



**HAL**  
open science

# INSTITUT COCHIN ENTITY DMP - GENOM'IC Core Facility

Franck Letourneur, Thomas Guilbert

► **To cite this version:**

Franck Letourneur, Thomas Guilbert. INSTITUT COCHIN ENTITY DMP - GENOM'IC Core Facility. Institut Cochin U1016. 2024. hal-04789242

**HAL Id: hal-04789242**

**<https://hal.science/hal-04789242v1>**

Submitted on 18 Nov 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

---

## INSTITUT COCHIN ENTITY DMP - GENOM'IC

Plan de gestion de données créé à l'aide de DMP OPIDoR, basé sur le modèle "Science Europe: structured template for research entities" fourni par DMP OPIDoR.

### Plan Details

<b>Plan title</b>	INSTITUT COCHIN ENTITY DMP - GENOM'IC				
<b>Deliverable</b>					
<b>Version</b>	First version				
<b>Plan purpose/scope</b>	<p>The aim of this DMP is to describe as precisely as possible the various research products produced by the services offered by the core facility.</p> <p>This DMP is part of a wider framework, that of making available all the DMPs of the 9 Institut Cochin core facilities (MOUSET'IC, PROTEOM'IC, PIV, PIME, METABOL'IC, CYBIO, IMAG'IC, GENOM'IC and HISTIM) as well as the DMP of the CID database dedicated to the transfer, storage, consultation and sometimes processing of most of the data produced by the Institut Cochin's core facilities.</p> <p>The purpose of this DMP entity is to provide drafting assistance for researchers wishing to draw up a DMP for their research projects and wishing to call on the expertise of our core facilities.</p>				
<b>Fields of science and technology (from OECD classification)</b>	1.6 Biological sciences				
<b>Language</b>	eng				
<b>Creation date</b>	2024-07-18				
<b>Last modification date</b>	2024-10-25				
<b>Identifier</b>					
<b>Identifier type</b>	DOI				
<b>License</b>	<table><tr><td><b>Name</b></td><td>Creative Commons Attribution 4.0 International</td></tr><tr><td><b>URL</b></td><td><a href="http://spdx.org/licenses/CC-BY-4.0.json">http://spdx.org/licenses/CC-BY-4.0.json</a></td></tr></table>	<b>Name</b>	Creative Commons Attribution 4.0 International	<b>URL</b>	<a href="http://spdx.org/licenses/CC-BY-4.0.json">http://spdx.org/licenses/CC-BY-4.0.json</a>
<b>Name</b>	Creative Commons Attribution 4.0 International				
<b>URL</b>	<a href="http://spdx.org/licenses/CC-BY-4.0.json">http://spdx.org/licenses/CC-BY-4.0.json</a>				

### Structure Details

<b>Entity's name</b>	Institut Cochin Core Facility - GENOM'IC
<b>Acronym</b>	GENOM'IC
<b>Identifier</b>	<a href="https://institutcochin.fr/plateformes/genomic">https://institutcochin.fr/plateformes/genomic</a>
<b>Identifier type</b>	URL
<b>Description</b>	<p>GENOM'IC is one of the 9 core facilities of the Cochin Institute. It offers its expertise in the field of transcriptome analysis. It brings together cutting-edge skills and technologies in molecular biology and bioinformatics that it makes available to academic and industrial scientific communities.</p> <p>GENOM'IC offers ISO9001 certified services combining the production of data from high-throughput sequencing and their statistical analysis. It specializes its activities in the study of gene expression and regulation (RNA-seq, ChIP-seq, ATAC-seq) on tissue, cellular (single cell RNA-seq) and spatial (spatial transcriptomics) approaches.</p> <p>GENOM'IC supports teams from assistance with experimental design and choice of the adapted technique to each project until scientific publication.</p> <p>GENOM'IC ensures constant technological monitoring allowing it to develop its technical approaches in order to offer them to research teams.</p>
<b>Creation date</b>	1995-01-01

### Research outputs :

1. raw and analysed Next Generation Sequencing data (Dataset)

