
**Genomics informatics — Reliability
assessment criteria for high-
throughput gene-expression data**

*Informatique génomique — Critères d'évaluation de la fiabilité des
données d'expression des gènes à haut débit*



This document is a preview generated by EUS



COPYRIGHT PROTECTED DOCUMENT

© ISO 2021

All rights reserved. Unless otherwise specified, or required in the context of its implementation, no part of this publication may be reproduced or utilized otherwise in any form or by any means, electronic or mechanical, including photocopying, or posting on the internet or an intranet, without prior written permission. Permission can be requested from either ISO at the address below or ISO's member body in the country of the requester.

ISO copyright office
CP 401 • Ch. de Blandonnet 8
CH-1214 Vernier, Geneva
Phone: +41 22 749 01 11
Email: copyright@iso.org
Website: www.iso.org

Published in Switzerland

Contents

Page

Foreword	iv
Introduction	v
1 Scope	1
2 Normative references	1
3 Terms and definitions	1
4 From sample to RNA	3
4.1 General	3
4.2 RNA integrity	3
4.3 RNA concentration	4
4.4 RNA purity	4
5 Expression profiling	4
6 Quality control metrics in RNA-seq analysis	4
6.1 General	4
6.2 Sequencing read	4
6.2.1 Total number of reads	4
6.2.2 Read length	5
6.2.3 Base call quality	5
6.2.4 GC content	5
6.2.5 Overrepresented sequence	5
6.2.6 Adapter residue	5
6.3 Alignment	5
6.3.1 Alignment ratio	5
6.3.2 Gene body coverage uniformity	5
6.3.3 Strand specificity	6
6.3.4 Insert size	6
6.3.5 Mismatch	6
6.3.6 Contamination from other sources	6
6.4 Expression	6
6.4.1 Expression distribution	6
6.4.2 Expressed genes	6
6.4.3 Saturation	7
6.4.4 Reproducibility	7
6.5 Differentially expressed genes	7
6.6 Biological interpretation of differentially expressed genes	7
6.7 Sample certificate of origin	8
6.8 Quality control of batch effects	8
7 Spike-in controls	8
8 Proficiency testing	8
8.1 General	8
8.2 RNA sources	9
8.3 Experimental design	9
9 Process management	9
Bibliography	10

Foreword

ISO (the International Organization for Standardization) is a worldwide federation of national standards bodies (ISO member bodies). The work of preparing International Standards is normally carried out through ISO technical committees. Each member body interested in a subject for which a technical committee has been established has the right to be represented on that committee. International organizations, governmental and non-governmental, in liaison with ISO, also take part in the work. ISO collaborates closely with the International Electrotechnical Commission (IEC) on all matters of electrotechnical standardization.

The procedures used to develop this document and those intended for its further maintenance are described in the ISO/IEC Directives, Part 1. In particular, the different approval criteria needed for the different types of ISO documents should be noted. This document was drafted in accordance with the editorial rules of the ISO/IEC Directives, Part 2 (see www.iso.org/directives).

Attention is drawn to the possibility that some of the elements of this document may be the subject of patent rights. ISO shall not be held responsible for identifying any or all such patent rights. Details of any patent rights identified during the development of the document will be in the Introduction and/or on the ISO list of patent declarations received (see www.iso.org/patents).

Any trade name used in this document is information given for the convenience of users and does not constitute an endorsement.

For an explanation of the voluntary nature of standards, the meaning of ISO specific terms and expressions related to conformity assessment, as well as information about ISO's adherence to the World Trade Organization (WTO) principles in the Technical Barriers to Trade (TBT), see www.iso.org/iso/foreword.html.

This document was prepared by Technical Committee ISO/TC 215, *Health informatics*, Subcommittee SC 1, *Genomics informatics*.

Any feedback or questions on this document should be directed to the user's national standards body. A complete listing of these bodies can be found at www.iso.org/members.html.

Introduction

High-throughput gene-expression profiling, including data generated from microarray, next-generation sequencing, and other forms of high-throughput technologies, is a revolutionary technology for genomic studies. It is a fast-moving field both in terms of innovation in measurement technology as well as advances on the data analysis side. High-throughput expression technology enables us to efficiently study complex biological systems and biological processes, mechanisms of diseases, and strategies for disease prevention and treatment. This technology is currently applied in the biomedical research community and industry, and plays an important role in disease characterization, drug development and precision medicine [1][2][3][4].

