Gene Lists from Large-Scale Experiments.”
45 *Nucleic Acids Research* 35 (Web Server issue): W193–200.
- 46 Risso, Davide, Fanny Perraudeau, Svetlana Gribkova, Sandrine Dudoit, and Jean-Philippe Vert.
47 2018. “A General and Flexible Method for Signal Extraction from Single-Cell RNA-Seq
48 Data.” *Nature Communications* 9 (1): 284.
- 49 Robinson, James T., Helga Thorvaldsdóttir, Wendy Winckler, Mitchell Guttman, Eric S. Lander,
50 Gad Getz, and Jill P. Mesirov. 2011. “Integrative Genomics Viewer.” *Nature Biotechnology*
51 29 (1): 24–26.
- 52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Ronen, Jonathan, and Altuna Akalin. 2018. "Network-Smoothing Based Imputation for Single Cell RNA-Seq." *F1000Research* 7 (January): 8.

"Software with Impact." 2014. *Nature Methods* 11 (February). Nature Publishing Group, a division of Macmillan Publishers Limited. All Rights Reserved.: 211.

Thorvaldsdóttir, Helga, James T. Robinson, and Jill P. Mesirov. 2013. "Integrative Genomics Viewer (IGV): High-Performance Genomics Data Visualization and Exploration." *Briefings in Bioinformatics* 14 (2): 178–92.

Zhang, Yong, Tao Liu, Clifford A. Meyer, Jérôme Eeckhoute, David S. Johnson, Bradley E. Bernstein, Chad Nusbaum, et al. 2008. "Model-Based Analysis of ChIP-Seq (MACS)." *Genome Biology* 9 (9): R137.