### Contributors

Name	Affiliation	Roles
Adoux Lucie - 0009-0004-8737-3877	Institut Cochin - 051sk4035	<ul style="list-style-type: none"> <li>• Data Collector</li> <li>• Data Collector</li> <li>• Data Collector</li> <li>• Project Member</li> </ul>
Guilbert Thomas - 0000-0001-5069-0730	Institut Cochin - 051sk4035	<ul style="list-style-type: none"> <li>• DMP manager</li> </ul>
Hamroune Juliette - 0000-0002-1812-2853	Institut Cochin - 051sk4035	<ul style="list-style-type: none"> <li>• Data Collector</li> <li>• Data Manager</li> <li>• Data Manager</li> <li>• Data Manager</li> <li>• Project Member</li> </ul>
Izac Brigitte - 0009-0004-5083-3681	Institut Cochin - 051sk4035	<ul style="list-style-type: none"> <li>• Data Collector</li> <li>• Data Collector</li> <li>• Data Collector</li> <li>• Project Member</li> </ul>
Jacques Sébastien - 0000-0003-0669-0586	Institut Cochin - 051sk4035	<ul style="list-style-type: none"> <li>• Data Collector</li> <li>• Data Collector</li> <li>• Data Collector</li> <li>• Project Member</li> </ul>
Letourneur Franck - 0000-0001-9465-4774	Institut Cochin - 051sk4035	<ul style="list-style-type: none"> <li>• Entity's Manager</li> <li>• Entity's Manager</li> <li>• Entity's Manager</li> <li>• Entity's Manager</li> <li>• Entity's Manager</li> <li>• Entity's Manager</li> <li>• Entity's Manager</li> <li>• Entity's Manager</li> <li>• Entity's Manager</li> <li>• Entity's Manager</li> </ul>
Martin Yoann	Institut Cochin - 051sk4035	<ul style="list-style-type: none"> <li>• Data Collector</li> <li>• Data Collector</li> <li>• Data Collector</li> <li>• Project Member</li> </ul>
Saintpierre Benjamin - 0009-0004-9440-5902	Institut Cochin - 051sk4035	<ul style="list-style-type: none"> <li>• Data Collector</li> <li>• Data Collector</li> <li>• Data Collector</li> <li>• Project Member</li> </ul>
Younsi Lilia - <a href="https://www.linkedin.com/in/hilia-y-484613173">linkedin.com/in/hilia-y-484613173</a>	Institut Cochin - 051sk4035	<ul style="list-style-type: none"> <li>• Data Collector</li> <li>• Data Collector</li> <li>• Data Collector</li> <li>• Project Member</li> </ul>

# INSTITUT COCHIN ENTITY DMP - GENOM'IC

---

## Data description and collection or re-use of existing data

### Research output description

<b>Name</b>	raw and analysed Next Generation Sequencing data
<b>Description</b>	<p>GENOM'IC provide its expertise to research teams in the transcriptomic analysis of their project. This is realized through production and analyses of Next generation sequencing data.</p> <p>For each new project, GENOM'IC define an experimental design to answer to the scientific question asked by the research team : number / quality of samples, type of library preparation and sequencing mode, bio-informatic methodology. GENOM'IC delivers then raw and analyzed data.</p> <p>We distinguish two steps in the new data generation :</p> <ul style="list-style-type: none"><li>• Raw data produced by the NGS Illumina sequencers (Nextseq 500, 2000 and Miseq) on biological samples : data relative to each sequencing run (composed of different samples) are first produced under a proprietary format (Bcl). Once obtained, software solutions convert it to generate data for each individual sample in a standardized open format (fastq).</li><li>• Analyzed data obtained from various processes of the raw data linked to the type of data produced and to the scientific question asked in the project. These processes generate intermediate and final analyzed data and need descriptive data of each sample obtained by the research team.</li></ul> <p>Whatever its nature, informations concerning each data (ex : raw data production date, type of technical treatment, versions of algorithms used for analysis, ...) are directly formalized during the production process (with a project file resuming all intermediary results and tracing all reagents used during the process) and the analysis process (with an analysis report resuming methodology, main results and tracing all versions used for the analysis).</p> <p>.GENOM'IC can also re-use existing data coming from</p> <ul style="list-style-type: none"><li>• previous partner project : minimal informations are provided by the research team.</li><li>• public database from litterature</li></ul>
<b>Type</b>	Dataset
<b>Workpackage</b>	
<b>Keywords</b>	<ul style="list-style-type: none"><li>• ADN (Thésaurus INRAE)</li><li>• ARN (Thésaurus INRAE)</li><li>• ARN messenger (Thésaurus INRAE)</li></ul>
<b>Keywords (free-text)</b>	
<b>Language</b>	eng
<b>Issued Date</b>	
<b>Persistent identifier</b>	The persistent identifier is the unique name of the project, indicating the year of production and a number.
<b>Identifier type</b>	Identifiant local
<b>May contain personal data?</b>	No
<b>May contain sensible data?</b>	No
<b>May take ethical issues into account?</b>	Yes