Challenges and pitfalls in the generation, analysis, and interpretation of high-throughput expression profiling data need to be addressed within the scientific community. Development of omics-based products that influence or improve patient health has been slower than expected. Studies attempting to reproduce findings of 53 papers in preclinical cancer research confirmed only 6 (11 %) of the results [5]. Misleading papers result in considerable expenditure of time, money and effort by researchers following false trails. This affects companies and investors, presenting yet another barrier for the translation of academic discoveries into new medicines by diverting funds away from real advances [6][7]. Irreproducible or inconsistent results could contribute to patient risk or death. As more and more irreproducible reports occur, some scientific journals reported the issue in 2014 [8][9]. The essential role of reproducibility of scientific research has been widely recognized [10].

There exist different reasons for low reproducibility in omics research. One possible reason is the complexity of omics data. The fact that the size of data is so massive that the manual inspection of data quality and analysis results is often impossible. Thus, quality control processes for high-throughput expression experiments are essential for the improvement of reproducibility of biological results.

The MicroArray and Sequencing Quality Control (MAQC/SEQC) consortia conducted three projects [11][12][13] to assess the reliability and reproducibility of genomics technologies, including microarrays, genome-wide association studies, and next-generation sequencing. This has led to the formation of the Massive Analysis and Quality Control Society (MAQC Society) [23], which is dedicated to quality control and analysis of massive data generated from high-throughput technologies for enhanced reproducibility and reliability [14]. It has provided a collection of quality metrics for expression data evaluation that corresponds to the reliability and reproducibility of high-throughput gene expression data for quality control, including (i) from sample to RNA, (ii) expression profiling, (iii) quality control metrics in RNA-seq, (iv) detecting differentially expressed genes, (v) biological interpretation, and (vi) spike-ins. Similar and complementary efforts have been reported elsewhere [15][16].

High-quality data are the foundation for deriving reliable biological conclusions from a gene-expression study. However, large differences in data quality have been observed in published data sets when the same platform was used by different laboratories. In many cases, poor quality of data was due not to the inherent quality problems of a platform but to the lack of technical proficiency of the laboratory that generated the data. Therefore, proficiency testing, an assessment of the overall competence performed through inter-laboratory comparisons, is introduced in this document to establish and monitor the quality of laboratory tests.

This document can be utilized to (i) enhance community's understanding of technical performance of high-throughput gene expression; (ii) benefit the interoperability of qualified gene-expression data by researchers, commercial entities and regulatory bodies, (iii) improve the application of high-throughput gene expression in industry and clinics, (iv) promote the acceptance of transparent reporting according to the FAIR (findable, accessible, interoperable, and reusable) data principles [17], and (v) contribute to the development of precision medicine.

Genomics informatics — Reliability assessment criteria for high-throughput gene-expression data

1 Scope

This document specifies reliability assessment criteria for high-throughput gene-expression data.

It is applicable to assessing the accuracy, reproducibility, and comparability of gene-expression data that are generated from microarray, next-generation sequencing, and other forms of high-throughput technologies.

This document identifies the quality-related data for the process of the next-generation sequencing of RNA (RNA-seq). The sequencing platform covered by this document is limited to short-read sequencers. The use of RNA-seq for mutation detection and virus identification is outside of the scope of this document.

This document is applicable to human health associated species such as human, cell lines, and preclinical animals. Other biological species are outside the scope of this document.

From a biological point of view, expression profiles of all genetic sequences including genes, transcripts, isoforms, exons, and junctions are within the scope of this document

2 Normative references

There are no normative references in this document.

3 Terms and definitions

For the purposes of this document, the following terms and definitions apply.

ISO and IEC maintain terminology databases for use in standardization at the following addresses:

- ISO Online browsing platform: available at <https://www.iso.org/obp>
- IEC Electropedia: available at <https://www.electropedia.org/>

3.1

adapter

short, chemically synthesized oligonucleotide that can be ligated to the ends of DNA or RNA molecules

3.2

alignment ratio

percentage of total sequenced reads aligned to an intended target region

Note 1 to entry: Alignment ratio is different according to the definition of the gene structure ("annotation") in that region, and is also affected by the origin of the RNA sample, library preparation, sequencing approach, and aligner.

3.3

batch effect

systematic technical variation in data unrelated to biological factors of interests and caused by processing samples in different batches