---

### Will existing data be reused?

<b>Justification</b>	<p>GENOM'IC can also re-use existing data coming from</p> <ul style="list-style-type: none"><li>• previous treated project by the platform with the agreement of the data owner research team</li><li>• public database from litterature in accordance with the license associated with the publication of the data concerned.</li></ul> <p>Data are the property of the supervisory authority on which the research team depends (INSERM, CNRS, University, ...).</p>
----------------------	--

---

### How new data will be collected or produced?

<b>Name of the method</b>	Production of sequence files and analysis results
<b>Description</b>	Raw sequencing data is produced by Illumina sequencers owned by GENOM'IC. Raw files are in .bcl binary format. Processed sequence files are in .fastq format, and produced by the Aozan tool, using bcl-convert software. Mapping files are in .bam format, and produced by STAR aligner. Count files are in .tsv format, and produced by RSEM tool. Files from the differential analysis are in various formats (.tsv, .png, .rds), all generated with the R package DESeq2.
<b>Data Nature</b>	Experimental Data
<b>Equipments, technical platforms used</b>	<ul style="list-style-type: none"> <li>GENOM'IC : <a href="https://cat.opidor.fr/index.php/GENOM'IC">https://cat.opidor.fr/index.php/GENOM'IC</a></li> </ul>
<b>Related references</b>	<ul style="list-style-type: none"> <li>Aozan : <a href="https://www.outils.genomique.biologie.ens.fr/aozan/index.html">https://www.outils.genomique.biologie.ens.fr/aozan/index.html</a></li> <li>STAR: ultrafast universal RNA-seq aligner : <a href="https://github.com/alexdobin/STAR">https://github.com/alexdobin/STAR</a></li> <li>RSEM: accurate transcript quantification from RNA-Seq data with or without a reference genome : <a href="https://github.com/deweylab/RSEM">https://github.com/deweylab/RSEM</a></li> <li>Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2 : <a href="https://bioconductor.org/packages/release/bioc/html/DESeq2.html">https://bioconductor.org/packages/release/bioc/html/DESeq2.html</a></li> </ul>

## Documentation and data quality

### What metadata and documentation (for example way of organising data) will accompany the data?

<b>Description</b>	<p>During all processes, original sample names are conserved and all associated elements are stored in a local SQL database, to retrieve information whenever needed.</p> <p>We identified two main types of metadata linked to the two main types of data (raw and analyzed). Each project is identified by a code "NGS-year of production - number". All metadata of a project are conserved in the directory file of the project.</p> <p>For raw data, we identified</p> <ul style="list-style-type: none"> <li>metadata coming from the scientific project (administrative informations, experimental design, objectives of the study, type and names of samples, organism studied, ...) : they are conserved in the project file filled by the team leader. Samples processed by GENOM'IC arrive with sample names given by the research lab they come from. GENOM'IC transcribes these names as filenames of the FASTQ files previously described.</li> <li>metadata coming from the production step (quality control of libraries and sequencing, fluorimetry measures, sequencing mode, ...) : they are conserved in the project directory in a follow-up file and in the systems ("cahier de laboratoire électronique") used to produce these data (instruments for QC control of libraries or RNA, sequencers).</li> </ul> <p>For analyzed data, we identified information about reference genome used for the analysis (reference files from the Ensembl database (<a href="https://ensembl.org">ensembl.org</a>)), and about the different tools used for the analysis, for which GENOM'IC provides bibliographical references and used versions of the tools, to match the FAIR recommendations. For each project, a report file resuming the different analysis steps and their results is given.</p>
<b>Related references</b>	<ul style="list-style-type: none"> <li>Metadata, FAIR principles, and their importance in genomics : <a href="https://genestack.com/assets/pdfs/The%20importance%20of%20metadata%20in%20genomics%20and%20the%20FAIR%20principles%20ebook.pdf">https://genestack.com/assets/pdfs/The%20importance%20of%20metadata%20in%20genomics%20and%20the%20FAIR%20principles%20ebook.pdf</a></li> <li>FAIR Genomes : <a href="https://fairgenomes.org/">https://fairgenomes.org/</a></li> <li>FAIR Genomes metadata schema promoting Next Generation Sequencing data reuse in Dutch healthcare and research : 10.1038/s41597-022-01265-x</li> </ul>
<b>Metadata/data standards</b>	<ul style="list-style-type: none"> <li>Genome Metadata : <a href="https://rdamsc.bath.ac.uk/api2/m19">https://rdamsc.bath.ac.uk/api2/m19</a></li> </ul>
<b>Metadata language code</b>	eng

### What methods will be used to ensure their scientific quality?

<b>Description</b>	<p>Each instrument used in the pipeline of the data production (quality control of libraries, sequencer) is annually controlled to ensure a perfect data quality. Each protocol used for data production follow processes described in our quality management system (certified Iso9001 V2015)</p> <p>For raw sequencing data, processed FASTQ files generated by Aozan are simultaneously checked for quality using the <a href="#">FastQC</a> and <a href="#">MultiQC</a> modules. Different metrics quality data are available (% of flowcell occupancy, % of Q30) and used as quality indicators in our ISO9001 quality management system.</p> <p>For analyzed data, the quality control is based on the expertise of the bioinformatician in charge of the analysis, using different tools. For example, the unsupervised approach (via PCA Principal Component Analysis) is used to confirm the belonging of each sample to its potential biological group. The alignment metrics are also one of the way to control quality data.</p>
<b>Related references</b>	<ul style="list-style-type: none"> <li>Aozan: an automated post-sequencing data-processing pipeline : <a href="https://www.outils.genomique.biologie.ens.fr/aozan/index.html">https://www.outils.genomique.biologie.ens.fr/aozan/index.html</a></li> <li>A quality control tool for high throughput sequence data. : <a href="https://www.bioinformatics.babraham.ac.uk/projects/fastqc/">https://www.bioinformatics.babraham.ac.uk/projects/fastqc/</a></li> </ul>

## Legal and ethical requirements, codes of conduct

### How will other legal issues, such as intellectual property rights and ownership, be managed? What legislation is applicable?

<b>Description</b>	The research organism financing the project is the owner of the data and results generated by the facility and of their intellectual property. GenomIC produce and analyze the data and store them (with metadata) for 6 months in its data server. These conditions are described in the general conditions of use of the facility.
--------------------	---

---

#### What ethical issues and codes of conduct are there, and how will they be taken into account?

<b>Description</b>	Main ethical issues are related to human studies. These data are sensible data because they include genetic informations usable for individual recognition. To prevent this, samples are pseudonymized. Data are stored in servers under the responsibility of the DSI INSERM. Higher storage protections are evaluated (double authentication access / isolated server from internet).
--------------------	---

---

## Data processing and analysis

#### How and with what resources will the data be processed / analyzed?

<b>Description</b>	.bcl file coming directly from the sequencers are converted in FASTQ files using BCL converter software integrated in Aozan software. .fastq files are then processed and analyzed as follows : <ul style="list-style-type: none"> <li>• alignment on the reference genome</li> <li>• counting genes</li> <li>• differential analysis</li> </ul> <p>The alignment step is done with the STAR aligner (see related references) using genome reference files from the Ensembl database. Output files are in BAM binary format. The counting of the genes is done with the RSEM tool. Output files are in .tsv format with gene IDs and associated raw and normalized values. Differential analysis is done with the DESeq2 R package. Various output files are produced : normalized data table in .tsv format, results data table with statistical output also in .tsv format, different graphs in .png format, and an R output file in .rds format, along with an HTML report. Pipelines are development tools from GenomIC with various coding languages like R or Python. They are regularly updated on the Git of the lab.</p>
--------------------	--

<b>Related references</b>	<ul style="list-style-type: none"> <li>• STAR: ultrafast universal RNA-seq aligner : <a href="https://github.com/alexdobin/STAR">https://github.com/alexdobin/STAR</a></li> <li>• RSEM: accurate transcript quantification from RNA-Seq data with or without a reference genome : <a href="https://github.com/deweylab/RSEM">https://github.com/deweylab/RSEM</a></li> <li>• Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2 : <a href="https://bioconductor.org/packages/release/bioc/html/DESeq2.html">https://bioconductor.org/packages/release/bioc/html/DESeq2.html</a></li> </ul>
---------------------------	---

<b>Equipments, technical platforms</b>	<ul style="list-style-type: none"> <li>• GENOM'IC : <a href="https://cat.opidor.fr/index.php/GENOM'IC">https://cat.opidor.fr/index.php/GENOM'IC</a></li> </ul>
--	--

---

## Storage and backup during the research process

#### How will data be stored and backed up during the research?

<b>Storage needs</b>	Raw data coming from sequencers are directly stored on storage servers at the end of their production in a dedicated directory. Actual data storage is based on two RAID servers (Nettapp serveur with 58To/67 - data from 2015 to 2020 - and a new Dell server with 51To/233 from 2021 to 2024). Raw data storage need is around 10-15 To per year. Storage requirements are indicated for an average year of data generation.
<b>Estimated volume of data</b>	12
<b>Unit</b>	TB
<b>Equipments, technical platforms</b>	<ul style="list-style-type: none"> <li>• CID : <a href="https://cat.opidor.fr/index.php/CID">https://cat.opidor.fr/index.php/CID</a></li> <li>• GENOM'IC : <a href="https://cat.opidor.fr/index.php/GENOM'IC">https://cat.opidor.fr/index.php/GENOM'IC</a> :</li> </ul>
<b>Measures taken for data security</b>	Experimental raw and analyzed data produced by GenomIC are stored on dedicated servers, located in the room servers of the institute Cochin. Data security is actually ensured by the security level of INSERM servers which host them. Access to the storage servers is reserved to GenomIC via a single authentication (login + pswd).

---

## Data sharing and long-term preservation

#### How will data be shared?

**Modalities of sharing**

At the end of the production phase, raw and analyzed data are shared with the team leader: raw data are loaded in the CID (<https://imagerie.cochin.inserm.fr/sis4web/login.php>) or sent by filesender (<https://filesender.renater.fr/>), or copied on a hard disk. At this point, the team leader is responsible of the long term storage of its data.

At the end of the publication phase, raw data need to be shared with the scientific community via their loading on dedicated databases (GEO, SRA).

**Reusability**

YES

**Data repository/catalogs**

- GEO : <https://www.ncbi.nlm.nih.gov/geo/> ()

---

**How will data be long-term preserved? Which data?****Justification**

Long term preservation of raw data is justified by their possible future use in a different scientific context of the one used to produce them. These data are fastq sequencing files extracted for each sample from the sequencing run. GenomIC hasn't the capacity to store all raw data during many years. So long term storage (unlimited store, depending of the storing tool) is under the responsibility of the team leader.

Beside storage on RAID servers, raw data are also backed up on external hard disks.

Once published, data are stored on dedicated database (GEO) and accessible to the scientific community.

**Estimated volume of data**

10

**Unit**

TB

**Start date****End date****Final dispositions**